

Université de Lorraine — École Nationale Supérieure des Mines de Nancy
Année universitaire 2024–2025 — Semestre 6 — Tronc commun scientifique

Inférence statistique^{*}

par Rémi PEYRE[†], maître de conférences

version du 8 mai 2025

*. Réf. 6JUCSN01

†. remi.peyre@mines-nancy.univ-lorraine.fr



Ce polycopié est mis à votre disposition selon les termes de la licence Creative Commons avec attribution et partage dans les mêmes conditions (CC BY-SA 4.0) : voir le texte de la licence sur

<https://creativecommons.org/licenses/by-sa/4.0/legalcode.fr>.

En résumé, vous êtes autorisé·e à copier et à distribuer ce document, ainsi qu'à le transformer ou à l'utiliser dans d'autres créations, y compris à des fins commerciales, selon les conditions suivantes :

- *Vous devez créditer l'auteur du document et indiquer si des modifications y ont été apportées ; cela, sans toutefois suggérer que l'auteur soutiendrait la façon dont vous utilisez le document. Vous devez également intégrer un lien vers la présente licence.*
- *Dans le cas où vous transformez ce document ou produisez une création à partir de celui-ci, l'œuvre ainsi produite devra être diffusée sous une licence compatible avec celle-ci.*

Comment lire ce polycopié

Comme vous pourrez le constater, ce polycopié regorge de remarques, exemples et autres commentaires, dont la fonction est d'aider à la pleine appréhension des concepts présentés, mais qui peuvent dérouter l'élève en phase d'apprentissage ou de révisions, de par le volume un peu effrayant d'information que tout cela représente... En fait, tout dans ce polycopié n'est pas à retenir par cœur, loin de là! ☺ Pour s'y retrouver, les points qu'il faut vraiment connaître (qui sont principalement les définitions essentielles de la science statistique) sont signalés dans la marge par une barre noire commençant par un point d'exclamation : on peut considérer qu'il s'agit là de ce qu'il faut *apprendre*, le reste du polycopié étant plutôt à *comprendre*. Certains passages mettant en avant des points tout particulièrement essentiels sont même marqués d'une barre noire plus épaisse, accompagnée d'un *double* point d'exclamation ! Au sein des passages non marqués dans la marge, l'importance des différents paragraphes est assez variable. J'ai typographié en petits caractères (et avec une marge agrandie) les passages dont j'estime qu'ils peuvent être totalement sautés sans nuire à la bonne acquisition du cours : je les ai même carrément exclus, ainsi que la plupart des notes de bas de page, de la version "allégée" du polycopié.

Exemple. La marque en marge de ce passage indiquera, dans la corps du polycopié, qu'il s'agit d'un passage qu'il faut connaître rigoureusement, notamment en vue de l'examen. !

Exemple. Dans le corps du polycopié, la marque en marge de ce passage signifiera qu'il s'agit d'un point du cours super-méga-important ! !!

Exemple. À l'inverse, ce passage-ci est un complément "culturel" qui peut être sauté sans nuire à la compréhension du cours en tant que tel, mais dont j'ai estimé que vous le trouveriez peut-être intéressant, notamment si vous souhaitez aller plus loin dans la science statistique au cours des années à venir. ☺

Remarque. À l'heure où de vifs débats interpellent la société contemporaine sur la transmission, consciente ou non, des stéréotypes de genre — une question particulièrement prégnante sein de la communauté mathématique, où les femmes sont aujourd'hui largement sous-représentées —, plusieurs personnes m'ont fait remarquer que des formulations comme « le futur ingénieur » ou « le statisticien » pouvaient être perçues comme excluant implicitement les femmes du champ de l'ingénierie ou de la statistique. D'un autre côté, écrire un tel polycopié en écriture inclusive aurait été stylistiquement assez lourd... J'ai donc fait le choix de restreindre l'usage de l'écriture inclusive à cette seule préface : dans la suite du texte, lorsque je devrai désigner « le lecteur », « le statisticien » etc., j'utiliserai un masculin grammatical ; mais celui-ci devra simplement être compris comme un neutre sémantique! ☺

Table des matières

I	Pré-requis de probabilités	9
1	Quelques notations utilisées dans ce polycopié	11
1.1	Divers	11
1.2	Booléens	12
1.3	Infinitésimaux	13
2	Distributions de probabilité et variables aléatoires	15
2.1	Distributions de probabilité	15
2.2	Variables aléatoires	18
2.3	Changement de variables multidimensionnel	20
3	Statistiques descriptives des distributions de probabilité	25
3.1	Fonction de répartition	25
3.2	Quantiles et intervalles de fluctuation	27
3.3	Mode	31
3.4	Quantités empiriques	32
3.5	Théorie de la décision optimale	33
4	Conditionnement et indépendance	37
4.1	Conditionnement	37
4.2	Indépendance	40
5	Théorèmes-limites	41
5.1	La loi des grands nombres	41
5.2	Le théorème-limite central	46
6	Distributions de probabilité remarquables	53
6.1	Familles de lois et paramétrages	53
6.2	Quelques lois remarquables discrètes	55
6.3	Les lois normales	62
6.4	Autres lois continues remarquables	66
II	Statistique bayésienne	71
7	Concept d'inférence statistique	73
7.1	Motivation : Deux exemples	73
7.2	Modèle d'inférence statistique	74
7.3	Notion d'inférence	78

7.4	Paramètres des modèles	79
7.5	Statistiques bayésienne et fréquentiste	81
7.6	Description d'un modèle statistique en langage ordinaire	83
7.7	Récapitulatif des principales notations	86
8	Le théorème de Bayes	89
8.1	Notion de postérieure	89
8.2	La formule de Bayes	92
8.3	Notion de priore impropre	94
8.4	Exemples de calculs	95
9	Méthodes statistiques bayésiennes	103
9.1	Importance des lois à postériori	103
9.2	Décision optimale en contexte bayésien	104
9.3	Probabilités à postériori	107
9.4	Intervalles de confiance et de prédiction bayésiens	112
9.5	Décisions et estimateurs	116
9.6	Conclusion sur l'analyse bayésienne	121
III	Méthodes fréquentistes	125
10	Notion de vraisemblance	127
10.1	Définition générale et type d'objet	127
10.2	Comment calculer la vraisemblance	132
10.3	Extensions du concept de vraisemblance	133
10.4	Deux propriétés remarquables de la vraisemblance	136
11	Estimation et prédiction en statistique fréquentiste	139
11.1	Prolégomènes	139
11.2	Fonctions de risque et de biais	140
11.3	Convergence des estimateurs	147
11.4	Estimateur du maximum de vraisemblance	150
11.5	Estimateurs empiriques	153
11.6	Méthode des moindres carrés	156
11.7	Estimation par tendance	158
11.8	Technique de substitution	159
12	Tests d'hypothèses	163
12.1	Motivation du concept	163
12.2	Formalisation du concept dans le cas unilatéral	167
12.3	Construction d'un test à partir d'une statistique	171
12.4	Recherche de statistiques de test pertinentes	177
12.5	Tests bilatéraux	182
12.6	Régime asymptotique des tests	184
13	La p-valeur	189
13.1	p -valeur associée à une statistique de test	189
13.2	Notion intrinsèque de p -valeur	193
13.3	Interprétation d'une p -valeur	194

13.4	Un exemple de test par p -valeur	196
13.5	Utilisation pratique d'un test par l'ingénieur	199
14	Intervalles de confiance et de prédiction fréquentistes	201
14.1	Intervalles de confiance	201
14.2	Intervalles de prédiction	208
IV	Du bon usage de la statistique	211
15	Choix de la priore	213
15.1	Problématique	213
15.2	Cas où la priore est non ambiguë	216
15.3	Postérieure d'hier, priore d'aujourd'hui	216
15.4	Traduction mathématique d'une expertise préalable	218
15.5	Priores non informatives	222
16	Fonctions d'utilité	227
16.1	Illustration du concept	227
16.2	Considérations sur la notion d'utilité	229
16.3	Un exemple en contexte industriel	232
17	Pièges et paradoxes des tests	235
17.1	Que signifie « accepter une hypothèse nulle » ?	235
17.2	Comment choisir l'hypothèse nulle	236
17.3	Trop de tests tuent le test !	243
17.4	Le paradoxe du test à droite	246
17.5	p -valeur vs taille d'effet	248
18	Visualisation des données	251
18.1	Représentation des distributions de probabilité	251
V	Modèles statistiques classiques	255
A	La régression linéaire	257
A.1	Régression linéaire simple	257

Première partie

Pré-requis de probabilités

Chapitre 1

Quelques notations utilisées dans ce polycopié

1.1 Divers

- Le *pour-cent* est une unité mathématique sans dimension, notée ‘%’, valant $\frac{1}{100}$. De même, le *pour-mille*, noté ‘‰’, vaut $\frac{1}{1000}$, et le *pour-dix-mille* (rarement utilisé), noté ‘‱’, vaut $\frac{1}{10000}$. Il ne faut rien entendre de plus dans les symboles ‘%’, ‘‰’ et ‘‱’ qu’une simple division par resp. 100, 1 000, et 10 000 : « 3,89 ‰ », par exemple, ne signifie rien de plus que « 0,003 89 » ! L’usage est cependant de réserver l’emploi des signes ‘%’ et compagnie à des quantités pouvant être interprétées comme des proportions ou des probabilités.
- Pour A et B deux ensembles disjoints, « $A \sqcup B$ » a la même signification que « $A \cup B$ », mais rappelle simplement que A et B sont disjoints. On utilise de même la notation « $\bigsqcup_{i \in I} A_i$ » pour l’union d’un nombre quelconque d’ensembles deux-à-deux disjoints.
- Ω étant un ensemble, pour $A \subseteq \Omega$, le complémentaire de A dans Ω est noté $\complement_{\Omega} A$. Bien souvent l’ensemble ambiant Ω est évident, et on note alors simplement « $\complement A$ ».
- \mathbb{N} désigne l’ensemble $\mathbb{N} \cup \{\infty\}$, muni de sa structure d’ordre et de sa topologie usuelles ; $\overline{\mathbb{R}}$ désigne l’ensemble $\mathbb{R} \cup \{-\infty, +\infty\}$, muni de sa structure d’ordre et de sa topologie usuelles.
- Pour a et b des réels, $a \wedge b$ désigne le minimum de a et b , tandis que $a \vee b$ désigne leur maximum.
- Pour a un nombre réel, a^+ et a^- désignent respectivement la partie positive et la partie négative de a , soit resp. $a \vee 0$ et $(-a) \vee 0$.
- J’utiliserai parfois la notation \vec{x}_I , pour ‘ x ’ un symbole et I un ensemble, comme raccourci pour $(x_i)_{i \in I}$: par exemple, le 9-uplet (a_1, a_2, \dots, a_9) pourra être abrégé en « $\vec{a}_{[1,9]}$ ».
- $\{\{a, b, \dots, z\}\}$ désigne le *multiensemble* contenant les valeurs a, b, \dots et z , c’est-à-dire la liste des valeurs a, \dots, z où l’ordre dans lequel sont présentées les valeurs n’est pas pris en compte, mais le *nombre d’occurrences* de cette valeur, si ! Par exemple, le multiensemble des lettres du mot « multiensemble », $\{\{m, u, l, t, i, e, n, s, e, m, b, l, e\}\}$ peut aussi être écrit $\{\{b, e, e, e, i, l,$

- $l, m, m, n, s, t, u\}$, mais ce n'est par contre par la même chose que $\{b, e, i, l, m, n, s, t, u\}$.
- « $(X) := (Y)$ » signifie que l'on définit (X) comme étant égal à (Y) ; de même, « $(X) = (Y)$ » signifie que l'on introduit la notation (Y) comme valant (X) . Par contre, « $(X) \stackrel{\text{déf}}{=} (Y)$ » signifie que l'égalité entre (X) et (Y) se *déduit* de la définition, soit de (X) , soit de (Y) .
 - Pour $A \subseteq \mathbb{R}^d$ de forme suffisamment régulière, $\text{vol}(A)$ désignera le volume (d -dimensionnel) de A (alias « longueur » en dimension 1 ou « aire » en dimension 2).
 - Pour \mathcal{X} un ensemble, l'ensemble des distributions de probabilité sur \mathcal{X} sera noté $\mathcal{M}_1(\mathcal{X})$.

1.2 Booléens

La fonction *négation*, notée « \neg » (parfois utilisée sans parenthèses en tant qu'opérateur, comme on le fait p. ex. pour la fonction \sin), est la fonction de $\{\text{VRAI}, \text{FAUX}\}$ dans lui-même définie par $\neg\text{VRAI} := \text{FAUX}$ et $\neg\text{FAUX} := \text{VRAI}$. Il existe également les opérateurs « et » et « ou », de $\{\text{VRAI}, \text{FAUX}\}^2$ dans $\{\text{VRAI}, \text{FAUX}\}$, dont les définitions sont évidentes : nous n'introduisons pas de notations pour ces opérateurs, qui seront écrits en toutes lettres. La *fonction indicatrice*, notée « $\mathbf{1}_\bullet$ », est la fonction de $\{\text{VRAI}, \text{FAUX}\}$ dans $\{0, 1\}$ définie par $\mathbf{1}_{\text{VRAI}} := 1$ et $\mathbf{1}_{\text{FAUX}} := 0$. Nous utiliserons aussi la fonction « co-indicatrice » $\mathbf{0}_A := \mathbf{1}_{\neg A} = 1 - \mathbf{1}_A$.

Une écriture comme « $a < b$ » ou « $f(x, y) \in A$ » est ce qu'on appelle une *proposition* : autrement dit, c'est un énoncé auquel, selon la valeur des paramètres libres de cet énoncé (lesdits paramètres étant a et b dans le premier cas, resp. x et y dans le second cas si nous supposons que f et A se réfèrent à une fonction et un ensemble fixés), on attribue un booléen qui dit si l'énoncé tient ou non, booléen qu'on qualifie de *valeur de vérité* de l'énoncé. Ainsi, une proposition peut être vue comme une application qui va de l'espace que peuvent prendre les paramètres libres dans l'ensemble $\{\text{VRAI}, \text{FAUX}\}$, associant à chaque valeur possible du jeu de paramètres la valeur de vérité de la proposition pour cette valeur-là. Dans le cadre de ce cours, où on se place dans un paradigme probabiliste, les variables libres des propositions seront des variables *aléatoires*, de sorte que la valeur de vérité d'une proposition sera une variable aléatoire à valeurs dans $\{\text{VRAI}, \text{FAUX}\}$: par exemple, si mon univers probabiliste est Ω et que $X := X(\omega)$ est une variable aléatoire à valeurs réelles, la proposition « $X > 0$ » est une variable aléatoire qui vaut VRAI sur les $\omega \in \Omega$ tels que $X(\omega) > 0$, resp. FAUX sur les ω tels que $X(\omega) \leq 0$. Et l'indicatrice $\mathbf{1}_{X > 0}$ de cette proposition sera la v.a. à valeurs dans $\{0, 1\}$ valant resp. 1 dans le premier cas et 0 dans le second.

Dans vos cours de probabilités de lycée et de S5, on vous a défini un évènement comme un sous-ensemble de l'univers probabiliste Ω . Dans ce cours, nous verrons plutôt les évènements comme des *propositions*, ce qui correspond bien mieux à la vision intuitive que nous en avons : par exemple, « le dé est tombé sur '4' » sera un évènement (en supposant que le résultat du dé soit modélisé par une variable aléatoire). En fait, là où vous voyiez auparavant un évènement comme un sous-ensemble $A \subseteq \Omega$, la contrepartie propositionnelle pour désigner cet évènement sera « $\omega \in \Omega$ » (qui est bien une variable aléatoire à valeurs booléennes) ; et inversement, pour P une proposition dont la valeur de vérité ne dépend que de l'éventualité ω , la vision

ensembliste de l'évènement P serait l'ensemble des ω pour lesquels P est vrai. Au niveau du formalisme, voir les évènements comme des propositions fera que, si E_0 et E_1 sont des évènements, on parlera de la conjonction de ces deux évènements en écrivant « E_0 et E_1 », et de leur disjonction en écrivant « E_0 ou E_1 », alors qu'en vision ensembliste vous auriez écrit resp. « $E_0 \cap E_1$ » et « $E_0 \cup E_1$ ». Au niveau notationnel, pour éviter de manipuler des guillemets dans les énoncés mathématiques, nous délimiterons si nécessaire une proposition comme « $f(x, y) \in A$ » par des accolades plutôt que par des guillemets, ce qui pourra donner les phrases comme, par exemple, « les évènements $\{X \geq 0\}$ et $\{|Y| \leq 1\}$ sont indépendants » : le choix de cette notation étant justifié par la correspondance implicite entre propositions et évènements.

1.3 Infinitésimaux

Deux notions légèrement distinctes d'intégration sont utilisées dans ce cours, que nous qualifierons respectivement d'intégration *algébrique* et *géométrique* :

Une intégrale algébrique s'écrit sous la forme :

$$\int_{x=a}^b f(x)dx. \quad (\text{AA})$$

Éventuellement on pourra omettre de rappeler la variable muette au niveau du signe d'intégration (« $\int_a^b f(x)dx$ »), voire se livrer à des simplification notationnelles encore plus radicales (jusqu'à « $\int_a^b f$ ») ; mais dans tous les cas, le domaine d'intégration est indiqué par des *bornes* figurant resp. en indice et en exposant du symbole d'intégration — corolairement, l'intégrale sera toujours en dimension 1^[*].

L'interprétation d'une telle intégrale est la suivante. On imagine que x se déplace de a à b par des pas infinitésimaux, se déplaçant successivement de $a =: x_0$ via x_1, x_2, \dots jusqu'à $x_N =: b$ ^[†], avec à chaque fois $|x_{i+1} - x_i| \ll 1$, le nombre N de pas étant évidemment très grand. Dans ce cas, « dx » représente l'*accroissement* $x_{i+1} - x_i$, qui est donc un *nombre réel* (infinitésimal) ; on note aussi cet accroissement infinitésimal δx lorsqu'on l'utilise en l'absence de symbole d'intégration associé.

À l'inverse, une intégrale géométrique s'écrit sous la forme :

$$\int_{x \in U} f(x) \text{vol}(dx). \quad (\text{AB})$$

(Des simplifications comme « $\int_U f(x) \text{vol}(dx)$ », voire « $\int_U f$ », seront tolérées). Cette fois-ci, le domaine d'intégration est écrit sous la forme d'un *ensemble* (et non de *bornes*), qui peut être de dimension quelconque. En outre, il apparaît dans l'intégrale une *mesure* (en général la mesure de volume^[‡]), qui est appliquée à dx .

Cette fois-ci, « dx » ne représente plus un nombre mais un *voisinage infinitésimal* (dans U) de x ; on notera aussi ce voisinage infinitésimal dx lorsqu'on l'utilise en

[*]. On peut en fait définir un avatar de l'intégrale de Riemann en dimension plus grande que 1 dans le cadre de la théorie des formes différentielles ; mais cela sort du champ de ce cours et même des compétences d'un ingénieur, fût-il mathématicien !

[†]. Noter que rien n'oblige ici la suite (x_i) à être monotone !

[‡]. Il sera *toléré* de sous-entendre la mesure lorsqu'il s'agit de la mesure de volume, écrivant alors « $\int_{x \in U} f(x)dx$ », mais cela est fortement déconseillé !

l'absence de symbole d'intégration associé (« dx » est alors une *variable* — à valeurs dans les parties infinitésimales de U). L'intégrale signifie alors qu'on découpe le domaine d'intégration U selon une *partition* en ensembles infinitésimaux dx_i , chaque dx_i étant repéré par un point x_i contenu dans ce voisinage (ou éventuellement à une distance infinitésimalement proche) ; et l'intégrale se réfère alors, formellement, à la quantité

$$\sum_{i=1}^N f(x_i) \text{vol}(dx_i). \quad (\text{AC})$$

Plus généralement, chaque fois que nous ferons intervenir une variable dx , il faudra comprendre que celle-ci désigne « un voisinage infinitésimal du point x » ! Par exemple, si j'écris quelque chose comme

$$\mathbb{P}(X \in dx) = f(x) \text{vol}(dx), \quad (\text{AD})$$

cela signifie que la loi de la variable aléatoire X admet, au point x , la densité $f(x)$ par rapport à la mesure de Lebesgue.

Remarque (AE). Dans le cas de la dimension 1, on se retrouve alors avec deux notions d'intégration. Celles-ci sont évidemment liées : pour $a \leq b$, on a $\int_{x=a}^b f(x)dx = \int_{x \in]a, b[} f(x) \text{vol}(dx)$, et pour $a > b$, $\int_{x=a}^b f(x)dx = -\int_{x=b}^a f(x)dx = -\int_{x \in]b, a[} f(x) \text{vol}(dx)$. Au passage, c'est parce que la formule prend une forme différente selon le sens de comparaison entre a et b qu'une valeur absolue intervient dans la formule de changement de variables en formalisme géométrique, mais pas en formalisme algébrique (confer remarque (BA)) : en effet, dans le théorème (AW), lorsque Φ est décroissante, Φ intervertit l'ordre des bornes de l'intervalle d'intégration algébrique, et il n'y a donc pas le même signe pour passer de l'intégrale algébrique à l'intégrale géométrique, d'où le fait que ce n'est pas Φ' mais $-\Phi'$ qui intervient dans ce cas : autrement dit, que Φ soit croissante ou décroissante, c'est $|\Phi'|$ qui intervient pour le changement de variable en formalisme géométrique. ♣

Chapitre 2

Distributions de probabilité et variables aléatoires

La statistique inférentielle s'appuie de façon essentielle sur la théorie des probabilités. Vous avez déjà vu diverses notions autour de cette théorie : d'abord au lycée et en classes préparatoires, sur un certain nombre de points épars, puis avec le cours de M. VILLEMONTAIS au semestre 5, grâce auquel vous avez pu regrouper et étendre ces notions au sein d'un formalisme rigoureux et cohérent.

Dans cette partie, nous allons rappeler les notions que vous avez vues, et introduire quelques extensions de celles-ci qui seront utilisées dans ce cours. Le formalisme sera moins rigoureux que dans vos cours précédents, car ici le but est de pouvoir manipuler de façon simple et intuitive les concepts nécessaires aux calculs statistiques, quitte à ce que le cadre mathématique sous-jacent comporte quelques imprécisions (en particulier : occultation des problématiques de mesurabilité, utilisation de notations à base infinitésimales, confusion des notions d'égalité et d'équivalence presque-sure).

2.1 Distributions de probabilité

Manipulation informelle des distributions de probabilité

Dans le cadre de ce cours, nous interpréterons les distributions de probabilités à densité comme un cas-limite des distributions discrètes. Plus précisément, nous considérerons qu'on peut découper \mathbb{R}^n en un très grand nombre de zones infinitésimales de la forme dx , où dx représente un « voisinage infinitésimal de x »^[*] dans \mathbb{R}^n , qu'on peut se représenter par exemple comme étant de la forme $\prod_{i=1}^n [x_i - \varepsilon, x_i + \varepsilon[$ avec ε très petit^[†]. Il est entendu que le découpage de \mathbb{R}^n en zones est fait

[*]. La notion de « voisinage » utilisée ici n'est pas exactement la même que celle utilisée en topologie : ici, quand nous disons que dx est un voisinage de x , cela signifie juste que la distance entre les points de dx et x est (au plus) du même ordre de grandeur que la distance entre deux points de dx . Par exemple, dans \mathbb{R} , lorsque $\varepsilon \searrow 0$, $[\varepsilon, 2\varepsilon]$ sera bien considéré comme un voisinage infinitésimal de 0 (même si en l'occurrence il ne contient pas 0 lui-même), mais pas $[\varepsilon^{1/2}, \varepsilon^{1/2} + \varepsilon]$ ne l'est pas, car la distance de cet intervalle à 0, qui est $\varepsilon^{1/2}$, est beaucoup plus grande que sa largeur, qui est ε .

[†]. Attention : Ici dx ne désigne donc pas un *accroissement infinitésimal* de x , qui serait un nombre (ou plus généralement un n -uplet), mais un *voisinage infinitésimal* de x , qui est donc une *zone infinitésimale* de \mathbb{R}^n . Accroissements et voisinages infinitésimaux sont donc notés de la même façon (ce choix est dû au fait qu'ils jouent des rôles très similaires en théorie de l'intégration),

de sorte que les zones constituent une *partition* de \mathbb{R}^n , c.-à-d. que tout point de \mathbb{R}^n appartienne à une zone et une seule. Dans ce cas-là, on pourra considérer que, pour une zone dx voisinage de x , c'est comme s'il y avait au point x une masse discrète très petite, dont la valeur serait proportionnelle au volume de dx . Ainsi, nous nous représentons notre distribution de probabilité à densité comme une myriade de minuscules masses discrètes. Le point crucial étant bien entendu que, lorsque « dx » devient réellement infinitésimal (c.-à-d. dans l'asymptotique idoine), l'objet qu'on obtient à la limite ne dépend pas de la façon précise dont on a maillé notre espace \mathbb{R}^n en zones.

Pour décrire notre distribution de probabilité à densité, il faudra alors dire quelle est la masse de chaque zone infinitésimale. Par exemple, on écrira : « La distribution de probabilité P est telle que, pour dx un voisinage infinitésimal de x , $P(dx) = f(x) \text{vol}(dx)$ », ce qui est simplement une façon intuitive de formaliser ce qu'on écrirait, en termes rigoureux, « P est la distribution sur \mathbb{R}^n dont la densité en x (par rapport à la mesure de Lebesgue) vaut $f(x)$ ». Bien entendu, l'expression dans le membre de droite doit toujours être proportionnelle au volume de dx .

Les grandes familles de distributions de probabilité... et les autres !

Distributions de probabilité ni discrètes, ni à densité Dans les cours de probabilité et de statistique, on manipule majoritairement des distributions de probabilité discrète et des distributions de probabilité à densité. Il pourrait dès lors être tentant, pour l'ingénieur, de croire que ce sont les seuls types de distributions de probabilité qu'on est amené à rencontrer dans un cadre industriel... Mais c'est FAUX ! Même si effectivement, certains cas très "exotiques" relèvent surtout de l'univers des mathématiciens "purs", il y a des cas très concrets qui sortent de la dichotomie entre mesures discrètes et à densité, et notamment des mesures "mixtes" qui "combinent" ces deux concepts.

Exemple (AF) (Données tronquées). Imaginez que vous êtes chargé·e de contrôler la qualité de la production d'une usine de pompes hydrauliques. Vous allez échantillonner des pompes sorties de la chaîne de production et les tester selon un certain protocole jusqu'à ce qu'elles tombent en panne ; votre objectif étant de dire « lorsque je prends une pompe sortie de l'usine, la probabilité que sa durée de vie soit comprise entre telle et telle durée est égale à telle valeur ». Disons qu'en général, les pompes tiennent de l'ordre de 30 jours [symbole du jour : 'd'] dans les conditions du test. Vous pourriez bien entendu tester vos pompes jusqu'à ce qu'elles soient toutes tombées en panne. Mais cela prendrait beaucoup de temps pour pas grand-chose... En pratique, on se fixe une durée maximale pour le test, disons 60 d. La durée de vie mesurée pour une pompe testée peut alors être de deux types : soit, par exemple « Cette pompe a tenu 20,8 d » ; soit « Cette pompe a tenu 60 d, c.-à-d. qu'elle était encore en fonctionnement à la fin du test ». La durée de vie mesurée pour une pompe prend donc ses valeurs dans l'intervalle $[0, 60]$ d, mais ce sera typiquement une distribution de probabilité *mixte* ayant une densité sur $[0, 60[$, et une masse ponctuelle au point 60, au sens où on décrira alors une telle distribution de probabilité par une formulation du type suivant : « Pour dx un voisinage infinitésimal de x , la masse attribuée par la distribution de probabilité à un voisinage infinitésimal dx de x vaut

mais le contexte permettra toujours de trancher quant au concept auquel la notation « dx » fait allusion.

- $(3,333 \times 10^{-2} \text{ d}^{-1}) \exp(-x / (30 \text{ d})) \text{ vol}(dx)$ si $x \in]0, 60[$;
- 13,534 % si $x = 60$.

». Notez bien qu'il y a un facteur « $\text{vol}(dx)$ » dans le premier cas, qui correspond à la partie à densité de notre distribution de probabilité, mais qu'il n'y en a *pas* dans le second, qui correspond à la partie discrète ! \clubsuit

Remarque (AG). Dans un tel cas, il est essentiel de ne surtout pas chercher à comparer une valeur de *densité* (par exemple, la densité la probabilité de la durée de vie autour de 7 d vaut $2,574 \times 10^{-2} \text{ d}^{-1}$) à une valeur de *probabilité* (par exemple, la probabilité d'une durée de vie de 60 j vaut 14,435 %) : cela reviendrait en substance à comparer une aire et une longueur, et d'ailleurs dans le cas présent ces deux quantités n'ont pas la même homogénéité physique ! \clubsuit

Exemple (AH) (Données inflatées en zéro). Si on cherche la plus grande faille dans un échantillon de matériau, la taille de cette plus grande faille suit une distribution sur \mathbb{R}_+ qui a une densité en général ; cependant il y a aussi une probabilité non nulle qu'il n'y ait *aucune* faille, de sorte que la distribution a aussi une masse non nulle en zéro... Dans ce cas, la loi de probabilité P suivi par la longueur de la plus grande faille sera telle que, pour dx un voisinage infinitésimal de x ,

$$P(dx) = \begin{cases} p_0 & \text{pour } x = 0 ; \\ f(x) \text{ vol}(dx) & \text{pour } x > 0, \end{cases} \quad (\text{AI})$$

où p_0 et f sont tels que $p_0 + \int_{x=0}^{\infty} f(x)dx = 1$. \clubsuit

Mentionnons enfin quelques cas de distributions ni discrètes ni à densité plus rares, mais qu'on peut tout de même parfois rencontrer dans les situations d'ingénierie :

Exemple (AJ) (Distribution uniforme sur la sphère). Si on cherche à décrire le lieu aléatoire d'un évènement se produisant à la surface de la Terre (par exemple, la chute d'une météorite), il est naturel de modéliser ladite surface de la Terre par une sphère. Mais alors, si on considère par exemple la densité uniforme sur cette sphère (au sens où la probabilité de tomber sur une certaine partie du globe est proportionnelle à la *surface* de celle-ci), on aura une mesure qui n'est ni discrète, ni à densité dans l'espace tridimensionnel ! (En l'occurrence on pourrait *paramétrer la surface de la Terre, par exemple par ses latitude et longitude, pour sa ramener à une distribution probabilité à densité en dimension 2* ; mais fondamentalement la "vraie" probabilité n'est pas à densité). Dans ce cas, la probabilité d'un voisinage infinitésimal de x est nulle si x n'est pas sur la sphère ; et si x est sur la sphère, elle est proportionnelle à la *surface* infinitésimale de l'*intersection* $dx \cap \mathcal{S}$ (où \mathcal{S} désigne notre sphère). Sachant que définir la notion de surface n'est pas si simple qu'il n'y paraît \ddagger ...

L'exemple ci-dessus rentre plus généralement dans le cadre des distributions de probabilité à densité sur une « sous-variété » de \mathbb{R}^d : en quelque sorte, il s'agit de

\ddagger . Pour ceux que cela intéresse, si M est une surface plongée dans l'espace tridimensionnel, x un point de M , et que dx est un voisinage infinitésimal de x dans la surface M de forme (presque) parallépipédique, avec pour côtés \vec{dx}_1 et \vec{dx}_2 (autrement dit, dx est l'ensemble des $x + \alpha \vec{dx}_1 + \beta \vec{dx}_2$ lorsque (α, β) décrit un ensemble de la forme $]a, a + 1[\times]b, b + 1[$), alors la surface de dx est définie comme le *volume* tridimensionnel du parallépipède de côtés \vec{dx}_1 , \vec{dx}_2 et $\vec{\Delta x}_3$, où $\vec{\Delta x}_3$ est un vecteur unitaire orthogonal à la fois à \vec{dx}_1 et \vec{dx}_2 — et ce volume peut d'ailleurs se calculer comme la valeur absolue du déterminant du triplet $(\vec{dx}_1, \vec{dx}_2, \vec{\Delta x}_3)$.

distributions qui sont portées par des espaces de dimension intermédiaire entre 0 et d .

2.2 Variables aléatoires

Loi d'une variable aléatoire

Notion de mesure image L'introduction du concept de loi est l'occasion de définir celui de *mesure image* :

! **Définition (AK)** (Mesure image). Soit P une distribution de probabilité sur un espace \mathcal{X} , \mathcal{Y} un autre espace, et $f: \mathcal{X} \rightarrow \mathcal{Y}$ une fonction (déterministe). Alors, si X est une variable aléatoire sur Ω à valeurs dans \mathcal{X} ayant pour loi P , la loi de $f(X)$ (qui est une variable aléatoire sur Ω à valeurs dans \mathcal{Y}) ne dépend que de P et de f (cf. proposition (AL)) : on appelle cette distribution de probabilité (qui porte sur \mathcal{Y}) la *mesure image de P par f* , et on la note $f_* P$. \heartsuit

Proposition (AL) (Formule de la mesure image). La mesure image $f_* P$ peut être caractérisée directement par la formule suivante : pour tout $B \subseteq \mathcal{Y}$,

$$(f_* P)(B) := P(f^{-1}(B)). \quad (\text{AM})$$

\diamond

Démonstration. Soit X une v.a. à valeurs dans X de loi P . Alors, pour $B \subseteq Y$,

$$\begin{aligned} (f_* P)(B) &\stackrel{\text{déf}}{=} (\text{Loi}(f(X)))(B) \stackrel{\text{déf}}{=} \mathbb{P}(f(X(\omega)) \in B) \\ &= \mathbb{P}(X(\omega) \in f^{-1}(B)) \stackrel{\text{déf}}{=} (\text{Loi}(X))(f^{-1}(B)) = P(f^{-1}(B)). \end{aligned}$$

\heartsuit

Remarque (AN). La notion de mesure image permet de re-formuler la définition d'une loi de façon plus compacte : la loi de la v.a. X est simplement la mesure-image de \mathbb{P} par X (vue comme une application de Ω dans \mathcal{X}). \clubsuit

Remarque (AO). La notion de mesure image est un outil puissant pour définir de nouvelles distributions de probabilité. Si par exemple P est la distribution uniforme sur $[0, 1]$ et $f: \mathbb{R} \rightarrow \mathbb{R}$ la fonction « carré », on peut définir en quelques mots la mesure image $f_* P^{[\S]}$: on pourra dire « la mesure image par la fonction carré de la mesure uniforme sur $[0, 1]$ », ou encore « la loi de X^2 pour $X \sim \text{Unif}^{\text{me}}(0, 1)$ » si on préfère parler en termes de variables aléatoires^[¶]. \clubsuit

Quelques abus de notations sur les lois Dans ce cours, nous utiliserons par commodité deux notations où (dans ces contextes précis) les lois sont assimilées à des variables aléatoires suivant ces lois :

[§]. Qui, en l'occurrence, n'est autre que la distribution bêta de paramètres $(1/2, 1)$.

[¶]. Parler en termes de variables aléatoires pourra s'avérer commode lorsqu'on cherche à créer une nouvelle loi à partir de *plusieurs* lois de départ : par exemple, il est plus intuitif de comprendre la définition « la loi de $(X+Y)$ pour $X \sim \text{Unif}^{\text{me}}(0, 1)$ et $Y \sim \text{Normale}(0, 1)$ avec X et Y indépendantes » que « la mesure image, par la fonction somme, du produit des lois $\text{Unif}^{\text{me}}(0, 1)$ et $\text{Normale}(0, 1)$ ».

— Une notation comme

$$\mathbb{P}(\text{Expon}^{\text{le}}(1) \geq 3) \quad (\text{AP})$$

désignera la probabilité qu'une variable aléatoire de loi *Exponentielle*(1) soit supérieure ou égale à 3 : cette probabilité ne dépend en effet pas de la variable aléatoire précise considérée mais bien uniquement de sa distribution, puisque ce n'est autre que $(\text{Expon}^{\text{le}}(1))(]3, \infty])$.

— Une notation comme

$$\text{Normale}\left(0, \frac{1}{\sqrt{2}}\right)^2 \quad (\text{AQ})$$

désignera la loi du carré d'une variable aléatoire de loi Normale $\left(0, \frac{1}{\sqrt{2}}\right)$: là encore, cela ne dépend pas de la variable aléatoire précise choisie mais uniquement de sa loi, puisque ce n'est autre que la mesure-image, par la fonction $x \mapsto x^2$, de la loi Normale $\left(0, \frac{1}{\sqrt{2}}\right)$.

Remarque (AR). Bien sûr, on peut combiner les deux notations ci-dessus : par exemple,

$$\mathbb{P}(|\text{Normale}(0, 1)| \geq 1,733) \quad (\text{AS})$$

désigne la probabilité que la valeur absolue d'une variable aléatoire de loi Normale(0, 1) dépasse 1,733. ♣

Remarque (AT). En revanche, attention : une expression comme « $\text{Unif}^{\text{me}}(0, 1) \times \text{Expon}^{\text{le}}(1)$ » n'aurait pas de sens ! En effet, de la seule information que X et Y suivent resp. les lois $\text{Unif}^{\text{me}}(0, 1)$ et $\text{Expon}^{\text{le}}(1)$, on ne peut pas déduire la loi de $X \times Y$: cela dépend en effet de la loi jointe de (X, Y) , comme nous le verrons plus loin. La notation assimilant une loi de probabilité à une variable aléatoire suivant cette loi ne fait sens que lorsqu'il y a *une seule* variable aléatoire impliquée ! (et donc, dans la formule, une seule distribution de probabilité). ♣

Remarque (AU). Pour continuer sur la remarque ci-dessus, en revanche, la loi de $X \times Y$ sera bien spécifiée dès lors qu'on impose que X et Y soient *indépendantes* : dans ce cas, la loi correspondante est la mesure-image, par l'application « produit (sur \mathbb{R}) », de la loi-produit $\text{Unif}^{\text{me}}(0, 1) \otimes \text{Expon}^{\text{le}}(1)$. Dans le cas où on considère la *somme* de deux variables indépendantes, cette notion de « somme de variables aléatoires indépendantes suivant les lois de » correspond en fait à un *produit de convolution* ^{[[]]} : on peut alors noter $P * Q$ la mesure-image de $P \otimes Q$ par la fonction « somme ». ♣

Remarque (AV). Pour P une distribution de probabilité, il faut bien prendre garde à ne pas confondre une expression comme « $P + a$ », qui désigne en vertu du formalisme précédent la mesure-image de la distribution P par l'application $x \mapsto x + a$ (distribution qui n'est autre que le produit de convolution $P * \delta_a$) avec une expression comme « $P + \delta_a$ », qui désigne la distribution de masse (ici ce n'est pas une distribution de probabilité, car la masse totale vaut 2) qui, à tout ensemble X , associe la masse $P(X) + \delta_a(X) = P(X) + \mathbf{1}_a \in X$. ♣

[[]]. Ici la notion de « produit de convolution » s'entend sur des *distributions de masse*, et non sur des *fonctions*. On utilise néanmoins le même nom, car les deux définitions coïncident lorsqu'on assimile une fonction à la distribution de masse ayant cette fonction pour densité.

2.3 Changement de variables multidimensionnel

La formule

!

Théorème (AW) (Formule du changement de variables multidimensionnel). Soient $d \in \mathbb{N}$, $M \in \mathbb{R}^d$, et $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ une application de classe C^1 au voisinage de M . Alors, pour dM un voisinage infinitésimal de M dans \mathbb{R}^d :

$$\text{vol}(\varphi(dM)) = |J_\varphi(M)| \text{vol}(dM), \quad (\text{AX})$$

où $J_\varphi(M)$ est le jacobien de φ évalué en M , c.-à-d. le déterminant la matrice de ses dérivées partielles :

$$J_\varphi(M) := \det \left(\left(\frac{\partial \varphi_i}{\partial x_j}(M) \right)_{ij} \right). \quad (\text{AY})$$

— Dans cette formule, on a désigné comme d’habitude par « φ_i » les composantes de la fonction φ , et par « x_j » les coordonnées canoniques de \mathbb{R}^d : le déterminant considéré est donc celui d’une matrice $d \times d$. \diamond

Remarque (AZ). Ne pas oublier la valeur absolue...! (d’ailleurs, un volume ne peut être que positif!). Noter que, lorsque le domaine de φ est connexe et que φ est bijective (ce qui sera systématiquement le cas lorsqu’on procèdera à un changement de variables), le jacobien sera forcément de signe constant, de sorte que la valeur absolue se manifestera simplement comme un éventuel changement de signe sur le résultat. \clubsuit

Remarque (BA). Noter qu’en dimension 1, la formule de changement de variable pour laquelle on tombe n’est *pas la même* que celle à laquelle vous êtes habitués : il y a ici une valeur absolue qui apparaît... Cela provient de la distinction entre le formalisme « algébrique » et le formalisme « géométrique » pour l’intégration (le premier étant celui que vous avez vu en classes préparatoires, le second celui qu’on utilise en théorie de la mesure) : voir aussi la remarque (AE). \clubsuit

Le théorème (AW) contient la formule essentielle — qui est, à mes yeux, la seule formule du cours qu’on soit plus ou moins obligé de retenir “bêtement”^[**] — ; mais il est pratiquement impossible de comprendre comment la mettre en œuvre sans avoir vu un exemple au préalable.

Exemples de mise en œuvre pour le calcul d’une mesure-image

Maintenant nous allons montrer comment la formule de changement de variables multidimensionnel peut être utilisée pour déterminer une mesure-image. Une première possibilité est d’utiliser la méthodes des fonctions-tests, ce qui nous ramène à la situation précédente avec les intégrales. Mais en fait, le formalisme infinitésimal que nous avons introduit permet de trouver *directement* la mesure-image sans passer par des intégrales, ce qui est plus simple — et limite du coup les risques d’erreur! \smile

[**]. Et encore... En fait, cette formule est plus ou moins la *définition* du déterminant : il faudrait donc dire, plutôt, que ce sont les formules pour calculer les déterminants qui sont à retenir “bêtement”! \smile

Exemple (BB) (Logarithme d'une loi exponentielle). Rappelons ici que la loi $\text{Expon}^{\text{le}}(\lambda)$, pour $\lambda \in \mathbb{R}_+^*$, est la loi décrite par la densité suivante (où dx est un voisinage infinitésimal de x) :

$$\mathbb{P}(\text{Expon}^{\text{le}}(\lambda) \in dx) = \mathbf{1}_{x > 0} \lambda e^{-\lambda x} \text{vol}(dx). \quad (\text{BC})$$

Cette loi ne prenant presque-surement que des valeurs strictement positives, on peut se demander quelle est sa mesure-image par la fonction \log_b ($\log_b(\bullet)$ est le logarithme en base b : $\log_b(x) = \log(x) / \log(b)$). Pour cela, demandons-nous quelle est masse que la mesure-image $\log_b \ast \text{Expon}^{\text{le}}(\lambda)$ associe à un voisinage dy de $y \in \mathbb{R}$. C'est la probabilité que, pour une variable aléatoire X suivant la loi $\text{Expon}^{\text{le}}(\lambda)$, $\log_b X$ soit dans dy ; mais puisque $\log_b(\bullet)$ est une bijection de \mathbb{R}_+^* dans \mathbb{R} dont la bijection réciproque est l'exponentielle de base b , dire que $\log_b X \in dy$ est équivalent à dire que $X \in b^{\text{dy}}$:

$$\mathbb{P}(\log_b \ast \text{Expon}^{\text{le}}(\lambda) \in dy) = \mathbb{P}(\text{Expon}^{\text{le}}(\lambda) \in b^{\text{dy}}). \quad (\text{BD})$$

Or b^{dy} est évidemment un voisinage infinitésimal de b^y . Son volume se déduit du volume de dy via la formule de changement de variable, en observant que l'application $y \mapsto b^y$ a pour dérivée $y \mapsto \ln b \times b^y$:

$$\text{vol}(b^{\text{dy}}) = |\ln b| b^y \text{vol}(dy). \quad (\text{BE})$$

Dès lors, la densité de la loi $\text{Expon}^{\text{le}}(\lambda)$ nous donne la probabilité recherchée ci-dessus : (notons que l'indicatrice n'apparaît pas ici, car b^y est automatiquement positif),

$$\mathbb{P}(\log_b \ast \text{Expon}^{\text{le}}(\lambda) \in dy) = \lambda e^{-\lambda b^y} \text{vol}(b^{\text{dy}}) = |\ln b| \lambda b^y e^{-\lambda b^y} \text{vol}(dy), \quad (\text{BF})$$

ce qui nous donne bien la densité recherchée puisque le membre de droite est un multiple explicite de $\text{vol}(dy)$. \clubsuit

Dans l'exemple suivant, nous allons regarder comment on peut procéder lorsqu'on est amené à considérer une mesure-image par une fonction *non* injective. L'idée, instanciée ci-dessous, est de séparer (localement) le changement de variables selon plusieurs sous-domaines pour nous ramener à des changements de variables (localement) bijectifs.

Exemple (BG) (Carré d'une loi normale décentrée). Dans cet exemple nous allons considérer le cas d'une mesure-image par une fonction *non* injective : notre objectif sera de déterminer la densité de la mesure-image de la loi Normale(1, 1) par la fonction « carré », sachant que la loi Normale(1, 1) a pour densité :

$$\mathbb{P}(\text{Normale}(1, 1) \in dx) = \frac{1}{\sqrt{2\pi}} e^{-(x-1)^2/2} \text{vol}(dx). \quad (\text{BH})$$

Un carré étant toujours positif, il est évident que la loi Normale(1, 1)² est portée par \mathbb{R}_+ ; en outre, la probabilité que Normale(1, 1)² vaille 0 est nulle, puisque $\mathbb{P}(\text{Normale}(1, 1) = 0)$ (dans la mesure où la loi Normale(1, 1) est à densité et donc diffuse) : on peut donc s'intéresser à la mesure-image uniquement sur \mathbb{R}_+^* . Soit donc $y \in \mathbb{R}_+^*$ et dy un voisinage infinitésimal de y . Nous cherchons $\mathbb{P}(X^2 \in dy)$, où X suit la loi Normale(1, 1). Or dire que $X^2 \in dy$ est équivalent à dire que $X \in \sqrt{dy} \cup -\sqrt{dy}$,

où \sqrt{dy} et $-\sqrt{dy}$ sont des voisinages despectifs de \sqrt{y} et $-\sqrt{y}$: ces voisinages étant disjoints, on a donc $\mathbb{P}(X \in \sqrt{dy} \cup -\sqrt{dy}) = \mathbb{P}(X \in \sqrt{dy}) + \mathbb{P}(X \in -\sqrt{dy})$. Puisque la dérivée de la fonction $y \mapsto \sqrt{y}$ est $y \mapsto 1/2\sqrt{y}$, on a par la formule de changement de variable que $\text{vol}(\sqrt{dy}) = \text{vol}(dy) / 2\sqrt{y}$; et puisque la dérivée de la fonction $y \mapsto -\sqrt{y}$ est $y \mapsto -1/2\sqrt{y}$ (qui est négative, donc la valeur absolue donnera un changement de signe), on a $\text{vol}(-\sqrt{dy}) = \text{vol}(dy) / 2\sqrt{y}$ (c'est le même volume que pour \sqrt{dy} , ce qui est logique dans la mesure où la réflexion par rapport à 0 est une isométrie et conserve donc le volume). Finalement en appliquant la densité de la loi Normale(1, 1), on trouve que

$$\begin{aligned} \mathbb{P}(\text{Normale}(1, 1)^2 \in dy) &= \frac{1}{\sqrt{2\pi}} e^{-(\sqrt{y}-1)^2/2} \frac{\text{vol}(dy)}{2\sqrt{y}} + \frac{1}{\sqrt{2\pi}} e^{-(-\sqrt{y}-1)^2/2} \frac{\text{vol}(dy)}{2\sqrt{y}} \\ &= \frac{1}{\sqrt{2\pi e}} \cosh(\sqrt{y}) y^{-1/2} e^{-y/2} \text{vol}(dy), \end{aligned}$$

ce qui nous donne donc la densité de la loi Normale(1, 1)² (en gardant à l'esprit que cette densité sera nulle pour $y \leq 0$). \clubsuit

Maintenant nous allons regarder un cas où la mesure-image vit dans un espace de dimension strictement plus petit que la mesure de départ, ce qui pose à priori une difficulté, vu que la formule de changement de variables requiert d'avoir la même dimension des deux côtés. L'idée dans ce cas est d'introduire une ou plusieurs coordonnées supplémentaires "artificielles" à l'arrivée afin de rendre les dimensions égales, puis d'intégrer selon ces coordonnées supplémentaires : voir ci-dessous.

Exemple (BI). Supposons que nous ayons deux variables aléatoires X et Y qui soient indépendantes et que la loi de X aussi bien que la loi de Y soit $\text{Unif}^{\text{me}}(0, 1)$ ^[††]; notre objectif est de déterminer la loi du produit $XY =: Z$.

Pour commencer, déterminons la densité de la loi jointe de $(X, Y) =: M$. À partir de la densité de la loi uniforme et de la propriété d'indépendance, on écrit que, pour dx un voisinage de x et dy un voisinage de y :

$$\begin{aligned} \mathbb{P}(X \in dx \text{ et } Y \in dy) &= \mathbb{P}(X \in dx) \mathbb{P}(Y \in dy) \\ &= \mathbf{1}_{x \in]0, 1[} \text{vol}(dx) \mathbf{1}_{y \in]0, 1[} \text{vol}(dy) = \mathbf{1}_{(x, y) \in]0, 1[^2} \text{vol}(dx) \text{vol}(dy). \end{aligned}$$

Or dire que $(X \in dx \text{ et } Y \in dy)$ signifie que $M \in dx \times dy$, où $dx \times dy$ est un voisinage infinitésimal (de forme produit) du point $(x, y) =: m$, voisinage que nous pouvons donner dm , et dont le volume est (vu que la mesure de Lebesgue multidimensionnelle est un produit de mesures de Lebesgue 1-dimensionnelles) $\text{vol}(dm) = \text{vol}(dx) \text{vol}(dy)$. Autrement dit ce que nous avons montré tout à l'heure est que, pour dm un voisinage de m de forme rectangulaire, on a

$$\mathbb{P}(M \in dm) = \mathbf{1}_{m \in]0, 1[^2} \text{vol}(dm) : \quad (\text{BJ})$$

et puisque nous avons une formule exprimant la densité pour les voisinages de forme rectangulaire, cette formule est également valable pour les voisinages de toute forme,

[††]. Comme nous avons précisé que X et Y étaient indépendantes, cette connaissance des lois individuelles de X et Y permet de déterminer complètement la loi jointe de (X, Y) : *confer* théorème (FE).

de sorte que nous avons identifié la densité de la loi M sur \mathbb{R}^2 , qui est $\mathbf{1}_{\bullet \in]0, 1]^2}$. (En fait, nous avons redémontré une propriété bien connue : la densité d'une mesure-produit est le produit des densités).

Maintenant, notons $P := (XY, Y)$: nous avons ajouté Y comme seconde coordonnée pour que P vive en dimension 2 comme M . Il est alors clair que P est la mesure-image de M par l'application $\Phi : (x, y) \mapsto (xy, y)$. L'application Φ réalise une bijection de $(\mathbb{R}_+^*)^2$ dans $(\mathbb{R}_+^*)^2$, avec $\Phi^{-1} : (z, y) \mapsto (z/y, y)$. Nous pouvons calculer le jacobien de Φ^{-1} (en notant $x := z/y$) :

$$\begin{vmatrix} \partial x / \partial z & \partial x / \partial y \\ \partial y / \partial z & \partial y / \partial y \end{vmatrix} = \begin{vmatrix} 1/y & -z/y^2 \\ 0 & 1 \end{vmatrix} = 1/y. \quad (\text{BK})$$

Dès lors nous pouvons écrire, pour $(x, y) := p \in (\mathbb{R}_+^*)^2$, notant dp un voisinage infinitésimal de p :

$$\mathbb{P}(P \in dp) = \mathbb{P}(\Phi(M) \in dp) = \mathbb{P}(M \in \Phi^{-1}(dp)). \quad (\text{BL})$$

Mais $\Phi^{-1}(dp)$ est un voisinage infinitésimal de $\Phi^{-1}(p) := m$, dont le volume, d'après la formule du changement de variable, est $\text{vol}(dp) / y$ (y étant toujours positif, on n'a pas à utiliser de valeur absolue), et donc

$$\mathbb{P}(P \in dp) = \mathbf{1}_{m \in]0, 1]^2} y^{-1} \text{vol}(dp) / y = \mathbf{1}_{0 < z < y < 1} y^{-1} \text{vol}(dp) \quad (\text{BM})$$

Nous avons donc trouvé la densité de P ; or Z est la première composante de P , donc sa loi se trouve en intégrant selon les différentes valeurs possibles pour la seconde composante : pour $z \in \mathbb{R}_+^*$ et dz un voisinage infinitésimal de z ,

$$\begin{aligned} \mathbb{P}(Z \in dz) &= \int_{y \in \mathbb{R}_+^*} \mathbb{P}(Z \in dz \text{ et } Y \in dy) = \int_{y \in \mathbb{R}_+^*} \mathbb{P}(P \in dz \times dy) \\ &= \int_{y \in \mathbb{R}_+^*} \mathbf{1}_{0 < z < y < 1} y^{-1} \text{vol}(dz \times dy) = \mathbf{1}_{z > 0} \left(\int_{y \in]z, 1[} y^{-1} \text{vol}(dy) \right) \text{vol}(dz) \\ &= \mathbf{1}_{z \in]0, 1[} \left(\int_{y=z}^1 y^{-1} dy \right) \text{vol}(dz) = \mathbf{1}_{z \in]0, 1[} \ln(1/z) \text{vol}(dz) : \end{aligned}$$

nous avons donc trouvé la densité du produit XY . ♣

Changement de variables inverse

Pour finir, mentionnons un corolaire du théorème (AW) qui peut s'avérer utile dans certaines circonstances :

Corolaire (BN) (Changement de variable inverse). *Soient U et V deux ouverts de \mathbb{R}^d et soit $\varphi : U \rightarrow V$ une bijection de classe C^1 entre U et V , dont nous supposons (ce qui est presque automatique^[††]) que la bijection réciproque est elle aussi de classe C^1 ^[*]. Soit $x \in U$ et notons $\varphi(x) := y$, de sorte que $x = \varphi^{-1}(y)$. Alors le jacobien de φ en x est égal à l'inverse du jacobien de φ en y :*

$$\det \left(\left(\frac{\partial \varphi_i}{\partial x_j}(x) \right) \right)_{ij} = \det^{-1} \left(\left(\frac{\partial \varphi_i^{-1}}{\partial y_j}(y) \right) \right)_{ij}. \quad (\text{BO})$$

◇

[††]. Il faut et il suffit pour cela que le jacobien de φ ne s'annule nulle part

[*]. On dit alors que φ est un C^1 -difféomorphisme entre U et V .

Remarque (BP). En dimension 1, on retrouve la formule pour la dérivée de la réciproque : $\varphi'(x) = 1/(\varphi^{-1})'(y)$. ♣

Démonstration. En fait, la formule du changement de variables ne nous permet de vérifier l'égalité que modulo valeur absolue. L'idée est tout simplement de considérer un volume infinitésimal dx autour du point x ; la formule du changement de variables nous assure alors qu'on a

$$\text{vol}(\varphi(dx)) = |\mathbf{J}_\varphi(x)| \text{vol}(dx). \quad (\text{BQ})$$

On pose alors $y := \varphi(x)$ et $dy := \varphi(dx)$ (qui est un voisinage infinitésimal de y) : et puisque φ est une bijection, on a alors $\varphi^{-1}(dy) = dx$, et la formule ci-dessus devient donc

$$\text{vol}(dy) = |\mathbf{J}_\varphi(x)| \text{vol}(\varphi^{-1}(dx)). \quad (\text{BR})$$

Mais d'autre part la formule du changement de variables appliquée à φ^{-1} donne que

$$\text{vol}(\varphi^{-1}(dy)) = |\mathbf{J}_{\varphi^{-1}}(y)| \text{vol}(dx), \quad (\text{BS})$$

d'où en combinant les deux formules :

$$\text{vol}(dy) = |\mathbf{J}_\varphi(x)| |\mathbf{J}_{\varphi^{-1}}(y)| \text{vol}(dy), \quad (\text{BT})$$

dont on tire finalement que $|\mathbf{J}_{\varphi^{-1}}(y)| = |\mathbf{J}_\varphi(x)|^{-1}$ après simplification par $\text{vol}(dy)$.

En fait, il n'est pas très difficile de démontrer directement le corolaire (BN). En effet, soient i, j deux indices ; on peut calculer la dérivée de x_i par rapport à x_j en écrivant, grâce à la règle de dérivation des fonctions composées, que

$$\frac{\partial x_i}{\partial x_j} = \frac{\partial}{\partial x_j}(\varphi_i^{-1} \circ \varphi) = \sum_k \frac{\partial \varphi_i^{-1}}{\partial y_k} \frac{\partial \varphi_k}{\partial x_j}, \quad (\text{BU})$$

où l'on reconnaît l'entrée d'indices (i, j) du produit matriciel $\mathbf{J}_{\varphi^{-1}} \times \mathbf{J}_\varphi$ ($\mathbf{J}_{\varphi^{-1}}$ et \mathbf{J}_φ étant évaluées respectivement en $\varphi(x)$ et en x). Comme d'autre part il est évident que $\partial x_i / \partial x_j \equiv \mathbf{1}_{i=j}$, on en déduit finalement que les matrices $\mathbf{J}_{\varphi^{-1}}(y)$ et $\mathbf{J}_\varphi(x)$ sont inverses l'une de l'autre, et donc que leurs déterminants sont eux aussi inverses l'un de l'autre. ♣

Le corolaire (BN) peut s'avérer utile lorsque l'application φ^{-1} a une expression plus sympathique que l'application φ elle-même, ce qui se produit souvent lorsqu'on considère des mesures-images, puisque dans ce cas-là c'est l'application φ^{-1} qui définit la mesure-image, alors qu'on doit utiliser l'application φ pour la formule du changement de variables.

Exemple (BV). Si nous reprenons par exemple le problème de l'exemple (BI), on avait $\Phi : (x, y) \mapsto (xy, y)$. On pouvait alors observer que le jacobien de Φ est

$$\begin{vmatrix} \partial z / \partial x & \partial z / \partial y \\ \partial y / \partial x & \partial y / \partial y \end{vmatrix} = \begin{vmatrix} y & x \\ 0 & 1 \end{vmatrix} = y, \quad (\text{BW})$$

d'où on aurait pu déduire grâce au corolaire (BN) que le jacobien de Φ^{-1} en (z, y) valait y^{-1} —ce qui est effectivement ce que nous avons trouvé en explicitant la fonction Φ^{-1} . ♣

Chapitre 3

Statistiques descriptives des distributions de probabilité

Décrire une distribution de probabilité quelconque (disons sur un espace \mathcal{X}) est quelque chose de compliqué, puisqu'il faudrait donner les probabilités de *tous* les $A \subseteq \mathcal{X}$... Il est donc bien utile d'introduire des quantités synthétiques qui permettent, en quelques chiffres, de se faire une idée du comportement de cette distribution. Les quantiles, l'espérance, la variance et la covariance sont de telles quantités, que nous allons définir dans ce chapitre. Notez qu'il s'agit là de quantités concernant uniquement les distributions de probabilité sur \mathbb{R} (pour les quantiles, l'espérance et la variance) ou sur \mathbb{R}^2 (pour la covariance), ou sur un sous-ensemble de ces espaces : il faudra donc, pour utiliser ces outils, considérer des variables aléatoires exprimables *numériquement*. — Mais fort heureusement, les données se présentent le plus souvent naturellement sous forme numérique! ☺

3.1 Fonction de répartition

Remarque (BX). Ci-dessous nous allons présenter les fonctions de répartition (puis les quantiles et les intervalles de fluctuation) dans le cas de distributions de probabilités portées par \mathbb{R} (quitte à ce qu'elles ne prennent leurs valeurs que dans un certain sous-ensemble de \mathbb{R}); mais en fait, on peut tout aussi bien définir ces concepts en remplaçant \mathbb{R} par n'importe quel ensemble totalement ordonné. Un cas d'ensemble totalement ordonné qu'on peut rencontrer pas si rarement est celui où les valeurs $-\infty$ et/ou $+\infty$ sont autorisées pour notre distribution de probabilité, qui devient alors une distribution sur $[-\infty, +\infty]$. Un autre cas est celui où on compare des éléments selon un critère discret en éventuels les exæquos selon un autre critère (de sorte qu'on obtient alors p. ex. une distribution de probabilité portant sur $\mathbb{Z} \times \mathbb{Z}$ muni de l'ordre lexicographique). Mais en fait, de telles structures d'ordres sont isomorphes à des structures d'ordre sur des sous-ensembles de \mathbb{R} (par exemple, l'ordre sur $[-\infty, +\infty]$ est isomorphe à l'ordre sur un segment), de sorte qu'on peut toujours se ramener au cas de \mathbb{R} malgré tout. (Il existe certes des structures d'ordre total qu'il est impossible de plonger dans \mathbb{R} ; mais on peut démontrer qu'il n'est pas possible de faire des probabilités de façon pertinente sur de telles structures). On peut aussi formuler les propriétés des fonctions de répartition et de quantile directement dans le cadre général d'un ensemble totalement ordonné : les résultats sont essentiellement les mêmes, modulo quelques subtilités supplémentaires. ☺

Préliminaires : Limites de fonctions

Définition (BY). Si $f: \mathbb{R} \rightarrow \mathbb{R}$ est une fonction telle que, pour tout $x \in \mathbb{R}$, les valeurs $\lim_{x' \rightarrow x, x' < x} f(x')$ et $\lim_{x' \rightarrow x, x' > x} f(x')$ existent dans \mathbb{R} (une telle fonction est qualifiée de *làglàd*, pour « limite à gauche, limite à droite »), alors on note, pour tout $x \in \mathbb{R}$:

$$f(x-) := \lim_{x' \nearrow x} f(x'); \quad (\text{BZ})$$

$$f(x+) := \lim_{x' \searrow x} f(x'); \quad (\text{CA})$$

moralement, $f(x-)$ représente donc la valeur que f a “juste avant x ”, et $f(x+)$ représente la valeur qu’elle a “juste après x ”. \heartsuit

On notera en particulier que

Proposition (CB). *Les fonctions monotones sont làglàd.* \diamond

Fonction de répartition

! **Définition (CC).** Soit P une distribution de probabilité sur \mathbb{R} . On appelle *fonction de répartition de P* “l’”application $F: \mathbb{R} \rightarrow [0, 1]$ définie par :

$$\begin{cases} F(x-) := \mathbb{P}(P < x); \\ F(x+) := \mathbb{P}(P \leq x). \end{cases} \quad (\text{CD})$$

\heartsuit

Remarque (CE). La définition ci-dessus sous-entend qu’il existe bel et bien une fonction làglàd F vérifiant les propriétés énoncées, ce qui est effectivement vrai (mais demanderait tout de même à être vérifié en toute rigueur). \clubsuit

Remarque (CF). En fait, la définition ci-dessus ne caractérise pas entièrement F , car il existe plusieurs fonctions vérifiant ces limites à gauche et à droite : si on veut vraiment voir F comme une fonction, il faudra donc la “normaliser” en fixant un choix parmi les différentes possibilités compatibles. La technique la plus courante pour normaliser une fonction définie par ses limites à gauche et à droite consiste à lui imposer des conditions de continuité, soit à gauche, soit à droite. En l’occurrence, cela aboutit aux deux possibilités de normalisation suivantes :

- Si on impose en outre à F d’être continue à droite, cela revient à poser $F(x) := F(x+) = \mathbb{P}(P \leq x)$: c’est le choix qui est fait dans tous les ouvrages de référence, dès lors qu’on a besoin de normaliser F ;
- À l’inverse, si on imposait à F d’être continue à gauche, cela reviendrait à poser $F(x) := F(x-) = \mathbb{P}(P < x)$. Néanmoins cette normalisation n’est jamais utilisée en pratique.

Vous vous demandez peut-être pourquoi j’ai absolument tenu à définir F par ses limites à droite et à gauche, puisque je viens d’expliquer qu’il existe une convention bien établie pour définir les valeurs de F elles-mêmes. La raison est qu’en fait, dans un sens mathématique profond, la “bonne” façon de définir une fonction de répartition F est plutôt de considérer qu’il s’agit seulement d’une “quasi-fonction” définie uniquement par ses limites à gauche et à droite ; et dès lors, il est naturel de refuser de choisir une convention !

Cette façon de procéder, quoique peu banale, a plusieurs avantages. Déjà, on n’aura plus à s’embêter à retenir si le cas d’égalité doit être regroupé avec les valeurs inférieures ou avec les valeurs supérieures \smile Mais cela va en fait bien plus loin : en s’interdisant de donner une valeur à $F(x)$, les résultats obtenus seront beaucoup plus symétriques (et donc plus élégants et faciles à retenir) que si on avait fixé une convention arbitrairement ; et aussi, la vision de la fonction de quantile comme « réciproque » de la fonction de répartition deviendra alors pleinement rigoureuse ! (modulo adaptation de la notion de réciproque pour les « quasi-fonctions »). \clubsuit

Remarque (CG). Au passage, en vertu du fait que seules les limites à gauche et à droite de la fonction de répartition sont réellement bien définies, lorsqu'on trace la fonction de répartition d'une distribution de probabilité discrète, il apparaît plus pertinent de *tracer* des segments verticaux là où la fonction de répartition fait des sauts : l'inverse de ce qu'on vous a toujours appris à faire...! ♣

On a les propriétés évidentes suivantes^[*] :

Proposition (CH).

- (i) Une fonction de répartition est croissante.
 (ii) Si F est la fonction de répartition d'une distribution de probabilité P sur \mathbb{R} , on a :

$$\begin{aligned} F(x) &\underset{x \searrow -\infty}{\rightarrow} 0; & \text{(CI)} \\ F(x) &\underset{x \nearrow +\infty}{\rightarrow} 1. & \diamond \end{aligned}$$

- (iii) Si P est une distribution à densité, ou plus généralement, si P ne donne de masse à aucun singleton (on dit alors que P est diffuse ; dans le cas contraire, les x tels que $\mathbb{P}(P = x) > 0$ sont appelés des atomes de P), la fonction de répartition de P est continue.

3.2 Quantiles et intervalles de fluctuation

Quantiles

Les quantiles sont en quelque sorte le concept réciproque de la fonction de répartition :

Définition (CJ). Soit P une distribution de probabilité sur \mathbb{R} , ayant pour fonction de répartition F . On définit la fonction de quantile (notée ici $Q(\bullet)$) de la distribution P comme “la” fonction làglàd ayant les limites à gauche et à droite suivantes : pour $p \in]0, 1[$,

$$\begin{cases} Q(p-) := \sup\{x \in \mathbb{R} \mid F(x) < p\}; \\ Q(p+) := \inf\{x \in \mathbb{R} \mid F(x) > p\}. \end{cases} \quad \text{(CK)}$$

En pratique, dans la quasi-totalité des cas rencontrés, $Q(p-)$ et $Q(p+)$ coïncident^[†] : on notera alors $Q(p)$ leur valeur commune, sans se poser plus de questions. ♡

[*]. Pour donner un sens rigoureux à la proposition (CH), il faudrait que je définisse précisément ce qu'on entend par les notions de croissance, de limites ou la continuité pour une « quasi-fonction » définie uniquement par ses limites à gauche et à droite, ce qui serait assez fastidieux... En fait, il vous suffit de retenir que les propriétés de la proposition (CH) sont vraies aussi bien pour la version continue à droite de la fonction de répartition que pour sa version continue à gauche — et plus généralement, pour tout choix de normalisation respectant la propriété de croissance.

[†]. Pour les fonctions de répartition, le fait que $F(x-)$ et $F(x+)$ puissent être différentes était vraiment important en pratique, car il se produit systématiquement dès lors que F est une distribution de probabilité discrète et que x correspond à une valeur possible de P . Mais pour les fonctions de quantile, il faudrait choisir p “exprès” à certaines valeurs très précises (dépendant de P) pour observer des valeurs distinctes pour $Q(p-)$ et $Q(p+)$: or en général les valeurs de p sont fixées de façon indépendante de la distribution P ; et dès lors ces cas particuliers n'ont aucune raison de se produire.

Remarque (CL). Dans la définition ci-dessus, je n'ai pas précisé quelle normalisation de la fonction de répartition je prenais : c'est parce que ce choix n'a aucune incidence^[‡] sur la définition de la fonction de quantile. ♣

Remarque (CM). Si Q est la fonction de quantile d'une distribution de probabilité portée par \mathbb{R} , $Q(p-)$ et $Q(p+)$ sont toujours finis pour $p \in]0, 1[$. En fait, on pourrait aussi définir les quantiles pour $p = 0$ et $p = 1$: dans ce cas, les quantiles de rangs $0-$ et $1+$ valent systématiquement resp. $-\infty$ et $+\infty$, tandis que les quantiles de rangs respectifs $0+$ et $1-$ correspondent respectivement au minimum et au maximum du support de P , c'est-à-dire de la plage où P prend effectivement ses valeurs^[§]. Par exemple, pour une distribution exponentielle, les quantiles de rang $0+$ et $1-$ sont respectivement 0 et $+\infty$; pour une distribution uniforme sur $]a, b[$, ce sont resp. a et b ; et pour une distribution normale (de variance non nulle), ce sont resp. $-\infty$ et $+\infty$. ♣

! *Remarque (CN).* Informellement, le quantile de rang p est l'abscisse qui sépare la masse de P entre une probabilité p à gauche, et $1 - p$ à droite : plus précisément, le quantile de rang p est le moment où s'effectue la transition entre le moment où on avait une probabilité inférieure à p de tomber à gauche de x , et celui où on se met à avoir une probabilité supérieure à p de tomber à gauche de x . ♣

Remarque (CO). Dans le contexte de la remarque ci-dessus, les notions de quantiles de rang $p-$ et $p+$ correspondent aux cas exceptionnels où il existe toute une plage de valeurs de x pour lesquelles on a une probabilité exactement p de tomber à gauche et une probabilité exactement $1 - p$ de tomber à droite : dans ce cas, l'extrémité gauche de cette plage correspond au quantile de rang $p-$, tandis que l'extrémité droite de la plage correspond au quantile de rang $p+$. ♣

Remarque (CP). Contrairement au cas des fonctions de répartition, où il existe une façon conventionnelle de définir $F(x)$ lorsque $F(x-)$ et $F(x+)$ diffèrent, il n'existe pas vraiment de convention universelle pour définir $Q(p)$ lorsque $Q(p-)$ et $Q(p+)$ diffèrent... Le choix le plus souvent fait est celui consistant à imposer à Q d'être continue à gauche, ce qui revient à poser $Q(p) := Q(p-)$. Mais une autre convention courante est de poser $Q(p) := (Q(p-) + Q(p+)) / 2$; et il y a aussi des auteurs qui définissent $Q(p) := Q(p-)$, *sauf* pour $p = 0$ où ils posent $Q(0) := Q(0+)$: bref, c'est un peu la pagaille... En fait, chacune de ces conventions présente ses avantages et ses inconvénients : car, de même que dans le cas des fonctions de répartition, le choix mathématiquement le plus pertinent est de ne *pas* fixer de convention et de ne définir $Q(\bullet)$ qu'à travers ses limites à gauche et à droite... ! Dans la suite, c'est en substance le choix que nous ferons : plus précisément, les énoncés sur « la » fonction de quantile seront tournés de façon à être valables pour toutes les façons possibles de normaliser la fonction de quantile^[¶]. ♣

! **Proposition (CQ).** *La fonction quantile est la « réciproque » de la fonction de répartition, au sens suivant : si on dessine le graphe de la fonction de répartition de P en traçant des segments verticaux pour les sauts éventuels de $F(\bullet)$, et qu'on échange les axes des abscisses et des ordonnées, alors on obtient le graphe de la fonction de quantile de P (là aussi, avec des segments verticaux pour les sauts éventuels de $Q(\bullet)$).* ◇

[‡]. Du moins sous réserve que la normalisation de $F(\bullet)$ soit prise croissante, ce qu'on supposera toujours.

[§]. Techniquement parlant, le *support* d'une distribution de probabilité sur un espace \mathcal{X} muni d'une topologie (polonaise), noté $\text{supp } P$, est défini comme le plus petit (au sens de l'inclusion) sous-ensemble fermé de \mathcal{X} dont la probabilité vaille 1 : on peut démontrer qu'un tel plus petit fermé existe bel et bien.

[¶]. Plus précisément, pour toutes les normalisations aboutissant à une fonction de quantile croissante : mais ce sera le cas pour tout choix raisonnable, incluant en particulier les différentes conventions mentionnées plus haut.

Proposition (CR). *La fonction de quantile est toujours croissante, et elle est strictement croissante si et seulement si la distribution P est diffuse (i.e. n'accorde de masse non nulle à aucun point). Plus généralement, les plateaux de la fonction de quantile correspondent aux sauts de la fonction de répartition, i.e. aux atomes de P , tandis que lessauts de Q correspondent aux plateaux de la fonction de répartition, c.-à-d. aux intervalles d'intérieur non vide auxquels P n'attribue aucune masse. \diamond*

Proposition (CS). *Soit P une distribution de probabilité sur \mathbb{R} ayant resp. pour fonction de répartition F et pour fonction de quantile Q . Alors, pour tous $x \in \mathbb{R}$, $p \in [0, 1]$, on a les équivalences suivantes :*

$$F(x+) < p \iff x < Q(p-); \tag{CT}$$

$$F(x-) \leq p \iff x \leq Q(p+). \tag{CU}$$

\diamond

Remarque (CV). Bien que j'aie indiqué cette proposition comme étant « à retenir », il vaut beaucoup mieux se contenter de retenir par cœur qu'il existe des relations qui expriment, moralement, le fait que se trouver avant le quantile de rang p est synonyme d'avoir une fonction de répartition plus petite que p (ce qui est justement l'idée derrière la définition des quantiles!), et s'aider de dessins pour retrouver un énoncé précis correct lorsqu'on a besoin de mettre des '+' et des '-' et/ou de savoir si les inégalités sont larges ou strictes... \clubsuit

Remarque (CW). On peut aussi essayer de retenir la proposition de la façon suivante : la première équivalence correspond aux cas où x est "carrément" avant le quantile de rang p : dans ce cas, on le compare au quantile de rang $p-$ (qui sera plus petit, en général, que le quantile de rang $p+$), et on impose que l'inégalité soit stricte ; et du côté des fonctions de répartition, on demande là encore que l'inégalité soit stricte, et on considère la valeur en $x+$ (qui est, en général, plus grande qu'en $x-$). Alors qu'à l'inverse, la seconde équivalence correspond à avoir un maximum de "tolérance" quant à l'affirmation de x est situé avant le quantile de rang p , resp. que la fonction de répartition en x est plus petite que p ... \clubsuit

Définition (CX) (q -Quantiles, Médiane). Pour d entier et $n \in]0, q[$, ce qu'on appelle le « n -ième d -quantile » de la distribution P est son quantile de rang nd . Pour d égal resp. à 3, 4, 5, 10 et 100, les d -quantiles sont appelés resp. « terciles », « quartiles », « quintiles », « déciles » et « centiles ». Pour d égal à 2, le quantile de rang $1/2$ est appelé *médiane* de la distribution P . \heartsuit

Théorème (CY). *Soient P une distribution de probabilité sur \mathbb{R} et $p \in]0, 1[$; et notons x le quantile de rang p de P . Soit $f: \mathbb{R} \rightarrow \mathbb{R}$ une fonction continue strictement croissante (resp. strictement décroissante). Alors le quantile de rang p (resp. de rang $(1-p)^{[**]}$) de $f_* P$ est égal à $f(x)$. \diamond*

Remarque (CZ). Dans certains cas, il peut être plus commode de "renverser l'ordre sur \mathbb{R} " : on définit alors la *fonction de répartition complémentaire* comme la probabilité d'être plus *grand* qu'une certaine valeur, resp. la *fonction de quantile complémentaire* comme la fonction qui, à p , associe un point séparant une probabilité p à sa droite d'une probabilité $1-p$ à sa gauche. Si nous notons $\check{F}(\bullet)$ et $\check{Q}(\bullet)$ ces fonctions complémentaires respectives, on a tout simplement $\check{F}(x) = 1 - F(x)$, resp. $\check{Q}(p) = Q(1-p)$. \clubsuit

[[]]. Ici on suppose que p est en fait de la forme « $a-$ » ou « $a+$ ».

[**]. En posant $1 - (a-) = (1-a)+$ et $(1 - (a+)) = (1-a)-$.

Intervalles de fluctuation

Définition (DA). Pour P une distribution de probabilité sur (un sous-ensemble de) \mathbb{R} , $\alpha \in]0, 1[$, un intervalle $I \subseteq \mathbb{R}$ est qualifié d'*intervalle de fluctuation* au niveau de risque α (ou « au niveau de confiance $1 - \alpha$ ») lorsque P donne une masse au moins égale à $1 - \alpha$ à l'intervalle I :

$$\mathbb{P}(P \notin I) \leq \alpha. \quad (\text{DB})$$

♥

Remarque (DC). Pour les applications pratiques, le niveau de risque α sera toujours pris $\leq 1/2$ ^[††]. C'est pourquoi il arrivera souvent qu'on dise simplement « au niveau [tant] », sans prendre la peine de préciser si on est en train de parler d'un niveau de *risque* ou d'un niveau de *confiance* : il suffira en effet de comparer la valeur du niveau à $1/2$ pour trancher ! Par exemple, un « niveau 10 % » se référera implicitement à un niveau de risque (soit $\alpha = 0,1$ en l'occurrence), tandis qu'un « niveau 95 % » se référera à un niveau de confiance (soit $\alpha = 0,05$ en l'occurrence). (Et dans le cas particulier d'un niveau $1/2$, « risque à 50 % » et « confiance à 50 % » sont en fait synonymes ! ☺). ♣

Remarque (DD). La définition pourrait aussi s'appliquer sans exiger que I soit un intervalle : dans un tel cas, on parlerait plutôt de « zone » de fluctuation. En particulier, le concept de zone de fluctuation garde du sens même lorsque l'ensemble mesurable sur lequel P porte n'est pas totalement ordonné. Mais dans ce cours nous nous limiterons aux intervalles sur \mathbb{R} , qui sont à la fois le cas le plus important et le plus fréquemment rencontré. ♣

Lemme (DE). Un intervalle de la forme $] -\infty, q]$ est un intervalle de fluctuation au niveau α pour P si et seulement si $\text{Répart}(P; q+) \geq 1 - \alpha$; ou, de manière équivalente, si et seulement si $q \geq \text{Qtile}(P; (1 - \alpha)-)$. Par conséquent, le plus petit intervalle de fluctuation au niveau α pour P ayant cette forme (il y en a bien un) est l'intervalle $] -\infty, \text{Qtile}(P; (1 - \alpha)-)]$.

De même, le plus petit intervalle de la forme $[q, \infty[$ qui soit un intervalle de fluctuation au niveau α pour P est l'intervalle $[\text{Qtile}(P; \alpha+), +\infty[$. ◇

Le lemme (DE) ci-dessus suggère la définition suivante :

! **Définition (DF).** Pour P une distribution de probabilité sur \mathbb{R} , $\alpha \in]0, 1/2]$, on parle de « l' »intervalle de fluctuation au niveau de risque α (ou au niveau de confiance $1 - \alpha$) de la distribution P pour désigner l'intervalle fermé compris entre ses quantiles de niveaux respectifs $(\alpha/2)+$ et $(1 - \alpha/2)-$ ^[‡‡], autrement dit, l'intervalle

$$[\text{Qtile}(P; \alpha/2), \text{Qtile}(P; 1 - \alpha/2)]. \quad (\text{DG})$$

En vertu du lemme (DE), cet intervalle constitue effectivement un intervalle de fluctuation : on le qualifie de « l' »intervalle (même si ce n'est pas le seul) dans la mesure où sa construction est dans un sens « canonique ». ♥

[††]. Il y aura une seule exception dans ce cours, qui sera lorsque nous introduirons le concept de *p-valeur* au chapitre 13 : et encore, l'utilisation de niveaux α plus grands que $1/2$ à ce moment-là servira uniquement à *expliquer* ce que signifie la *p-valeur* ; mais dans tous les cas, vous ne serez jamais amenés à parler d'« intervalle à 75 % de risque » ! ☺

[‡‡]. En pratique, on ne s'ennuiera pas avec les fioritures '+' et '-', confer la remarque figurant à la fin de la définition (CJ).

Remarque (DH). Pour que la formule ci-dessus soit toujours valide, il est très important de prendre des bornes *fermées* à l'intervalle : en effet, dès lors que P est une mesure discrète (ou, de manière générale, qu'il y a une masse non nulle aux bornes de l'intervalle de fluctuation canonique), la formule devient le plus souvent fausse... ! C'est pourquoi, de manière générale, nous écrirons systématiquement les intervalles de fluctuation (et aussi, du même coup, les intervalles de confiance et de prédiction) sous forme fermée. ♣

Remarque (DI). En fait, du point de vue mathématique, la construction la plus pertinente pour construire un intervalle de fluctuation canonique consiste à définir « l' » intervalle de fluctuation (au niveau de risque α) pour P comme l'ensemble des x tels que

$$\alpha/2 < \text{Répart}(P; x+) \quad \text{et} \quad \text{Répart}(P; x-) < 1 - \alpha/2. \quad (\text{DJ})$$

Cette construction, que je qualifierai au besoin d'« intervalle canonique *strict* », donne presque le même intervalle que la construction de la définition (DF), à ceci près que l'intervalle qu'elle fournit pourra, dans certains cas, être ouvert au niveau d'une ou deux de ses bornes... Du point de vue pratique, néanmoins, ce raffinement n'apporterait essentiellement rien : en effet, on pourrait montrer que, dans les cas où l'intervalle strict diffère de l'intervalle « fermé » de la définition (DF), la probabilité que la loi P tombe précisément au niveau d'un des points qui diffèrent entre ces deux intervalles... est nulle ! Dès lors, dans la mesure où le seul véritable intérêt de la construction de l'intervalle de fluctuation strict est de rendre certains résultats mathématiques plus rigoureux, autant ne pas s'encombrer avec cette subtilité, et en rester à la plus simple définition (DF) : et ce, d'autant plus que prendre l'habitude d'écrire tous les intervalles de fluctuation sous forme fermée évitera de risquer de commettre certaines erreurs fâcheuses ! ☺ (cf. remarque (DH) ci-dessus). ♣

Remarque (DK). Notez que je n'ai pas affirmé ci-dessus que l'intervalle de fluctuation canonique serait minimal parmi les intervalles (fermés) de fluctuation de niveau α pour P , car ce n'est pas vrai en général. Par exemple, si P est la distribution de probabilité uniforme sur $\{-1, 0, +1\}$, alors l'intervalle de fluctuation (canonique) de niveau 40 % pour P est $[-1, +1]$, alors que pourtant il existe d'autres intervalles de fluctuation strictement plus petits à ce même niveau : notamment $[-1, 0]$ et $[0, +1]$. ♣

Remarque (DL). On peut se demander ce qui se passe si l'on prend $\alpha = 0$ dans la définition des intervalles de fluctuation, autrement dit, qu'on souhaite trouver un intervalle contenant *toute* la masse de P : dans ce cas-là, plutôt que de parler d'un « intervalle de fluctuation à 100 % de confiance », on dit que l'intervalle I « supporte » la distribution P . En fait, dans la cas $\alpha = 0$, toutes les formules évoquées restent parfaitement valables, avec néanmoins la possibilité que certains quantiles deviennent infinis. Ce cas n'est néanmoins pas très intéressant du point de vue de la statistique, puisque la philosophie même de la statistique est qu'on renonce à avoir des certitudes *absolues* en vue d'obtenir des informations plus intéressantes...

Il existe toujours un plus petit intervalle (non nécessairement fermé) supportant la mesure P , qu'on appelle parfois l'*étendue* de P .

Plus intéressant, si on considère plus généralement les *zones* (pas forcément des intervalles) supportant P , on peut démontrer qu'il existe toujours un plus petite zone (au sens de l'inclusion) *fermée* supportant P , qu'on appelle alors le *support* de P . Le support de P a ceci d'intéressant qu'on peut continuer à le définir y compris lorsque la mesure P porte sur un autre espace que \mathbb{R} : il suffit juste qu'il s'agisse d'un espace muni d'une structure topologique « polonaise », c'est-à-dire équivalente (du point de vue des ouverts qu'elle définit) à une topologie métrique complète et séparable. ♣

3.3 Mode

Définition (DM) (Mode d'une distribution de probabilité discrète). Pour P une distribution de probabilité sur un espace discret \mathcal{X} , le *mode* de P , que je noterai

!

parfois $\text{Mode}(P)$, est la valeur de x ^[*] pour laquelle $P(\{x\})$ est maximale : autrement dit, la valeur qu'il y a le plus de chance qu'une v.a. de loi P prenne. \heartsuit

Proposition (DN). Si P est une distribution de probabilité sur l'espace discret \mathcal{X} et que $f: \mathcal{X} \rightarrow \mathcal{Y}$ est une application injective, alors

$$\text{Mode}(f_* P) = f(\text{Mode}(P)) : \quad (\text{DO})$$

autrement dit, le mode d'une distribution discrète est préservé par la mesure-image (injective). \diamond

On parle aussi de mode pour les distributions à densité, mais c'est plutôt par abus de langage :

Définition (DP) (Mode d'une distribution de probabilité à densité). Pour P une distribution de probabilité à densité sur un espace \mathbb{R}^d , on appelle par extension « mode » de P la valeur x pour laquelle la densité $P(dx)/\text{vol}_d(dx)$ de P par rapport à la mesure de Lebesgue est maximale. Informellement, cela correspond à la valeur pour laquelle il y a le plus de chance qu'une v.a. de loi P tombe “juste à côté”. \heartsuit

Remarque (DQ). Attention, il n'y a pas d'analogie de la proposition (DN) pour les distributions de probabilité à densité : le mode d'une distribution à densité n'est pas préservé par reparamétrisation en général! \clubsuit

3.4 Quantités empiriques

Définition (DR) (Distribution empirique). Pour x_1, \dots, x_n ($n \geq 1$) des données à valeurs dans un espace \mathcal{X} , la *distribution empirique* des x_i (ou plus exactement, du multi-ensemble $\{\{x_1, \dots, x_n\}\}$) est la distribution de probabilité P sur \mathcal{X} définie par

$$\forall A \subseteq \mathcal{X} \quad P(A) := \frac{\sum_{i=1}^n \mathbf{1}_{x_i \in A}}{n}; \quad (\text{DS})$$

ou autrement dit,

$$P := \sum_{i=1}^n \delta_{x_i}. \quad (\text{DT})$$

Remarque (DU). Il existe également une version *pondérée* des la notion de distribution empirique : si la i -ième donnée x_i est attachée au poids α_i , où les poids sont positifs et non identiquement nuls, la *distribution empirique pondérée* de $(x_i)_{1 \leq i \leq n}$ pour les poids en question sera la distribution de probabilité

$$\frac{\sum_{i=1}^n \alpha_i \delta_{x_i}}{\sum_{i=1}^n \alpha_i}. \quad (\text{DV})$$

Cependant nous n'aurons pas besoin d'utiliser les versions pondérées dans le cadre de ce cours. \clubsuit

Définition (DW) (Quantité empirique). Pour « schmilblick » un concept s'appliquant aux distributions de probabilité sur \mathcal{X} (fonction de répartition, quantiles, espérance, variance, covariance, ...), et x_1, \dots, x_n des valeurs de \mathcal{X} , le *schmilblick empirique* de $\{\{x_1, \dots, x_n\}\}$ est tout simplement le schmilblick de la distribution empirique de $\{\{x_1, \dots, x_n\}\}$. Dans le cas de l'espérance, l'espérance empirique est plus simplement appelée *moyenne*. \heartsuit

[*]. Que nous supposons ici unique pour simplifier, ce qui est le cas en général \heartsuit

Définition (DX). Dans ce cours, nous noterons resp. $\text{moy}(x_1, \dots, x_n)$, $\text{var}_{\text{emp}}(x_1, \dots, x_n)$, $\text{var}_{\text{emp}}^{1/2}(x_1, \dots, x_n)$, $\text{cov}_{\text{emp}}((x_1, y_1), \dots, (x_n, y_n))$ les moyenne, variance empirique, écart-type empirique et covariance empirique des valeurs considérées. Notez que nous utilisons des minuscules lorsque nous nous référons à des quantités empiriques, qui s'appliquent donc à un *jeu de données*, et des majuscules lorsque nous nous référons aux concepts sous-jacents, qui s'appliquent donc à une *distribution de probabilité*. \heartsuit

Remarque (DY). Le fait que nous mentionnions un indice « emp » pour la variance, l'écart-type et la covariance empiriques sert à ne risquer aucune confusion avec les variantes de ces concepts prenant en compte la *correction de Bessel* (cf. annexe ??). Dans ce cours, on pourra se dispenser de cet indice, qui sera sous-entendu par défaut ; mais il faut bien garder à l'esprit que d'autres sources, et le logiciel *R* en particulier, considèrent au contraire que le concept par défaut est celui incluant la correction de Bessel. \clubsuit

Proposition (DZ). On a en particulier :

$$\text{moy}(x_1, \dots, x_n) := \frac{x_1 + \dots + x_n}{n} ; \quad (\text{EA})$$

$$\text{var}_{\text{emp}}(x_1, \dots, x_n) := \text{moy}(x_i - \text{moy}(x_j)_{1 \leq j \leq n})_{1 \leq i \leq n}^{[\ddagger]} \quad (\text{EB})$$

$$= \text{moy}(x_i^2)_{1 \leq i \leq n} - (\text{moy}(x_i)_{1 \leq i \leq n})^2. \quad (\text{EC})$$

\diamond

Remarque (EE). Il est essentiel de bien comprendre que les variance empirique & C^{ie} ne sont *pas* des quantités aléatoires : $\text{var}_{\text{emp}}(x_1, \dots, x_n)$ est simplement une *fonction* des données^[\ddagger], tout ce qu'il y a de plus déterministe : par exemple, $\text{var}_{\text{emp}}^{1/2}(2, 4, 10, 4) = 3$, fin de la discussion ! Attention toutefois : le fait est qu'on considère souvent les variances empiriques dans le cas où X_1, \dots, X_n sont des (réalisations de) variables aléatoires ; et dans ce cas évidemment $\text{var}_{\text{emp}}(X_1, \dots, X_n)$ sera une variable aléatoire, mais c'est simplement parce que c'est une fonction (déterministe) de variables aléatoires... ! \clubsuit

3.5 Théorie de la décision optimale

Motivation

Le but de ce chapitre est d'introduire une notion probabiliste conceptuellement très importante en ingénierie : la prise de décision optimale en contexte incertain.

Concernant la notion de « prise de décision optimale », son lien avec l'ingénierie est évident : si la mission de votre équipe d'ingénieurs, par exemple, est de concevoir

[†]. Soit, explicitement :

$$\frac{1}{n} \sum_{i=1}^n \left(x_i - \frac{\sum_{j=1}^n x_j}{n} \right)^2. \quad (\text{ED})$$

[‡]. Notons au passage qu'il y a un léger abus de langage à noter de la même façon la variance empirique quel que soit le nombre de données auquel elle s'applique : en réalité bien sûr, il y a une fonction $\mathbb{R} \rightarrow \mathbb{R}$ quand on a un jeu de 1 donnée, une fonction de $\mathbb{R}^2 \rightarrow \mathbb{R}$ quand on a un jeu de 2 données, &c.

un viaduc, vous allez devoir décider du type d'ouvrage, des dimensions de la voie, du nombre de piles ou de câbles, des matériaux utilisés, etc., pour optimiser un certain ensemble de critères ou de contraintes : cout de l'ouvrage, sécurité, capacité de la voie, durabilité, facilité d'entretien, considérations écologiques, etc. Vous avez déjà vu, dans vos cours du semestre 5 (et certains d'entre vous étudieront le sujet plus à fond ultérieurement), comment trouver le jeu de paramètres permettant d'optimiser une fonction donnée.

Cependant, dans de nombreux cas, la décision à prendre participe d'un pari sur l'avenir... C'est, en fait, une situation très fréquente dans la vie quotidienne, même lorsqu'on n'est pas ingénieur !

Mettons, par exemple, que vous soyez une étudiante (hétérosexuelle, célibataire) éprise du séduisant Quentin. Pour avoir une chance de le convaincre de vous choisir pour partenaire, vous devez vous faire mieux connaître de lui en lui consacrant du temps, par exemple en se rendant à cette soirée où vous savez qu'il sera présent. Si Quentin s'avère effectivement intéressé par vous, le gain potentiel est fabuleux, représentant des décennies de bonheur à deux ! Sauf que, si ça se trouve, même en vous présentant sous votre meilleur jour, Quentin n'aura pas ce "feeling" qui l'attirera vers vous... La seule chose qui dépende de vous, c'est le fait d'aller ou non à cette soirée ; mais le résultat de cette décision dépendra des préférences intimes de Quentin, que vous ne pouvez pas connaître au moment de prendre votre décision ! Vous devez donc décider du plan de votre soirée sans savoir de quoi il en retourne : on est ainsi dans une situation de prise de décision en contexte incertain.

Formalisation du problème

Formellement, la question est de prendre une certaine décision d , à valeurs dans un espace noté \mathcal{D} , sachant que le résultat de cette décision dépendra aussi d'une certaine quantité G , à valeurs dans un espace noté \mathcal{G} , que l'on ne connaît pas encore au moment de prendre la décision, et dont la valeur future sera aléatoire : l'enjeu étant que notre décision ayant des conséquences aussi « bonnes » que possible.

Pour pouvoir procéder à un traitement mathématique, nous avons donc besoin de *quantifier* en quoi une conséquence peut s'avérer plus ou moins « bonne » ! L'outil approprié pour ce faire est celui de *fonction de perte* :

Définition (EF). Dans le contexte de ce chapitre, une *fonction de perte* est une application $\ell : \mathcal{G} \times \mathcal{D} \rightarrow \mathbb{R}$ dont l'interprétation est la suivante :

- (i) La valeur $\ell(g, \hat{d})$ nous dit à quel point est-ce que la décision \hat{d} amènera des conséquences nuisibles si la quantité d'intérêt s'avère valoir g : plus cette valeur est petite, moins les conséquences sont nuisibles à nos yeux ;
- (ii) En outre, entre quatre valeurs de pertes possibles $\ell_0, \ell'_0, \ell_1, \ell'_1$ avec $\ell_0 < \ell'_0$ et $\ell_1 < \ell'_1$, le ratio $(\ell'_1 - \ell_1) / (\ell'_0 - \ell_0)$ nous dit *de combien de fois pire* est le fait de passer de la situation de perte ℓ_1 à la situation de perte ℓ'_1 , par rapport au fait de passer de la situation de perte ℓ_0 à la situation de perte ℓ'_0 . ♡

Remarque (EG). L'opposé d'une fonction de perte est qualifié de *fonction d'utilité* (traditionnellement notée ' u ') : c'est exactement le même principe, sauf que cette fois-ci nous avons une préférence pour les valeurs les plus *grandes* ! Les économistes travaillent plus volontiers avec des fonctions d'utilité ; mais dans le cadre de la statistique, il est préférable d'utiliser des fonctions de perte. ♣

Remarque (EH). Le choix d'une fonction de perte comprend une composante fondamentalement *subjective*. En effet, la notion de « conséquence nuisible » inclut la prise en compte de plusieurs aspects qui ne sont pas objectivement commensurables : considérations financières, aspects éthiques (conditions de travail des employés ; impact sur la santé des consommateurs ; possibilités de nuisances démocratiques ; ...), dégradation de l'environnement, ... La fonction de perte intégrant tous ces aspects au sein d'un même nombre, la pondération de ces différents aspects reflètera de l'échelle de valeur de celui qui choisira la fonction de perte! ♣

Décision optimale

Une fois le contexte formalisé, la notion de décision optimale en contexte incertain devient très simple :

Principe (EI). *Dans un contexte où la quantité G , inconnue au moment de prendre la décision, sera réalisé selon une loi aléatoire P , une décision d est d'autant meilleure que l'espérance*

$$\mathbb{E}^{g \sim P}(\ell(d, g)) \tag{EJ}$$

est petite. ◇

Chapitre 4

Conditionnement et indépendance

4.1 Conditionnement

Définition (EK) (Conditionnement par rapport à un évènement de probabilité non nulle). Si A est un évènement de probabilité non nulle, la probabilité *conditionnée à la réalisation de A* , notée $\mathbb{P}(\bullet \mid A)$, est la distribution de probabilité sur Ω telle que pour tout évènement $B \subseteq \Omega$:

$$\mathbb{P}(B \mid A) = \frac{\mathbb{P}(B \cap A)}{\mathbb{P}(A)}. \quad (\text{EL})$$

(Théorème : il s'agit bien d'une distribution de probabilité). Les espérance, variance, loi etc. sous cette probabilité conditionnelle sont alors notées resp. $\mathbb{E}(\bullet \mid A)$, $\text{Var}(\bullet \mid A)$, $\text{Loi}(\bullet \mid A)$, etc. ♥

Remarque (EM). Intuitivement, $\mathbb{P}(\bullet \mid A)$ décrit comment se répartissent *entre elles* les différentes possibilités pour lesquelles A est réalisé. ♣

Définition (EN) (Conditionnement par rapport à un évènement de probabilité nulle). Lorsque A est un évènement de probabilité nulle, on peut, en général, définir *quand même* $\mathbb{P}(\bullet \mid A)$ par passage à la limite^[*]. Typiquement, si X est une variable aléatoire à valeurs dans \mathbb{R}^d et x un élément de \mathbb{R}^d , on assimilera

$$\mathbb{P}(B \mid X = x) \quad (\text{EO})$$

à

$$\mathbb{P}(B \mid X \in dx), \quad (\text{EP})$$

pour dx un voisinage infinitésimal de x . Sous certaines hypothèses de régularité dont nous ne préoccuperons pas ici, on peut montrer que cette distribution de probabilité a bien une limite lorsqu'on fait tendre l'intervalle infinitésimal dx vers le singleton $\{x\}$, et que cette limite est bien une distribution de probabilité sur Ω ; en outre, sous des hypothèses assez générales, cette distribution de probabilité conditionnelle obtenue ainsi par passage à la limite peut être définie pour (presque-)toute valeur de x . ♥

Théorème (EQ). Soient X et Y deux variables aléatoires resp. à valeurs dans \mathbb{R}^n et \mathbb{R}^m , dont la loi jointe a une densité $f(x, y)$ par rapport à la mesure de Lebesgue sur \mathbb{R}^{n+m} (x désignant un point de \mathbb{R}^n et y un point de \mathbb{R}^m). Alors la loi de Y sachant « $X = x_0$ » est à densité aussi, et cette densité est donnée par :

$$\mathbb{P}(Y \in dy \mid X = x_0) = \frac{f(x_0, y)}{\int_{y' \in \mathbb{R}^m} f(x_0, y') \text{vol}_m(dy')} \text{vol}_m(dy). \quad (\text{ER})$$

◇

Démonstration. Pour dx_0 un voisinage infinitésimal de x_0 dans \mathbb{R}^n , on écrit que $\mathbb{P}(Y \in dy \mid X = x_0)$ est formellement la même chose que $\mathbb{P}(Y \in dy \mid X \in dx_0)$, qui vaut par définition de la probabilité conditionnelle :

$$\frac{\mathbb{P}(Y \in dy \text{ et } X \in dx_0)}{\mathbb{P}(X \in dx_0)}. \quad (\text{ES})$$

Notons que « $Y \in dy$ et $X \in dx_0$ » peut se réécrire comme « $(X, Y) \in dx_0 \times dy$ », où $dx_0 \times dy$ est un voisinage infinitésimal de (x_0, y) dans \mathbb{R}^{n+m} . Par conséquent, par définition de la densité, le numérateur de notre fraction vaut

$$f(x_0, y) \text{vol}(dx_0 \times dy), \quad (\text{ET})$$

où $\text{vol}(dx_0 \times dy) = \text{vol}(dx_0) \times \text{vol}(dy)$. Concernant le dénominateur, partitionnons \mathbb{R}^m en un très grand nombre de zones de dimensions infinitésimales, de la forme dy' , où y' décrit un maillage très dense de \mathbb{R}^m . On a alors

$$dx_0 \times \mathbb{R}^m = \bigsqcup_{y'} dx_0 \times dy', \quad (\text{EU})$$

d'où

$$\mathbb{P}(X \in dx_0) = \sum_{y'} \mathbb{P}(X \in dx_0 \text{ et } Y \in dy'). \quad (\text{EV})$$

Mais, d'après le travail que nous avons fait pour le numérateur,

$$\mathbb{P}(X \in dx_0 \text{ et } Y \in dy') = f(x_0, y') \text{vol}(dx_0) \text{vol}(dy'), \quad (\text{EW})$$

d'où

$$\mathbb{P}(X \in dx_0) = \text{vol}(dx_0) \times \left(\sum_{y'} f(x_0, y') \text{vol}(dy') \right) = \left(\int_{y'} f(x_0, y') dy' \right) \text{vol}(dx_0), \quad (\text{EX})$$

où la dernière égalité n'est autre que la définition informelle de l'intégrale (multi-dimensionnelle). Le résultat vient finalement après simplification du numérateur et du dénominateur par $\text{vol}(dx_0)$ (ce qui était attendu, puisque ce dx_0 n'était qu'un artifice technique). ◇

!

Théorème (EY). Soient μ une distribution de probabilité sur \mathcal{X} et, pour tout $x \in \mathcal{X}$, ν_x une distribution de probabilité sur \mathcal{Y} . Alors il existe une unique distribution de probabilité π sur $\mathcal{X} \times \mathcal{Y}$ telle que, si (X, Y) a pour loi (jointe) π , X suit la loi μ et, pour tout $x \in \mathcal{X}$, Y suit la loi ν_x sous $\mathbb{P}(\bullet \mid X = x)$. En particulier, si $\mathcal{X} = \mathbb{R}^n$, $\mathcal{Y} = \mathbb{R}^m$, que μ est la distribution de densité $f(x)$ et que ν_x est la distribution de densité $g_x(y)$, π sera la distribution de densité $f(x)g_x(y)$ sur \mathbb{R}^{n+m} . ◇

[*]. La définition rigoureuse est beaucoup plus technique et subtile ; nous ne nous en soucions pas ici, mais garderons juste à la tête que la définition ci-dessus pourrait poser problème en toute généralité, mais pas dans les cas que nous rencontrerons en pratique.

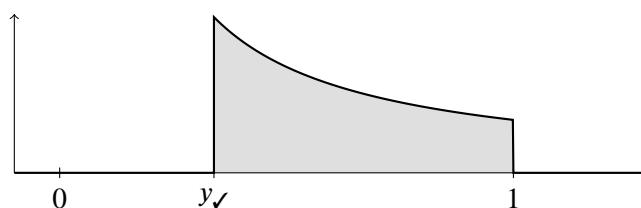


FIGURE 4.1 – Représentation de la X sachant la valeur de Y (on a pris ici $y_0 = 0,34$) dans l'exemple (EZ).

Exemple (EZ). À titre d'exemple, considérons la situation suivante. On tire une variable aléatoire $X(\omega)$ uniformément entre 0 et 1, puis on tire une variable Y uniformément entre 0 et $X(\omega)$. La première question est, quelle est la loi jointe du couple (X, Y) ? Puisque nous avons décrit la loi de Y *une fois connue la valeur de X* , cela revient à dire qu'on connaît la loi conditionnelle de Y sachant la valeur de X : plus précisément, pour tout $x \in [0, 1]$,

$$\text{Loi}(Y \mid X = x) = \text{Unif}^{\text{me}}(0, x). \quad (\text{FA})$$

Or la loi $\text{Unif}^{\text{me}}(0, 1)$ a pour densité $x \mapsto \mathbf{1}_{x \in [0, 1]}$, tandis que la loi $\text{Unif}^{\text{me}}(0, x)$ a pour densité $y \mapsto \mathbf{1}_{y \in [0, x]}x^{-1}$. D'après le théorème (EY), on en déduit que la densité jointe de (X, Y) est

$$(x, y) \mapsto \mathbf{1}_{x \in [0, 1]} \times \mathbf{1}_{y \in [0, x]}x^{-1} = \mathbf{1}_{0 \leq y \leq x \leq 1}x^{-1}. \quad (\text{FB})$$

Continuons notre analyse et supposons qu'on nous ait révélé la valeur de Y , *mais pas celle de X* , et qu'on se demande alors ce que peut bien valoir X compte tenu de cette information : autrement dit, on cherche la loi de X *sachant que $Y = y_0$* , où y_0 est la valeur effectivement observée pour la variable Y (y_0 apparaît donc ici comme un *paramètre* — non spécifié — et non comme une variable). Cette fois-ci, c'est donc le théorème (EQ) qu'il faut utiliser : on trouve que la loi de X sachant que $\{Y = y_0\}$ a pour densité

$$x \mapsto Z^{-1}\mathbf{1}_{x \in [y_0, 1]}x^{-1} \quad (\text{FC})$$

(voir figure 4.1), où

$$Z = \int_{x=0}^1 \mathbf{1}_{x \in [y_0, 1]}x^{-1}dx \quad (\text{FD})$$

(cette dernière valeur ne dépend pas de x — la variable ' x ' dans l'expression de Z étant liée). Notez qu'en fait, l'expression de Z comme intégrale est automatique, puisque l'intégrale d'une densité de probabilité doit valoir 1. Par ailleurs, la valeur précise de Z (qui, en l'occurrence, vaut $\ln(1/y_0)$) ne nous importe (en général) pas vraiment (et d'ailleurs, dans bien des cas elle ne serait pas calculable) : en effet, décrire une densité de probabilité à *constante multiplicative près* caractérise complètement cette densité, vu la contrainte de normalisation à 1 de la masse totale.

Remarquez au passage l'aspect contre-intuitif du résultat obtenu : alors qu'on aurait pu se dire que le fait de savoir $Y = 0,34$ nous permettait seulement d'en conclure de $x \geq 0,34$, on s'aperçoit qu'en fait, *du fait de l'information obtenue*, les valeurs de X les plus petites ont plus de chance d'avoir eu lieu que les plus grandes, par exemple, alors que *a priori* X avait autant de chances d'être comprise entre 0,34 et 0,67 qu'entre 0,67 et 1, *à posteriori* (c'est-à-dire, *sachant* l'information sur Y) le premier cas a environ 1,77 fois plus de chances de se produire que le second! ♣

4.2 Indépendance

Lois produits

! **Théorème (FE)** (Loi produit). *Si les variables aléatoires X et Y sont indépendantes, alors la loi jointe de X et Y ne dépend que des lois de X et Y : on l'appelle la loi produit de $\text{Loi}(X)$ et $\text{Loi}(Y)$, notée $\text{Loi}(X) \otimes \text{Loi}(Y)$. De manière générale, la loi produit d'une distribution μ sur \mathcal{X} et d'une distribution ν sur \mathcal{Y} est caractérisée par :*

$$\forall A \subseteq \mathcal{X}, B \subseteq \mathcal{Y} \quad (\mu \otimes \nu)(A \times B) = \mu(A)\nu(B). \quad (\text{FF})$$

◇

Remarque (FG). Ce théorème n'est en fait pas valide lorsqu'on travaille sous la théorie des probabilités finiment additives : pour le coup, on a vraiment besoin de la théorie de Kolmogorov pour parler de lois produits. ♣

Proposition (FH) (Produit de lois discrètes). *Si les variables aléatoires X et Y ne prennent leurs valeurs que dans des ensembles discrets resp. \mathcal{X} et \mathcal{Y} et sont indépendantes, alors la variable jointe (X, Y) suit une loi discrète à valeurs dans $\mathcal{X} \times \mathcal{Y}$, avec pour tout $(x, y) \in \mathcal{X} \times \mathcal{Y}$:*

$$\mathbb{P}((X, Y) = (x, y)) = \mathbb{P}(X = x) \mathbb{P}(Y = y) \quad (\text{FI})$$

◇

Démonstration. Immédiat en observant que $\{(X, Y) = (x, y)\}$ peut se réécrire $\{X = x \text{ et } Y = y\}$. ♣

Proposition (FJ) (Produit de lois à densité). *Si les variables aléatoires X et Y sont à valeurs dans resp. \mathbb{R}^n et \mathbb{R}^m avec des lois de densités respectives $x \mapsto f(x)$ et $y \mapsto g(y)$, et que ces variables aléatoires sont indépendantes, alors la variable jointe (X, Y) a une loi de densité $(x, y) \mapsto f(x)g(y)$.* ◇

Démonstration. Cherchons la densité de la loi en (x, y) . On considère un voisinage infinitésimal de 0 dans \mathbb{R}^{n+m} de forme produit $dx \times dy$: alors

$$\begin{aligned} \mathbb{P}((X, Y) \in (x, y) + dx \times dy) &= \mathbb{P}(X \in x + dx \text{ et } Y \in y + dy) \\ &\stackrel{\text{ind}^{\text{ce}}}{=} \mathbb{P}(X \in x + dx) \mathbb{P}(Y \in y + dy) = f(x) \text{vol}(dx)g(y) \text{vol}(dy) \\ &= f(x)g(y) \text{vol}(dx \times dy), \end{aligned}$$

ce qui donne bien la densité recherchée. ♣

Remarque (FK). Bien entendu, on peut aussi considérer des produits de lois discrètes par des lois continues, ou même des choses encore plus compliquées... L'idée de base pour les calculs des deux propositions ci-dessus continuera à être valable, mais ne conduira alors pas à une "belle" expression de loi qui serait discrète ou à densité. ♣

Chapitre 5

Théorèmes-limites

Nous présentons à présent deux théorèmes fondamentaux de la théorie des probabilités qui permettent de comprendre le comportement à la limite de la somme de variables aléatoires indépendantes : comme nous le verrons dans la suite ce cours, l'importance de ces théorèmes en modélisation statistique est primordiale ! Dans ce polycopié, j'ai opté volontairement pour des énoncés assez informels, afin de ne pas nous encombrer avec les technicités qu'on appelle « uniforme intégrabilité », « convergence en probabilité » ou « convergence en loi »... Par ailleurs, notez que les énoncés qu'on appelle « loi des grands nombres » et « théorème-limite central » *stricto sensu* concernent en fait seulement le cas où les variables aléatoires sont identiquement distribuées : les résultats ci-dessous en sont donc en toute rigueur des généralisations, auxquelles les puristes donneraient des noms différents.

5.1 La loi des grands nombres

Commençons par rappeler l'énoncé de la loi des grands nombres telle qu'on l'introduit en général dans les cours de probabilités :

Théorème (FL) (Loi (faible) dans grands nombres, version stricte). *Soit P une distribution de probabilité L^1 sur \mathbb{R} et X_1, X_2, \dots une suite i.i.d. de variables de loi P . Alors*

$$\text{Loi}\left(n^{-1} \sum_{i=1}^n X_i\right) \xrightarrow{n \rightarrow \infty} \delta_{\mathbb{E}(P)}. \quad (\text{FM})$$

◇

Remarque (FN). Noter que la loi qui apparaît dans le membre de gauche de la relation de convergence est directement déterminée par P et par n : le fait d'introduire les X_i est juste un artifice pour exprimer cette loi plus facilement. ♣

Remarque (FO). La loi « forte » des grands nombres ajoute que la convergence de la moyenne des X_i vers la constante $\mathbb{E}(P)$ a également lieu presque-surement (et aussi dans L^1) ; mais cela ne comporte aucun intérêt du point de vue de la modélisation, donc nous en ferons abstraction. ♣

Remarque (FP). Les deux défauts de la loi des grands nombres dans sa version « classique » sont, d'une part, qu'on a besoin que les X_i soient tous de même loi ; d'autre part, que rien n'est dit sur la vitesse de convergence vers la limite, et donc du contrôle qu'on peut obtenir pour une valeur fixée de n . Or ces deux points sont

cruciaux en modélisation ; d'où le besoin d'une généralisation, que nous donnons ci-dessous. Cette généralisation sera donnée d'abord en termes informels (la seule à retenir, car de toutes façons, on n'a jamais de contrôle formel sur le qu'on fait en matière de modélisation), puis une version plus formelle sera présentée. ♣

! **Théorème (FQ)** (Loi des grands nombres, version généralisée informelle). *Soit une variable aléatoire $X = X_1 + \dots + X_n$ s'écrivant comme une somme de variables aléatoires, où :*

- (i) Les X_i sont indépendantes ;
- (ii) X_i a des fluctuations de l'ordre de ℓ_i ;
- (iii) Les fluctuations des X_i ne sont pas trop "sauvages" (à l'échelle ℓ_i) ;
- (iv) Le nombre de X_i nécessaires pour contribuer substantiellement à la somme des ℓ_i est très grand,

alors, à l'échelle $\sum_i \ell_i$, X a une très grande probabilité d'être presque égale à son espérance ; autrement dit on a une probabilité presque égale à 1 que $|X - \mathbb{E}(X)|$ soit beaucoup plus petit que $\sum_i \ell_i$. ◇

Remarque (FR). Notez que, par linéarité de l'espérance, l'espérance de X n'est autre que la somme des espérances des X_i . ♣

Remarque (FS). Par « les X_i ont des fluctuations de l'ordre de ℓ_i », on entend que "l'essentiel" des valeurs possibles que X_i peut prendre sont comprises dans une plage dont la longueur est "de l'ordre de" ℓ_i . Une façon possible de rendre concrète cette définition informelle pourrait être de dire que la *distance interquartiles* de $\text{Loi}(X_i)$ (à savoir, la différence entre les troisième et premier quartiles de cette loi) doit être inférieure à ℓ_i . Cependant, en fonction du contexte, il peut y avoir d'autres façons pertinentes d'exprimer l'idée de « fluctuations de l'ordre de ℓ », de sorte qu'une définition floue est en fait préférable ici.

Sur la notion que les fluctuations ne soient pas trop "sauvages" : intuitivement, l'idée est de dire que, pour chaque i , la probabilité que X_i prenne des valeurs nettement plus éloignées que ℓ_i de sa valeur la plus typique (disons sa médiane, ou son espérance) est très petite. Plus précisément, la condition technique requise relève en fait que ce qu'on appelle l'*uniforme intégrabilité* : la façon décroît dont la probabilité que X_i prenne des valeurs très éloignées de ses valeurs typiques doit, d'une part, avoir toujours la même forme, et d'autre part, assurer que $|X_i - \text{médiane}(\text{Loi}(X_i))|$ soit intégrable (avec, dès lors, une intégrale qui sera bornée par un multiple de ℓ_i). Le théorème ?? vu plus loin donne une façon précise de formuler cela.

Si la définition générale de « fluctuations pas trop sauvages à l'échelle ℓ_i » est compliquée, il se trouve cependant qu'heureusement, on peut donner une condition *suffisante* pour cette notion dans le cadre du théorème (FQ), consistant à simplement contrôler l'*écart-type* de X_i par ℓ_i :

$$\forall i \in \{1, \dots, n\} \quad \text{Var}^{1/2}(X_i) \leq \ell_i. \quad (\text{FT})$$

♣

! *Point (FU).* Concernant la remarque ci-dessus, dans le cadre de ce cours, pour les exemples que vous aurez à traiter, soit cette dernière forme de la condition de contrôle uniforme des fluctuations sera vérifiée, soit on sera dans un cas où les fluctuations des X_i n'auront clairement pas de contrôle "civilisé" à l'échelle

désirée. Mais dans tous les cas, il ne vous sera jamais demandé d'expliciter ce qu'on entend par « fluctuations pas trop sauvages » : je vous demande juste de rappeler qu'il y a une condition de « contrôle des fluctuations à l'échelle ℓ_i » à vérifier, et éventuellement de justifier informellement que cela semble vrai, du point de vue de l'idée générale, pour la situation à laquelle vous êtes confrontés. ♣

Remarque (FV). Toujours au sujet de la remarque (FS), notez qu'il y a des distributions de probabilité d'écart-type infini, dont les fluctuations seraient donc d'échelle infinie au sens proposé dans cette remarque. Et de fait, il existe des lois pour lesquelles on n'observe pas ce phénomène de « concentration » vers une loi presque constante à l'échelle n : par exemple, si les X_i suivent toutes la loi Cauchy(0, 1), déterminée par la densité

$$\mathbb{P}(\text{Cauchy}(0, 1) \in dx) := \frac{1}{\pi} \frac{1}{1+x^2} \text{vol}_1(dx), \quad (\text{FW})$$

alors on pourrait montrer que $X_1 + \dots + X_n$ suit la loi $n \times \text{Cauchy}(0, 1)$. Ainsi, alors que les X_i semblent avoir des fluctuations de l'ordre de grandeur de 1 (et toutes de la même forme, évidemment), leur somme n'est pas du tout constante à l'ordre de grandeur de n ... Cependant, lorsque nous mesurons l'échelle des fluctuations au sens de la remarque (FS), on s'aperçoit que cela ne cause pas de contradiction avec la loi des grands nombres : en effet, l'écart-type d'une loi de Cauchy est infini^[*], de sorte qu'il n'existe aucune échelle qui contrôle les fluctuations des X_i ! En revanche, si nous avons mesuré les fluctuations au sens de, par exemple, la distance interquartiles, on aurait trouvé pour les X_i des fluctuations d'ordre 2, tandis que X aurait des fluctuations d'ordre 2000, ce qui n'est évidemment pas négligeable devant 1000×2 ^[†]. Cela montre donc, d'une part, qu'on ne peut pas considérer n'importe quel concept de « fluctuations » lorsqu'on souhaite énoncer rigoureusement le théorème (FQ) ; et d'autre part, que le concept de fluctuations retenu devra nécessairement donner à certaines distributions de probabilité un « échelle de fluctuations » infinie. ♣

Remarque (FX). L'autre point un peu compliqué dans la définition ci-dessus est que « le nombre de X_i nécessaires pour contribuer substantiellement à la somme des ℓ_i est très grand ». Qu'entend-on par là ? En fait, si tous les ℓ_i étaient égaux, on pourrait simplement écrire « n est très grand ». Mais comme l'énoncé autorise les ℓ_i à être inégaux, il faut éviter un cas dégénéré comme celui où, par exemple, X_1 suivrait la loi $\text{Unif}^{\text{me}}(0, 1)$, n serait égal à 10^9 , et chacun des X_i pour $i > 1$ suivrait la loi $\text{Unif}^{\text{me}}(0, 10^{-19})$: dans ce cas, la somme des X_i serait *en pratique* quasiment égale à X_1 , et c'est donc *comme si* on avait une unique variable aléatoire... ! En fait, ce qui est demandé ici, c'est que pour tout ensemble I dont le cardinal ne serait pas « très grand », la somme des X_i pour $i \in I$ ne représenterait qu'une fraction négligeable de la somme totale des X_i : et plus précisément (dans la mesure où la notion de « négligeable » n'est pas claire pour des sommes de variables aléatoires) on veut que la somme des ℓ_i pour $i \in I$ soit négligeable devant la somme totale des ℓ_i . En fait, une manière simple, quoique pas forcément aussi claire à première vue, d'exprimer cette idée est de demander que, pour tout i , le ratio entre ℓ_i et la

[*]. En fait, même l'espérance d'une loi de Cauchy ne peut pas être définie, tellement cette loi a des queues lourdes ! On peut d'ailleurs montrer que dans le cas où les X_i sont i.i.d. et qu'on considère le comportement de leur somme à l'échelle n dans l'asymptotique $n \rightarrow \infty$, il est nécessaire et suffisant pour que la loi des grands nombres soit valide qu'on puisse définir une espérance aux X_i : c'est en fait l'énoncé classique de la loi des grands nombres tel qu'il est donné dans la plupart des cours de probabilité.

[†]. En fait, on peut même trouver des exemples similaires (avec les X_i i.i.d.) où les choses se passent d'une manière encore pire, au sens où les fluctuations de X (mesurées au sens de la distance interquartiles) croissent de façon *plus rapide* que n lorsque n augmente ! (aussi étonnant que cela paraisse à première vue). C'est par exemple le cas si les X_i suivent une loi Pareto(1, k) pour $k \in]0, 1[$.

somme de tous les ℓ_i' soit très grand. C'est cette formulation qui est utilisée dans le théorème ?? ci-après. \clubsuit

Remarque (FY). Il est clair que les conditions (i) et (iv) sont nécessaires :

- En ce qui concerne la condition (i), considérons X_1 suivant n'importe quelle distribution de probabilité dont les fluctuations sont d'échelle $\ell \in]0, \infty[$, et posons $X_2 = \dots = X_n = X_1$. (Il n'y a donc absolument pas indépendance!). Alors X sera égale à nX_1 , et aura donc des fluctuations d'échelle $n\ell$; alors que la loi des grands nombres requerrait des fluctuations beaucoup plus petites que $n\ell$.
- En ce qui concerne la condition (iv), prenons $n = 1$ (l'indépendance sera alors automatiquement satisfaite), avec X_1 suivant n'importe quelle distribution de probabilité dont les fluctuations sont d'échelle $\ell \in]0, \infty[$: dans ce cas, $X = X_1$, de sorte les fluctuations de X seront d'échelle ℓ , ce qui n'est évidemment pas beaucoup plus petit que $n \times \ell = 1 \times \ell = \ell$.

\clubsuit

Exemple (FZ). Soient X_1, \dots, X_{999} 999 variables indépendantes avec $X_1 \sim \text{Bernoulli}(0,001)$, $X_2 \sim \text{Bernoulli}(0,002)$, \dots , $X_{999} \sim \text{Bernoulli}(0,999)$; et soit $X := \sum_{i=1}^{999} X_i$. Les X_i sont indépendantes par hypothèse, et leur nombre peut manifestement être qualifié de « très grand ». Par ailleurs, puisque chaque X_i est à valeurs dans $\{0, 1\}$, l'interprétation intuitive de la notion de « contrôle des fluctuations » montre que l'ordre de grandeur des fluctuations des X_i est uniformément contrôlé par 1^[‡]. Par conséquent, à l'échelle $999 \times 1 = 999$, X devrait être « presque égale à son espérance »

$$\mathbb{E}(X) = \sum_{i=1}^{999} \mathbb{E}(X_i) = \sum_{i=1}^{999} \frac{i}{1000} = 499\frac{1}{2}; \quad (\text{GA})$$

autrement dit, l'archi-majorité des valeurs de X devraient être éloignées de la valeur $499\frac{1}{2}$ d'une distance beaucoup plus petite que 999. Vérifions cela empiriquement sous R :

```
> # Pour illustration : La commande ci-dessous génère, à chaque appel,
> # un nouveau vecteur aléatoire composé de 9 réalisations de variables
> # de Bernoulli indépendantes, de paramètres respectifs 0.1, ..., 0.9.
> # (En vrai, il faudra remplacer « 9 » par « 999 », et les paramètres par
> # 0.001, ..., 0.999).
> rbinom(9, 1, seq(from=.1, to=.9, by=.1))
[1] 0 0 0 0 1 1 0 0 1
> # Nous allons maintenant faire 100 tirages de ce type,
> # calculer à chaque fois la somme, et afficher les 100 résultats obtenus.
> resultats = numeric(100)
> for(i in 1:100)
+ resultats[i] = sum(rbinom(999, 1, seq(from=.001, to=.999, by=.001)))
> resultats
[1] 480 481 504 503 496 496 505 504 500 526 499 506 486 507 486 504 497 497
[19] 502 516 491 495 493 498 487 516 512 495 505 541 512 512 501 488 489 508
[37] 485 510 498 510 499 492 514 501 501 476 503 483 510 516 477 515 497 481
[55] 501 507 500 525 476 505 507 506 504 490 517 499 493 516 502 517 500 516
[73] 498 520 483 498 503 499 492 513 481 506 479 513 507 487 515 524 491 481
[91] 506 528 504 500 506 490 495 530 492 501
> # Trions maintenant les résultats pour voir comment se répartissent
> # les valeurs obtenues :
> sort(resultats)
[1] 476 476 477 479 480 481 481 481 481 483 483 485 486 486 487 487 488 489
```

[‡]. Avec la formulation rigoureuse proposée, on pourrait observer que pour tout i , on a $\text{Var}^{1/2}(X_i) \leq 1/2$.

[19] 490 490 491 491 492 492 492 493 493 495 495 495 496 496 497 497 497 498
 [37] 498 498 498 499 499 499 499 500 500 500 500 501 501 501 501 501 502 502
 [55] 503 503 503 504 504 504 504 504 505 505 505 506 506 506 506 506 507 507
 [73] 507 507 508 510 510 510 512 512 512 513 513 514 515 515 516 516 516 516
 [91] 516 517 517 520 524 525 526 528 530 541
 > # En 100 tirages, nous n'avons trouvé aucune valeur qui fût
 > # plus éloignée de 499.5 que ne l'est 541, soit une distance maximale
 > # de 41.5 ; ce qui est effectivement beaucoup plus petit que 999 !

♣

Voyons, pour information, un énoncé précis qui quantifie le théorème informel
Théorème (GB) (Loi des grands nombres, version généralisée formelle). *Soit $k > 1$. Alors il existe une fonction (explicitable) $\varphi_k(\bullet, \bullet) : \mathbb{R}_+^* \times]1, \infty[\rightarrow \mathbb{R}$, telle que pour tout $\varepsilon > 0$, on ait $\varphi_k(\varepsilon, \bar{n}) \xrightarrow{\bar{n} \rightarrow \infty} 0$, pour laquelle on a la propriété suivante :*

Pour $n \in \mathbb{N}^$, X_1, \dots, X_n une famille de variable aléatoires indépendantes toutes intégrables^[§], et $(\ell_i)_{1 \leq i \leq n}$ des nombres réels strictement positifs tels qu'on ait*

$$\forall i \in]1, n] \quad \forall z \in]1, \infty[\quad \mathbb{P}(|X_i - \mathbb{E}(X_i)| \geq z\ell_i) \leq \frac{1}{z^k}, \quad (\text{GC})$$

alors pour tout $\varepsilon > 0$,

$$\mathbb{P}\left(\left|\sum_i X_i - \sum_i \mathbb{E}(X_i)\right| \geq \varepsilon \sum_i \ell_i\right) \leq \varphi_k\left(\varepsilon, \frac{\sum_i \ell_i}{\max_i \ell_i}\right). \quad (\text{GD})$$

◇

Remarque (GE). En fait, il n'est pas nécessaire que n soit fini pour appliquer le théorème : tout ce qui compte est que la somme des ℓ_i le soit (et que la somme des $\mathbb{E}(X_i)$ soit absolument convergente), comme on le vérifie aisément par continuité. De même, on pourrait éventuellement autoriser certains des ℓ_i à valoir 0, ce qui reviendrait simplement à ajouter des v.a. constantes. ♣

Remarque (GF). Dans la condition sur le contrôle des fluctuations de X_i :

$$\forall z \in]1, \infty[\quad \mathbb{P}(|X_i - \mathbb{E}(X_i)| \geq z\ell_i) \leq \frac{1}{z^k}, \quad (\text{GG})$$

il n'est en fait pas très important que ce soit l'espérance de X_i qui intervienne. En effet, s'il existe un m_i tel qu'on ait le même contrôle sur $\mathbb{P}(|X_i - m_i| \geq z\ell_i)$, la condition d'origine est vérifiée, en remplaçant simplement ℓ_i par $(2k-1)/(k-1) \times \ell_i$. ♣

Remarque (GH). Pour $k = 3$, il est facile d'expliciter une fonction φ_k convenable : dans ce cas en effet, on a $\text{Var}(X_i) \leq 3\ell_i^2$, et donc on peut contrôler la variance de $\sum X_i$, puis appliquer l'inégalité de Bienaymé-Tchebychev pour en déduire, via l'inégalité de Hölder, le résultat, avec $\varphi_3(\varepsilon, \bar{n}) = 1/3\varepsilon^2\bar{n}$. (Un argument similaire fonctionne en fait dès lors que $k > 2$; par contre, pour $k \leq 2$ le théorème est plus difficile).

Le cas où on contrôle directement les écarts-types des X_i , quant à lui, implique les hypothèses du théorème pour $k = 2$. ♣

Remarque (GI). Le théorème tel qu'il est écrit ne redonne pas la forme classique de la loi des grands nombres. Pour cela, il faut aller encore plus loin et remplacer la condition de contrôle des queues en puissance de k par une fonction de la forme

$$\forall z \in]1, \infty[\quad \mathbb{P}(|X_i - \mathbb{E}(X_i)| \geq z\ell_i) \leq \kappa(z), \quad (\text{GJ})$$

où $\kappa(z)$ doit être une fonction (la même pour tous les i) vérifiant $\int_{z=1}^{\infty} \kappa(z) dz < \infty$. (D'où le fait, quand on se restreint au cas particulier $\kappa(z) := 1/z^k$, qu'on doit avoir $k > 1$). Dans ce cas la fonction φ_k doit être bien entendu remplacée par une certaine fonction $\varphi_{\kappa(\bullet)}$. Avec cette variante, on obtient bien un résultat qui généralise complètement le théorème ?? ♣

[§]. Cette condition est en fait automatique à cause du contrôle sur les fluctuations des X_i .

5.2 Le théorème-limite central

J'ai fait ici le choix de traiter le théorème-limite central d'une façon aussi similaire possible à la loi des grands nombres : comme vous allez le constater, les similarités sont en effet frappantes ! Commençons par l'énoncé "strict" ce théorème :

Théorème (GK) (Théorème-limite centre, version stricte). *Soit P une distribution de probabilité L^2 sur \mathbb{R} et X_1, X_2, \dots une suite i.i.d. de variables de loi P . Alors*

$$\text{Loi}\left(\frac{\sum_{i=1}^n X_i - n \mathbb{E}(P)}{n^{1/2}}\right) \xrightarrow{n \rightarrow \infty} \text{Normale}(0, \text{Var}(P)). \quad (\text{GL})$$

◇

Remarque (GM). On peut également écrire le théorème-limite central dans un cadre multidimensionnel, les X_i étant alors dans \mathbb{R}^d , et les $\text{Var}(X_i)$ devenant alors des matrices de covariances (dans $\mathbb{R}^{d \times d}$) : tout fonctionne exactement de la même façon. ♣

Remarque (GN). En fait, on peut même voir le TLC multidimensionnel comme un corolaire du TLC unidimensionnel, à cause du résultat suivant : une famille de distributions de probabilité $(P_n)_{n \in \mathbb{N}}$ sur \mathbb{R}^d converge en loi vers P_∞ si et seulement si, pour toute projection linéaire $\pi: \mathbb{R}^d \rightarrow \mathbb{R}$, $\pi_* P_n$ converge en loi vers $\pi_* P_\infty$. ♣

Le théorème-limite central, dans sa version stricte, souffre exactement des deux mêmes défauts que la loi des grands nombres : il requiert que les X_i aient toutes la même loi, et ne permet pas de quantifier la vitesse de convergence. Mais, de même que pour la loi des grands nombres, on peut en donner une version généralisée, extrêmement similaire à la version généralisée de la loi des grands nombres, tant dans son énoncé formel que dans son énoncé informel :

! **Théorème (GO)** (Théorème-limite central, version généralisée informelle). *Soit une variable aléatoire $X = X_1 + \dots + X_n$ s'écrivant comme une somme de variables aléatoires, où :*

- (i) Les X_i sont indépendantes ;
- (ii) X_i a des fluctuations de l'ordre de ℓ_i ;
- (iii) Les fluctuations des X_i ne sont pas trop "sauvages" (à l'échelle ℓ_i) ;
- (iv) Le nombre de X_i nécessaires pour que la somme des ℓ_i^2 contribue substantiellement à la somme des $\text{Var}(X_i)$ est très grand,

alors la loi de X est presque égale à la distribution normale ayant les mêmes espérance et écart-type que X ; autrement dit :

$$\text{Loi}(X) \approx \text{Normale}(\mathbb{E}(X), \text{Var}(X)). \quad (\text{GP})$$

◇

Remarque (GQ). Comme pour la loi des grands nombres, on a par linéarité de l'espérance que $\mathbb{E}(X) = \sum_i \mathbb{E}(X_i)$. En outre, puisque les X_i sont indépendantes, la variance de leur somme est égale à la somme de leurs variances :

$$\text{Var}(X) = \sum_{i=1}^n \text{Var}(X_i). \quad (\text{GR})$$

♣

Remarque (GS). On notera que la variance de X_i est elle aussi une façon d'exprimer les fluctuations de $\text{Loi}(X_i)$, de sorte qu'en fait, les ℓ_i ne pourront pas être n'importe quoi par rapport aux $\text{Var}^{1/2}(X_i)$. Plus précisément, dans le contexte du théorème-limite central, on peut montrer que les seules façons convenables d'interpréter la notion de « fluctuations pas trop sauvages » imposent qu'on ait $\ell_i \geq \text{Var}^{1/2}(X_i)$. ♣

Remarque (GT). Notez que, cette fois-ci, le condition qui impose qu'il y ait réellement un grand nombre de variables qui interviennent dans la somme est un peu plus compliquée, puisqu'elle compare les ℓ_i^2 à la somme totale des $\text{Var}(X_i)$. En fait, on pourrait écrire le théorème en imposant que les ℓ_i soient pris égaux aux $\text{Var}^{1/2}(X_i)$, ce qui rendrait alors la condition d'éparpillement des variables plus simple (et plus similaire à celle du théorème-limite central) ; mais cela empêcherait alors de considérer le cas où on a une poignée de variables dont le comportement serait très « sauvage » à l'échelle $\text{Var}^{1/2}(X_i)$, mais néanmoins « gentil » à une échelle ℓ_i supérieure dont la contribution serait faible : or ce genre de situations peut réellement se rencontrer en modélisation, comme nous le verrons d'ailleurs plus loin. ♣

Remarque (GU). La condition de « gentillesse » des fluctuations dans le théorème (GO) n'est, là encore, pas évidente à expliciter : ce qui se cache derrière est le concept d'*uniforme intégrabilité quadratique*... On peut néanmoins y écrire une condition suffisante, en disant par exemple que le moment quatrième (centré) de X_i doit être majoré par ℓ_i^4 :

$$\forall i \in \{1, \dots, n\} \quad \mathbb{E}((X_i - \mathbb{E}(X_i))^4) \leq \ell_i^4. \quad (\text{GV})$$

Il serait tentant de penser que ces conditions (iii) et (iv) ne sont que des subtilités techniques qui ne poseront jamais souci en pratique ; mais malheureusement, les contre-exemples (GZ) et (HA) ci-dessous montrent qu'il existe des situations tout à fait concrètes où le théorème-limite central échoue à cause d'une de ces deux conditions ! ♣

Point (GW). Concernant la remarque précédente, pas de panique : de façon similaire au paragraphe « dans le cadre de ce cours » qui concluait la remarque (FS), dans le cadre de ce cours, vous n'êtes pas censés savoir écrire à quoi correspondent précisément les conditions (iii) et (iv) (et encore moins les vérifier) : tout ce que je vous demande, c'est de rappeler qu'il faut avoir ces deux conditions (en en donnant les énoncés heuristiques tels qu'ils apparaissent dans le théorème (GO)), et de vous convaincre à gros traits qu'elles sont « moralement » valides dans la situation qui vous intéresse. (Si jamais je devais vous présenter une situation où une de ces conditions n'était pas vérifiée, ce serait toujours de façon flagrante). ♣

Exemple (GX). Nous reprenons ici la situation de l'exemple (FZ) : $X := X_1 + \dots + X_{999}$, les X_i étant indépendantes et de lois respectives Bernoulli($i/1000$). Nous avons déjà dit dans l'exemple (FZ) que les conditions (i) et (iv), qui sont identiques dans les théorèmes (FQ) et (GO), sont vérifiées. Nous voulons maintenant vérifier la condition (??), consistant à dire qu'il n'y a pas de X_i dont l'écart-type soit nettement plus grand que pour la plupart des autres X_i . Ici l'écart-type de X_i est $\sqrt{i(1000-i)}/1000$, qui est maximal pour X_{500} où il vaut alors 0,5, et dont les valeurs sont $\geq 0,4$ pour tout $i \in \{200, \dots, 800\}$, ce qui représente bien une large partie des X_i : l'hypothèse (iv) est donc bien vérifiée, et les écarts-type sont pour la plupart de l'ordre de 0,5. Maintenant, l'hypothèse (iii) consiste alors à dire que la probabilité qu'un des X_i s'éloigne de son espérance de nettement plus de 0,5 devient rapidement très petite. Mais ici tous les X_i sont à valeurs dans $\{0, 1\}$, donc quel que soit i , il est strictement impossible d'avoir $|X_i - \mathbb{E}(X_i)| > 1 = 2 \times 0,5$: l'hypothèse

de contrôle des fluctuations est donc vérifiée, avec $\ell_i = 1$. [¶]. On s'attend donc à ce que la loi de X soit très proche d'une loi normale ayant des paramètres compatibles avec l'espérance et l'écart-type de X .

L'espérance de X , nous l'avons calculé dans l'exemple (FZ), vaut 499,5. En ce qui concerne la calcul de l'écart-type, on utilise que la variance d'une somme indépendante [§] de variables aléatoires est la somme des variances; et puisqu'un calcul simple montre que $\text{Var}(X_i) = \frac{i}{1000} \left(1 - \frac{i}{1000}\right)$, on en déduit que

$$\text{Var}(X) = \sum_{i=1}^{999} \frac{i}{1000} \left(1 - \frac{i}{1000}\right) \stackrel{[**]}{=} \frac{999 \times 1001}{6 \times 1000} = 166,6665, \quad (\text{GY})$$

d'où $\text{Var}^{1/2}(X) = \sqrt{166,6665} \approx 12,910$: on doit donc avoir que la loi de X est très proche de la loi Normale(499,5; 12,910). Pour vérifier cela, nous voulons comparer les fonctions de répartition de la loi de X et de la loi Normale(499,5; 12,910), c.-à-d. les fonctions $x \mapsto \mathbb{P}(X \leq x)$ et $x \mapsto \mathbb{P}(\text{Normale}(499,5; 12,910) \leq x)$.

En ce qui concerne la deuxième fonction, nous pouvons la calculer grâce à la fonction `pnorm` de R . Pour la première, malheureusement il n'y a pas d'expression analytique de la fonction de répartition de X ; nous sommes donc obligés de l'approcher numériquement par la méthode de Monte-Carlo : c.-à-d. que nous simulons un grand nombre de fois (disons 2^{16}) la loi de X de façon indépendante, et que nous considérons que la mesure empirique des valeurs ainsi obtenues est très proche de la véritable loi de X . La code R est donné ci-dessus.

Le résultat graphique est donné sur la figure 5.1 : effectivement, la loi de X est très bien approchée par la loi normale que nous avons prédite! ☺

```
> # On calcule les moyenne et écart-type de X et de la loi normale
> # censée l'approximer, appelés m et s.
> m = 499.5
> (s = sqrt (166.6665))
[1] 12.90994
> # N est le nombre de simulations pour Monte-Carlo
> (N = 2^16)
[1] 65536
> # On simule N fois la loi de X (indépendamment à chaque fois bien sûr),
> # et on stocke les résultats dans "realisations.de.X". Notez que
> # la ligne de commande pour simuler la loi de X
> # est reprise du code précédent.
> realisations.de.X = numeric (N);
> for (i in 1:N)
+ realisations.de.X[i] = sum (rbinom (999, 1,
seq (from=.001, to=.999, by=.001)))
> # À l'aide de la commande "str", on peut regarder à quoi ressemble
```

[¶]. Par contre, si on avait pris "naïvement" $\ell_i \propto \text{Var}^{1/2}(X_i)$, cette fois-ci cela n'aurait pas fonctionné : en effet, si on prend par exemple $\ell_i = 8 \text{Var}^{1/2}(X_i)$, on a que $\mathbb{E}(|X_i - \mathbb{E}(X_i)|^4) \leq (8 \text{Var}^{1/2}(X_i))^4$ pour la grande majorité des i ... mais pas dans les 30 cas (sur 999) où $i \leq 15$ ou $i \geq 985$. Et c'est là qu'on comprend l'intérêt de rendre ℓ_i potentiellement différent de $\text{Var}^{1/2}(X_i)$: car certes, les variables aléatoires en question ont des fluctuations "sauvages" à l'échelle de leur écart-type; mais elles ont malgré tout des fluctuations bien contrôlées à l'échelle 1 : or, pour 30 variables seulement, remplacer les écarts-types par 1 n'a pas grand impact sur la somme des ℓ_i^2 ! Ainsi la condition d'éparpillement est bien vérifiée.

[§]. Ne pas oublier l'hypothèse d'indépendance ici !

[**]. On utilise ici l'identité $\sum_{i=1}^{N-1} i(N-i) = (N-1)N(N+1)/6$, que je ne démontre pas (même si c'est assez simple à vérifier), dans la mesure où cela n'est pas nécessaire à la compréhension de l'exemple.

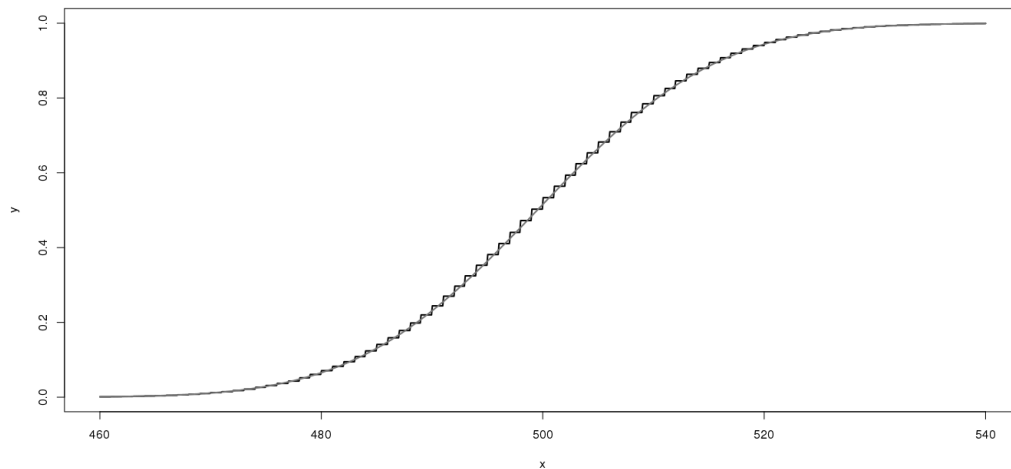


FIGURE 5.1 – La loi de X , dont la fonction de répartition est tracée en noir, est pratiquement égale à la loi normale de mêmes espérance et écart-type, tracée en gris.

```

> # "realisations.de.X".
> str (realisations.de.X)
 num [1:65536] 473 479 504 500 499 492 506 505 503 526 ...
> # Grâce aux conversions implicites entre booléens et numériques,
> # la commande « mean (simulations.de.X <= x) » nous donne la proportion
> # des simulations de X qui ont été inférieures ou égales à x, c.-
> # à-d.
> # (l'estimation par Monte-Carlo de) la probabilité que X soit
> # inférieure ou égale à x, alias « la valeur en x de
> # la fonction de répartition de X ».
> mean (realisations.de.X <= 470)
[1] 0.0123291
> mean (realisations.de.X <= 490)
[1] 0.2442474
> mean (realisations.de.X <= 510)
[1] 0.8062286
> # Maintenant on va tracer la fonction de répartition de X,
> # pour x variant entre 460 et 540. Pour ce faire, on appelle "x"
> # un vecteur contenant toutes les valeurs en abscisse qu'on échantillonne
> # entre 460 et 540, et "y" sera le vecteur des valeurs correspondantes
> # de la fonction de répartition.
> x = seq (from=460, to=540, length.out=2^10)
> y = numeric (length (x))
> for (i in 1 : length (x))
+ y[i] = mean (realisations.de.X <= x[i])
> # On trace la fonction de répartition de X (en noir).
> plot (x, y, type = "l", col = "black", lwd = 2)
> # On ajoute à ce tracé la courbe de la fonction de répartition
> # de notre loi normale (en gris).
> curve (pnorm (x, m, s), xname = "x", from = 460, to = 540, add = TRUE,
+ col = "gray", lwd = 2)

```

⊞

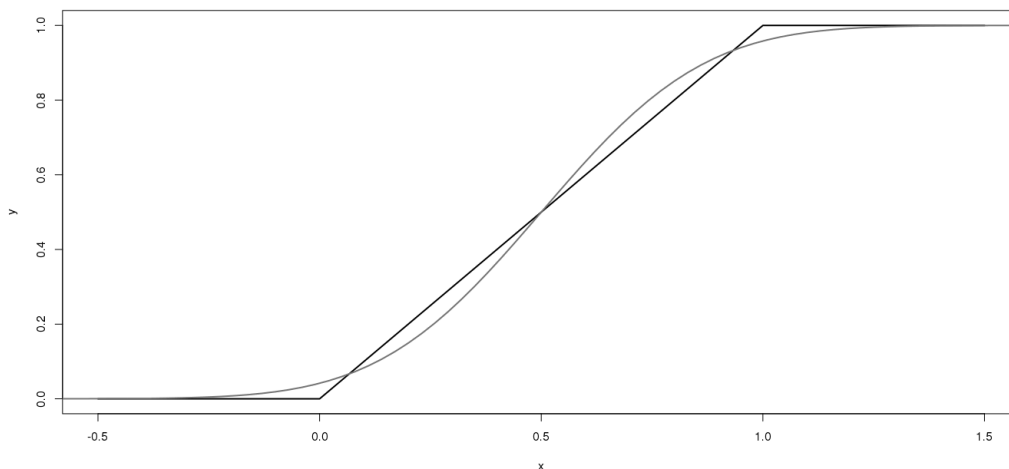


FIGURE 5.2 – Comparaison entre les fonctions de répartition d'une loi uniforme (en noir) et de la loi normale de mêmes espérance et variance (en gris) : la différence est sensible !

Exemple (GZ). Le contre-exemple suivant montre en quoi est-ce que la condition (??), selon laquelle aucune variable ne doit prédominer devant les autres, est nécessaire. Soient X_1, \dots, X_{1000} des variables indépendantes pour lesquelles X_i est distribuée uniformément sur $\{0, 10^{-i}, 2 \cdot 10^{-i}, \dots, 9 \cdot 10^{-i}\}$; et soit $X := \sum_i X_i$. Alors X est un nombre décimal de la forme $0, x_1 x_2 \dots x_{1000}$ dont tous les chiffres ont été tirés indépendamment uniformément entre 0 et 9; autrement dit, c'est un nombre aléatoire tiré uniformément dans l'ensemble $\{0, 10^{-1000}, 2 \cdot 10^{-1000}, \dots, 1 - 10^{-1000}\}$, ce qui est pratiquement la même chose que la loi $\text{Unif}^{\text{me}}(0, 1)$: X ne suit donc pas du tout une loi normale ! (voir fig. 5.2). Quand on regarde ce qui fait que le théorème (GO) ne s'applique pas, on s'aperçoit que les conditions (i), (iv) et (iii) sont bien satisfaites, mais que la condition (??) est violée, puisque l'écart-type de X_i est proportionnel à 10^{-i} , de sorte que la v.a. X_1 , notamment, a un écart-type beaucoup plus important que toutes les autres. ♣

Exemple (HA). Le contre-exemple suivant montre en quoi est-ce que la condition (iii), selon laquelle les fluctuations extrêmes des variables doivent être très rares, est nécessaire. Soient X_1, \dots, X_{1500} des v.a.i.i.d. suivant la loi Bernoulli(10^{-3}). Alors la somme $X := \sum_i X_i$ suit la loi Binom^{le}(1500, 10^{-3}), dont vous savez depuis les classes préparatoires qu'elle est très proche de la loi Poisson(1,5) (la valeur 1,5 provenant du calcul de 1500×10^{-3}). Ce qui n'est absolument pas une loi normale ! (voir fig. 5.3). Cette fois-ci les conditions (i), (iv) et (??) sont bien vérifiées, mais c'est la condition (iii) qui coince : en effet, la loi Bernoulli(10^{-3}), dont l'espérance est pratiquement nulle et dont la variance vaut environ 0,03, a une probabilité de 1 ‰ de fluctuer de plus de 30 écarts-types (cela se produit lorsqu'elle prend la valeur 1); ce qui est bien trop élevé pour une fluctuation aussi importante, car correspondant à un kurtosis très important d'environ 998 ! ♣

De même que pour la loi des grands nombres, nous pouvons formaliser notre énoncé du théorème-limite central en un résultat plus rigoureux :

Théorème (HB) (Théorème-limite central, version généralisée formelle). *Soit $k > 2$. Alors il existe une fonction (explicitable) $\psi_k(\bullet) : [1, \infty[\rightarrow \mathbb{R}$, telle que pour tout $\varepsilon > 0$, on ait $\varphi_k(\bar{n}) \xrightarrow{\bar{n} \rightarrow \infty} 0$, pour laquelle on a la propriété suivante :*

Pour $n \in \mathbb{N}^$, X_1, \dots, X_n une famille de variable aléatoires indépendantes toutes de carrés intégrables^[††], et $(\ell_i)_{1 \leq i \leq n}$ des nombres réels strictement positifs tels qu'on*

[††]. Cette condition est en fait automatique à cause du contrôle sur les fluctuations des X_i .

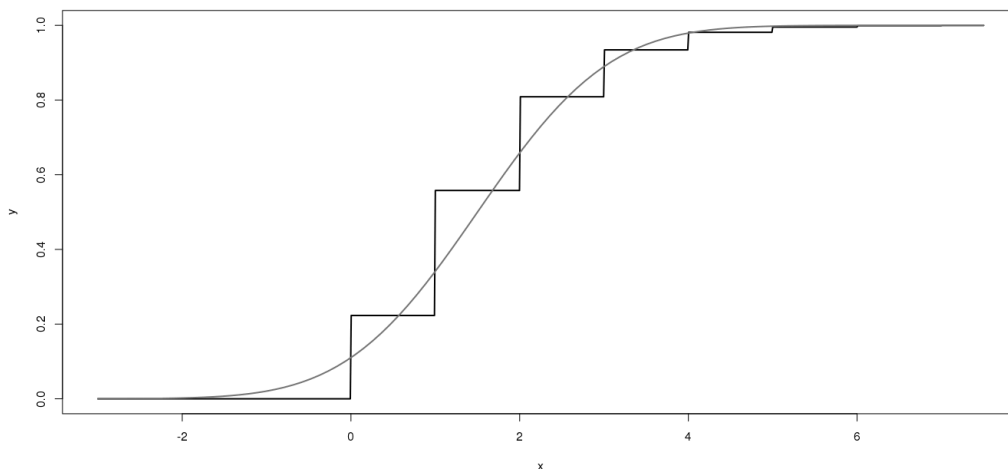


FIGURE 5.3 – Comparaison entre les fonctions de répartition d’une loi Poisson(1,5) (en noir) et de la loi normale de mêmes espérance et variance (en gris) : la différence est très sensible !

ait

$$\forall i \in \llbracket 1, n \rrbracket \quad \forall z \in]1, \infty[\quad \mathbb{P}(|X_i - \mathbb{E}(X_i)| \geq z\ell_i) \leq \frac{1}{z^k}, \quad (\text{HC})$$

$$d_{\text{KS}}(\text{Loi}(X), \text{Normale}(\mathbb{E}(X), \text{Var}(X))) \leq \psi_k \left(\frac{\sum_i \text{Var}(X_i)}{\max_i \ell_i^2} \right), \quad (\text{HD})$$

où $d_{\text{KS}}(P, Q)$ désigne la distance de Kolmogorov-Smirnov entre deux distributions de probabilité P et Q définies sur \mathbb{R} , définie comme la norme du supremum de la différence entre leurs fonctions de répartition :

$$d_{\text{KS}}(P, Q) := \sup_{x \in \mathbb{R}} |\mathbb{P}(Q \leq x) - \mathbb{P}(P \leq x)| : \quad (\text{HE})$$

le point ici étant en particulier que la distance de Kolmogorov-Smirnov “métrise”, dans un certain sens, la convergence sur les distributions de probabilité : en particulier, une suite de lois de probabilité $(P_n)_{n \in \mathbb{N}}$ converge vers une loi-limite P_∞ dès lors que $d_{\text{KS}}(P_n, P_\infty) \xrightarrow{n \rightarrow \infty} 0$, la réciproque étant également vraie dès lors que la mesure P_∞ est diffuse^[††]. \diamond

Remarque (HF). On notera que, contrairement à l’énoncé analogue pour la loi des grands nombres où on requerrait d’avoir $k > 1$, ici on a besoin d’un contrôle plus fort sur les fluctuations, puisqu’on a besoin d’avoir $k > 2$: cela reflète en fait l’hypothèse, dans les énoncés classiques, qu’on avait besoin d’avoir une intégrabilité L^2 plutôt que L^1 . \clubsuit

Remarque (HG). La condition de contrôle “gentil” des fluctuations que nous avons énoncée plus haut implique les hypothèses du théorème pour $k = 4$. \clubsuit

[††]. En outre, la distance de Kolmogorov-Smirnov est invariante par les bijections croissantes, et donc en particulier par l’opération consistant à recentrer et normaliser la variance de X : par conséquent le résultat énoncé implique bien la convergence vers la loi normale standard comme on l’écrit habituellement.

Chapitre 6

Distributions de probabilité remarquables

Préambule Certaines distributions de probabilités remarquables possèdent des noms : par exemple, les lois uniformes, les lois de Poisson, les lois gamma, ... Ce chapitre en présente quelques-unes, ainsi que leurs propriétés de base. Dans le cadre de ce cours, une partie des définitions, repérées dans la marge, doivent être connues par cœur. Les propriétés de ces distributions, elles, ne sont le plus souvent pas censées être apprises à la lettre : néanmoins, dans certains cas j'attends de vous que vous soyez capables de les retrouver à partir des définitions (par exemple, retrouver la densité d'une loi exponentielle à partir de sa caractérisation). Ainsi, dans le cadre de ce chapitre, lorsqu'une propriété d'une loi est marquée dans la page comme « à retenir », cela doit être compris comme « ... ou à savoir retrouver au besoin » ; en outre, le cas échéant, les démonstrations permettant de retrouver les propriétés énoncées seront imprimées en taille normale (contrairement aux preuves fournies à pur titre de complément, pour lesquelles j'utilise les petits caractères).

Une attention plus particulière sera apportée au cas des lois normales, pour lesquelles, nonobstant ce qui précède, un certain nombre de propriétés non triviales, mais essentielles, sont à connaître par cœur.

6.1 Familles de lois et paramétrages

Paramètres spécifiant une loi au sein d'une famille La quasi-totalité des lois que nous allons aborder ci-dessous viennent par *familles* : ainsi, même si l'on parle parfois par abus de langage de « la » loi normale (unidimensionnelle), il y a en réalité une famille de lois normales, une loi normale au sein de cette famille étant spécifiée par deux paramètres. Selon la famille de lois, il peut y avoir besoin de 0, 1, 2 ou 3 paramètres, voire plus, pour spécifier la loi.

Typologie des paramètres Pour les lois de probabilité sur \mathbb{R} , on a coutume de distinguer trois types de paramètres parmi ceux identifiant une distribution de probabilité au sein d'une famille donnée : les paramètres de *forme*, les paramètres d'*échelle* (ou, éventuellement, d'échelle à une certaine puissance) et les paramètres de *position* :

Définition (HH). Supposons une loi de la forme $\text{Famille}(\lambda, \mu_1, \mu_2, \dots)$, où λ et les μ_i sont les paramètres de la loi. Il se peut que certaines transformations affines laissent

la famille de lois stable, au sens suivant : pour T une telle transformation affine, pour tout choix de paramètres $(\lambda, \mu_1, \mu_2, \dots)$, la mesure-image de Famille $(\lambda, \mu_1, \mu_2, \dots)$ par la transformation T est à nouveau une distribution de la famille Famille (mais de paramètres différents bien sûr si $T \neq \text{id}$) : notons cette mesure-image Famille $(\lambda', \mu'_1, \mu'_2, \dots)$. Selon la façon dont λ et λ' sont liés par la transformation T , on utilisera des adjectifs spécifiques pour qualifier le paramètre λ :

- Lorsque $\lambda' = T(\lambda)$ (pour toute transformation affine T et tout choix de paramètres $(\lambda, \mu_1, \mu_2, \dots)$, s'entend), on dit que λ est un paramètre de *position* : informellement, λ repère une certaine position dans la distribution de probabilité $(\lambda, \mu_1, \mu_2, \dots)$, et “bouge” donc de conserve avec la distribution de probabilité lorsqu'on lui fait subir une transformation affine.
- Lorsque, pour une transformation affine s'écrivant $T(x) = ax + b$, la relation liant λ' à λ est que $\lambda' = a\lambda$, alors λ est qualifié de « paramètre d'échelle ». Moralement, λ dit donc à quel point la distribution est étalée.
- Variante du cas précédent : lorsque la relation liant λ' à λ est de la forme « $\lambda' = a^k \lambda$ » pour un certain $k \neq 0$, λ est qualifié de « paramètre d'échelle à la puissance k ». Ainsi, pour $k = -1$, λ sera proportionnel à l'inverse de l'«étalement» de la distribution ; pour $k = 2$, il sera proportionnel au carré de son étalement ($k = 2$), etc.
- Lorsque, quelle que soit la transformation affine T , on a $\lambda' = \lambda$, alors λ est qualifié de *paramètre de forme* : moralement, cela signifie que λ ne dépend que de la forme de la distribution et reste donc invariant lorsqu'on lui applique une transformation affine. ♡

Cette définition sera éclairée par les exemples que nous verrons dans la suite du chapitre.

Attention aux choix de conventions Il est très important de savoir que le paramétrage des lois n'est pas toujours standardisé ! Un des exemples les plus marquants est celui des lois normales : certains (dont ce cours) les paramètrent en précisant leur moyenne et leur *variance*, tandis que d'autres (notamment de nombreux logiciels, comme R) les paramètrent en indiquant leur moyenne et leur *écart-type*. Le cas des lois exponentielles donne également lieu à plusieurs paramétrages concurrents : certains (dont ce cours) les paramètrent par leur taux, d'autres par leur demi-vie, d'autres encore par leur espérance... Le bon ingénieur, toujours soucieux d'éviter les quiproquos lors de ses collaborations^[*], ne parlera donc pas de « la loi normale de paramètres $(1, 4)$ », mais plutôt de « la loi normale *d'espérance 1 et de variance 4* »...!^[†] Un souci similaire peut aussi être rencontré avec certains choix de conventions : par exemple, certains considèrent que les lois géométriques comptent le nombre *total* d'essais requis pour obtenir un succès (ce sont donc des lois portées par \mathbb{N}^*), tandis que d'autres (dont ce cours) considèrent qu'elles comptent les

[*]. Un des exemples les plus célèbres est celui de la sonde *Mars Climate Orbiter*, qui en 1999 s'est désintégrée dans l'atmosphère martienne en raison de l'interfaçage entre les travaux d'une équipe qui travaillait en kilomètres et ceux d'une autre équipe qui travaillait en miles... Une gaffe à 600 M\$! (valeur ajustée de l'inflation pour 2025).

[†]. Cela dit, au sein d'un même travail, il n'est bien entendu pas nécessaire de re-préciser les conventions employées pour une certaine famille de loi à chaque occurrence de ladite famille ! ☹ Ainsi, dans le cadre de ce cours, les conventions que j'utiliserai pour les lois classiques seront fixées une fois pour toutes au sein du présent chapitre ; et dès lors, je pourrai me contenter de parler de « la loi Expon^{le}(λ) » sans me sentir obligé d'ajouter à chaque fois « le paramètre utilisé pour spécifier une loi exponentielle étant le taux de cette dernière »...

échecs avant un succès... Donc, là aussi, il faut préciser la convention quand il y a lieu.

Notation des familles de lois Au lycée, on vous a introduit un certain nombre de notations très abrégées pour certaines distributions de probabilité, comme ‘ \mathcal{G} ’ pour les distributions géométriques par exemple. Je vous demande explicitement de ne pas utiliser ces abréviations dorénavant (avec une tolérance éventuelle pour la loi normale), car en statistique, le nombre de distributions “classiques” est tellement grand que cela prêterait à confusion : ‘ \mathcal{G} ’ désigne-t-il une loi géométrique ? une loi gamma ? une loi de Gumbel ?... Pour éviter les confusions, on notera donc explicitement « Géométrique(p) », ou éventuellement une abréviation non ambiguë comme « Géom^{que}(p) ».

6.2 Quelques lois remarquables discrètes

Remarque (HI). Les familles de lois portées par des structures discrètes comme \mathbb{N} ne peuvent évidemment pas être invariantes par un groupe continu de transformations affines (puisque en général, l’image affine d’un point de \mathbb{N} ne serait pas dans \mathbb{N} !) : dès lors, dans le cas discret il n’y a pas vraiment lieu de distinguer paramètres de position, d’échelle et de forme. ♣

Loi uniforme sur un ensemble fini

Définition (HJ) (Loi uniforme sur un ensemble fini). Pour \mathcal{X} un ensemble fini, la *loi uniforme* sur \mathcal{X} , notée $\text{Unif}^{\text{me}}(\mathcal{X})$, est la distribution de probabilité portée par \mathcal{X} qui donne une masse $1/|\mathcal{X}|$ à chaque point de \mathcal{X} . ♡

Remarque (HK). Nonobstant la notation utilisée, les lois uniformes sur les ensembles finis ne forment pas réellement une « famille » de lois, dans la mesure où l’ensemble \mathcal{X} ne peut pas se ramener à la description d’un nombre fini fixé de réels, et n’est donc pas vraiment un « paramètre ». ♣

Un cas particulièrement fréquent est le cas où $\mathcal{X} = \llbracket a, b \rrbracket$ pour $a, b \in \mathbb{Z}$ (et $a \leq b$). Dans ce cas, on peut calculer explicitement l’espérance et la variance de la loi. Dans la mesure où l’espérance de la loi découle directement de la symétrie de la loi $\text{Unif}^{\text{me}}(\llbracket a, b \rrbracket)$ par rapport au point $(a+b)/2$, j’attends de vous que vous soyez capable de la retrouver de vous-mêmes ; pour la variance en revanche, elle vous sera redonnée le cas échéant (ou éventuellement fera l’objet d’un exercice guidé).

Théorème (HL). Pour $a, b \in \mathbb{Z}$ avec $a \leq b$,

$$\mathbb{E}(\text{Unif}^{\text{me}}(\llbracket a, b \rrbracket)) = \frac{a+b}{2}. \quad (\text{HM})$$

◇

Théorème (HN). Sous les mêmes hypothèses que le théorème (HL),

$$\text{Var}(\text{Unif}^{\text{me}}(\llbracket a, b \rrbracket)) = \frac{1}{12}(b-a)^2 + \frac{1}{6}(b-a). \quad (\text{HO})$$

◇

Remarque (HP). Attention à ne pas confondre la loi $\text{Unif}^{\text{me}}(\llbracket a, b \rrbracket)$, qui est la loi *discrète* sur l’ensemble d’entiers $\{a, a+1, \dots, b-1, b\}$, avec la loi $\text{Unif}^{\text{me}}(a, b)$, qui est la loi *continue* sur l’intervalle réel $]a, b[$ (confer § 6.4). ♣

Un cas particulier de loi uniforme sur un ensemble fini est lorsque cet ensemble est $\{-1, +1\}$, autrement dit, lorsqu'on tire un *signe* aléatoire uniforme. La loi $\text{Unif}^{\text{me}}(\{-1, +1\})$ d'un tel signe uniforme est également appelée « loi de Rademacher ».

L'espérance et la variance d'une loi de Rademacher se trouvent immédiatement par calcul direct, et doivent donc pouvoir être recalculées au besoin :

Proposition (HQ).

$$\begin{aligned} \mathbb{E}(\text{Unif}^{\text{me}}(\{\pm 1\})) &= 0; & (\text{HR}) \\ \text{Var}(\text{Unif}^{\text{me}}(\{\pm 1\})) &= 1. & \diamond \end{aligned}$$

Lois de Bernoulli

! **Définition (HS)** (Loi de Bernoulli). Pour $p \in [0, 1]$, la *loi de Bernoulli* de paramètre p , notée $\text{Bernoulli}(p)$, est la loi portée par $\{0, 1\}$ qui attribue la masse p à 1 et la masse $(1 - p)$ à 0. \heartsuit

Dans la mesure où les espérance et variance d'une loi de Bernoulli s'obtiennent immédiatement par calcul direct ^[‡], vous êtes censés pouvoir les retrouver par vous-mêmes :

! **Proposition (HT).**

$$\begin{aligned} \mathbb{E}(\text{Bernoulli}(p)) &= p; & (\text{HU}) \\ \text{Var}(\text{Bernoulli}(p)) &= p(1 - p). & \diamond \end{aligned}$$

Remarque (HV). Il peut être néanmoins être pratique de retenir par cœur la formule pour la variance de la loi de Bernoulli. Un moyen mnémotechnique pour ce faire est d'observer que, puisque la loi de Bernoulli devient une masse de Dirac lorsque $p = 0$ ou $p = 1$, sa variance doit s'annuler dans ce cas, ce qui explique l'apparition du facteur « $p(1 - p)$ » ; et dès lors, il faut juste retenir qu'il n'y a rien d'autre à ajouter ! \smile \clubsuit

Lois binomiales

! **Définition (HW)** (Loi binomiale). Pour $n \in \mathbb{N}$ et $p \in [0, 1]$, la *loi binomiale* de paramètres n et p , notée $\text{Binom}^{\text{le}}(n, p)$, est la loi de la somme de n v.a.i.i.d. $\text{Bernoulli}(p)$: en d'autres termes, c'est la mesure-image par l'application « somme » (de \mathbb{N}^n dans \mathbb{N}) de $\text{Bernoulli}(p)^{\otimes n}$. Cette loi décrit donc le nombre de succès lorsqu'on procède à n expériences indépendantes ayant chacune une probabilité p de réussir. \heartsuit

De la définition de la loi binomiale découle immédiatement la proposition suivante, très utile :

! **Proposition (HX).** Pour $n, m \in \mathbb{N}$ et $p \in [0, 1]$, si X et Y sont deux variables indépendantes avec $X \sim \text{Binom}^{\text{le}}(n, p)$ et $Y \sim \text{Binom}^{\text{le}}(m, p)$, alors $(X + Y)$ suit la loi $\text{Binom}^{\text{le}}(n + m, p)$. [En termes plus savants, on formule parfois cela en disant qu'on a la formule de convolution ^[§] : $\text{Binom}^{\text{le}}(n, p) * \text{Binom}^{\text{le}}(m, p) = \text{Binom}^{\text{le}}(n + m, p)$]. \diamond

[‡]. Pour la variance, le plus simple est d'utiliser que $\text{Var}(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$, et que pour X suivant une loi de Bernoulli, $X^2 = X$: on trouve alors que la variance vaut $p - p^2$, soit $p(1 - p)$.

[§]. Pour P et Q deux lois de probabilité sur \mathbb{R} , le *produit de convolution* de P et Q , noté $P * Q$, est la mesure-image de la loi-produit $P \otimes Q$ (qui est une loi sur \mathbb{R}^2) par l'application « somme » (qui va de \mathbb{R}^2 dans \mathbb{R}). Ce produit de convolution décrit donc la loi de la somme $(X + Y)$ lorsque X et Y sont deux v.a. indépendantes de lois respectives P et Q .

Remarque (HY). Ne surtout pas oublier l'hypothèse d'indépendance, qui est essentielle! Ne pas oublier non plus que les seconds paramètres des deux lois binomiales doivent être égaux! ♣

La fonction de masse d'une loi de Bernoulli se retrouve immédiatement à partir de sa définition (nous rappelons le raisonnement ci-dessous), et on obtient alors la proposition suivante :

Proposition (HZ). *La loi Binom^{le}(n, p) est portée par $\llbracket 0, n \rrbracket$; et pour $k \in \llbracket 0, n \rrbracket$:*

$$\mathbb{P}(\text{Binom}^{\text{le}}(n, p) = k) = \binom{n}{k} p^k (1-p)^{n-k}. \quad (\text{IA})$$

◇

Démonstration. En effet, soient X_1, \dots, X_n i.i.d. Bernoulli(p) et S leur somme. La probabilité que S vaille k est la probabilité que k exactement des X_i valent 1 et que les $(n-k)$ autres valent zéro. Cet évènement se décompose en $\binom{n}{k}$ évènements disjoints, selon quels sont les indices i pour lesquels X_i vaut 1 (le nombre de façons de choisir cet ensemble d'indices correspond bien à la sélection de k éléments parmi n). Maintenant, quelle est la probabilité que les k variables X_i ayant les bons indices valent 1 et que les $(n-k)$ autres variables valent 0? Ces variables étant indépendantes, il faut faire le produit des probabilités pour chaque possibilité individuelle : on a ainsi k facteurs p (pour les X_i devant valoir 1) et $(n-k)$ facteurs $(1-p)$ (pour ceux devant valoir 0), ce qui conduit bien au résultat recherché (en observant que nos $\binom{n}{k}$ évènements disjoints se trouvent ici être équiprobables). ♣

Remarque (IB). En revanche, la formule du coefficient binomial $\binom{n}{p} = n! / p!(n-p)!$ n'est pas exigible dans ce module, et sera rappelée au besoin. (Mais la *définition* du coefficient binomial, à savoir le nombre de façons de choisir p éléments parmi n , doit être sue!). ♣

En utilisant la définition de la loi binomiale à partir de lois de Bernoulli (et en utilisant resp. la linéarité de l'espérance et la formule d'additivité des variances pour des v.a. indépendantes), on a la proposition suivante, dont vous devez savoir retrouver les résultats :

Proposition (IC). *Pour $n \in \mathbb{N}$, $p \in [0, 1]$:*

$$\mathbb{E}(\text{Binom}^{\text{le}}(n, p)) = np; \quad (\text{ID})$$

$$\text{Var}(\text{Binom}^{\text{le}}(n, p)) = np(1-p). \quad \diamond$$

Lois de Poisson

Définition (IE) (Loi de Poisson). Pour $\theta \in \mathbb{R}_+$, la *loi de Poisson* de paramètre θ , notée Poisson(θ), est la loi du nombre d'évènements qui se produisent lorsqu'on dispose d'un très grand nombre d'évènements possibles, que ces évènements sont indépendants, que chacun de ces évènements a une probabilité individuelle très petite de se produire, et que l'espérance du nombre d'évènements se produisant vaut θ . On a donc $\mathbb{E}(\text{Poisson}(\theta)) = \theta$ par *définition*. ♣

Remarque (IF). Autant la *formule* de la loi de Poisson, qui est compliquée et ne peut pas être retenue par un “truc” simple, n’est pas à connaître par cœur dans le cadre de ce cours, autant j’estime qu’il est vraiment important de connaître la définition informelle ci-dessus : en effet, la pertinence qu’il peut y avoir (ou pas) à utiliser une loi de Poisson dans une modélisation découlera du fait qu’on peut considérer qu’on se trouve (ou pas) dans la situation de la définition (IE) ! ♣

Remarque (IG). On définit parfois la loi Poisson(θ) comme la limite de la loi Binom^{le}($n, n^{-1}\theta$) lorsque n tend vers l’infini. Cette façon de voir les choses, qui a l’avantage d’être plus rigoureuse, revient en fait à considérer, dans la définition (IE) ci-dessus, le cas où les événements individuels sont tous équiprobables. Cependant, en réalité la loi de Poisson est bien plus générale que cela, puisqu’elle apparaît même en l’absence d’équiprobabilité : tout ce qui importe, c’est que tous les événements individuels aient chacun une probabilité minuscule ! ♣

De la définition de la loi de Poisson découle l’importante propriété suivante (analogue à la proposition (HX)). (Et, comme dans le cas de la proposition (HX)), on veillera à ne pas oublier l’hypothèse d’indépendance !) :

! **Proposition (IH).** Pour $\theta, \lambda \in \mathbb{R}_+$, si X et Y sont des variables aléatoires indépendantes de lois resp. Poisson(θ) et Poisson(λ), alors $X+Y$ suit la loi Poisson($\theta+\lambda$). Ou encore, en termes de convolution : Poisson(θ) * Poisson(λ) = Poisson($\theta+\lambda$). ◇

La fonction de masse d’une loi de Poisson peut être exprimée analytiquement, même si je ne vous demande pas de retenir ce résultat :

Théorème (II). La loi Poisson(θ) est portée par \mathbb{N} ; et pour tout $n \in \mathbb{N}$:

$$\mathbb{P}(\text{Poisson}(\theta) = n) = e^{-\theta} \frac{\theta^n}{n!}. \quad (\text{IJ})$$

Démonstration. Notons $(A_i)_{i \in \mathcal{J}}$ l’ensemble des événements considérés, $(p_i)_{i \in \mathcal{J}}$ leurs probabilités respectives, et posons $\varepsilon := \max_{i \in \mathcal{J}} p_i$. (Rappelons qu’avec nos hypothèses, on a $|\mathcal{J}| \gg 1$, $p_i \ll 1 \forall i \in \mathcal{J}$, et $\varepsilon \ll 1$). Notons $S := \sum_{i \in \mathcal{J}} \mathbf{1}_{A_i}$ la variable aléatoire qui compte combien de A_i sont réalisés, et fixons-nous $k \in \mathbb{N}$: notre objectif est alors de calculer (aux approximations près découlant des hypothèses) la probabilité $\mathbb{P}(S = k)$.

Dans la suite de cette démonstration, nous noterons $\mathcal{P}_k(\mathcal{J})$ l’ensemble des parties de \mathcal{J} de cardinal k , et similairement $\mathcal{P}_{<k}(\mathcal{J})$ l’ensemble des parties de cardinal $< k$. Avec ces notations, l’évènement $\{S = k\}$ se décompose comme l’union disjointe d’évènements de la forme

$$\{A_i \text{ pour tout } i \in J \text{ et non-}A_i \text{ pour tout } i \in \mathcal{J} \setminus J\} \quad (\text{IK})$$

pour $J \in \mathcal{P}_k(\mathcal{J})$.

Or, par indépendance, la probabilité de l’évènement ci-dessus vaut

$$\prod_{i \in J} p_i \times \prod_{i \in \mathcal{J} \setminus J} (1 - p_i) = \prod_{i \in J} \frac{p_i}{1 - p_i} \times \prod_{i \in \mathcal{J}} (1 - p_i), \quad (\text{IL})$$

ce qui est (quasiment) égal à $\prod_{i \in J} p_i \times e^{-\theta}$: en effet, d’une part le produit $\prod_{i \in J} (1 - p_i)^{-1}$ est quasiment égal à 1 (du fait que tous les $(1 - p_i)^{-1}$ sont très proches de 1, et qu’on a multiplié seulement k facteurs) ; et d’autre part, le produit $\prod_{i \in \mathcal{J}} (1 - p_i)$ s’approxime par $\prod_{i \in \mathcal{J}} e^{-p_i} = \exp(-\sum_{i \in \mathcal{J}} p_i) \stackrel{\text{dét}}{=} e^{-\theta}$. Finalement, on a déduit que, aux approximations près,

$$\mathbb{P}(S = k) = e^{-\theta} \sum_{J \in \mathcal{P}_k(\mathcal{J})} \prod_{i \in J} p_i. \quad (\text{IM})$$

Ci-dessus, nous avons une somme qui porte sur des *ensembles* de taille k . Il serait plus pratique d’avoir une somme portant sur des k -uplets : en effet, si (i_0, \dots, i_{k-1}) est un k -uplet vérifiant $\{i_0, \dots, i_{k-1}\} = J$, on pourra ré-écrire $\prod_{i \in J} p_i$ comme $\prod_{q=0}^{k-1} p_{i_q}$. Mais attention, si nous faisons la somme sur tous les k -uplets de \mathcal{J}^k , nous allons devoir prendre en compte deux subtilités :

- D'une part, que chaque ensemble J sera encodé par $k!$ k -uplets distincts, selon l'ordre dans lequel on listera ses éléments ;
- D'autre part, que les k -uplets comportant au moins une paire de composantes identiques ne correspondront à aucun J de \mathcal{P}_k .

La prise en compte de ces deux subtilités nous conduit alors à la formule :

$$n!e^\theta \mathbb{P}(S = k) = \sum_{\vec{i} \in \mathcal{I}^k} \prod_{q=0}^{k-1} p_{i_q} - \sum_{\vec{i} \in \Delta} \prod_{q=0}^{k-1} p_{i_q}, \quad (\text{IN})$$

où la notation ' \vec{i} ' désigne un k -uplet (i_0, \dots, i_{k-1}) , et où Δ désigne l'ensemble des k -uplets de \mathcal{I}^k ayant au moins deux composantes identiques.

Dans la formule ci-dessus, la première somme du membre de droite se factorise en

$$\prod_{q=0}^{k-1} \sum_{i_q \in \mathcal{I}} p_{i_q} \stackrel{\text{d'éf}}{=} \prod_{q=0}^{k-1} \theta = \theta^k, \quad (\text{IO})$$

ce qui est précisément ce dont nous avons besoin pour arriver à notre formule de la loi de Poisson : ainsi, il ne reste plus qu'à démontrer que la seconde somme est négligeable pour conclure la démonstration !

Nous souhitions donc, pour finir, contrôler $\sum_{\vec{i} \in \Delta} \prod_{q=0}^{k-1} p_{i_q}$. Pour $(i_0, \dots, i_{k-1}) =: \vec{i}$ un élément de Δ , notons $E(\vec{i}) := \{i_0, \dots, i_{k-1}\}$: du fait que nous avons supposé $\vec{i} \in \Delta$, nous avons $|E(\vec{i})| < k$, et donc (puisque tous les p_i sont $\leq \varepsilon$, et que le produit en version k -uplet compte au moins un facteur en double)

$$\prod_{q=0}^{k-1} p_{i_q} \leq \varepsilon \prod_{j \in E(\vec{i})} p_j. \quad (\text{IP})$$

En sommant terme à terme, nous obtenons ainsi que

$$\sum_{\vec{i} \in \Delta} \prod_{q=0}^{k-1} p_{i_q} \leq \varepsilon \sum_{\vec{i} \in \Delta} \prod_{j \in E(\vec{i})} p_j. \quad (\text{IQ})$$

Dans la somme ci-dessus, nous avons sommé " \vec{i} par \vec{i} " ; nous allons maintenant la ré-écrire pour la sommer " $E(\vec{i})$ par $E(\vec{i})$ ". Notons déjà que tous les $E(\vec{i})$ de la somme ci-dessus sont dans $\mathcal{P}_{<k}(\mathcal{I})$. Néanmoins, pour $J \in \mathcal{P}_{<k}(\mathcal{I})$, il peut y avoir plusieurs k -uplets \vec{i} tels que $E(\vec{i}) = J$... Mais nous savons qu'il y en a au plus $(k-1)^k$, puisqu'on tel \vec{i} doit nécessairement avoir toutes ses composantes dans J ! On a donc la majoration suivante :

$$\sum_{\vec{i} \in \Delta} \prod_{j \in E(\vec{i})} p_j \leq (k-1)^k \sum_{J \in \mathcal{P}_{<k}(\mathcal{I})} \prod_{j \in J} p_j. \quad (\text{IR})$$

Mais ici nous observons que, grâce à l'identité (IM) établie quelques lignes plus haut, la somme de produits qui apparait dans le membre de droite n'est autre que $e^\theta \mathbb{P}(S < k)$, et se trouve donc en particulier être inférieure à e^θ ! En fin de compte, nous avons obtenu que $\sum_{\vec{i} \in \Delta} \prod_{q=0}^{k-1} p_{i_q} \leq \varepsilon(k-1)^k e^\theta$, qui est bien négligeable devant 1 (dans la mesure où les valeurs k et θ , elles, sont supposées modérées) : cela conclut notre démonstration \smile

La variance d'une loi de Poisson peut être retrouvée par un argument particulièrement court quoique astucieux. Je ne vous demande pas de savoir la retrouver de but en blanc, mais éventuellement à l'aide d'une petite indication :

Théorème (IS). *Pour tout $\theta \in \mathbb{R}_+$,*

$$\text{Var}(\text{Poisson}(\theta)) = \theta. \quad (\text{IT})$$

\diamond

Démonstration. Nous avons dit qu'une v.a. de Poisson était égale à la somme d'un grand nombre de v.a. de Bernoulli indépendantes de petit paramètre. Or pour une loi de Bernoulli de paramètre p , l'espérance vaut p tandis que la variance vaut $p(1-p)$: par conséquent, pour p très proche de 0, l'espérance est pratiquement égale à la variance. Or l'espérance d'une somme de v.a. (indépendantes ou pas) est égale à la

somme de leurs espérances, tandis que la variance d'une somme de v.a. indépendantes est égale à la somme de leurs variances. Dès lors, puisqu'ici nous sommes partis de v.a. dont les espérances étaient égales aux variances respectives, et que ces variables sont indépendantes, l'espérance de la somme (donc, de la loi de Poisson) sera égale à la variance de la somme. D'où le résultat annoncé. \diamond

Lois géométriques

Notation (IU). Attention ! Dans ce cours (ainsi que certains logiciels, par exemple R), la définition de la loi géométrique sera différente de celle que vous utilisiez jusqu'à présent : alors qu'on vous avait défini la loi géométrique de paramètre p comme le nombre d'*essais* nécessaires avant d'obtenir un succès lorsque chaque essai réussit avec probabilité p , ici $\text{Géom}^{\text{que}}(p)$ désignera le nombre d'*échecs* avant le premier succès ! Autrement dit, si l'on désigne par ' \mathcal{G} ' la loi géométrique "version prépa", nous prendrons ici $\text{Géom}^{\text{que}}(p) := \mathcal{G}(p) - 1$. (au sens de « mesure-image de la loi $\mathcal{G}(p)$ par la fonction "soustraire 1" »).

(Cependant, dans la mesure où ce choix crée un risque d'erreur particulièrement élevé, je rappellerai systématiquement la convention retenue si des lois géométriques interviennent dans des exercices de travaux dirigés ou des problèmes d'examen ☺).

! **Définition (IV)** (Loi géométrique). Pour $p \in]0, 1]$, la *loi géométrique* de paramètre p , notée $\text{Géom}^{\text{que}}(p)$, est le nombre d'échecs avant d'obtenir un succès, lorsqu'on retente indéfiniment et indépendamment une expérience ayant une probabilité de succès p . Autrement dit, lorsque les $(X_i)_{i \in \mathbb{N}}$ sont des v.a.i.i.d. Bernoulli(p), c'est la loi de $\min\{i \in \mathbb{N} \mid X_i = 1\}$. \heartsuit

À partir de cette définition, vous devez être capables de retrouver rapidement l'expression de la fonction de masse d'une loi géométrique :

! **Proposition (IW).** La loi $\text{Géom}^{\text{que}}(p)$ est portée par \mathbb{N} ; et pour tout $n \in \mathbb{N}$, on a :

$$\mathbb{P}(\text{Géom}^{\text{que}}(p) = n) = (1 - p)^n p. \quad (\text{IX})$$

\diamond

Démonstration. Une manière de démontrer cette formule est de commencer par observer qu'on a

$$\mathbb{P}(\text{Géom}^{\text{que}}(p) \geq n) = (1 - p)^n : \quad (\text{IY})$$

En effet, pour X suivant la loi géométrique, dire que $\{X \geq n\}$ revient à dire que les n premières expériences ont été des échecs : chacun de ces expériences étant indépendante et ayant une probabilité d'échec de $1 - p$, le calcul est immédiat. On en déduit ensuite la probabilité que X vaille n en observant que $\{X = n\} = \{X \geq n\} \setminus \{X \geq n + 1\}$, le second ensemble du membre de droite étant inclus dans le premier.

Une autre preuve, peut-être plus directe, consiste à remarquer que notre variable géométrique vaut n si et seulement si les n premières expériences sont ratées *et* que la $(n + 1)$ -ième est réussie. \heartsuit

Remarque (IZ). Il peut être utile de retenir qu'il y a une formule particulièrement simple (et simple à démontrer) pour la probabilité qu'une loi géométrique dépasse une certaine valeur, d'autant qu'il y a une formule extrêmement similaire (et avec une preuve similaire) pour la loi exponentielle. \clubsuit

Remarque (JA). On notera que, puisque p a été supposé strictement positif, la probabilité que les n premières expériences échouent toutes tend vers 0 lorsque n tend vers l'infini, autrement dit qu'il y aura forcément (au sens de « presque-surement ») une expérience qui finira par réussir — ce qui est au demeurant intuitivement évident — : ainsi, la loi géométrique prend toujours une valeur finie et est donc bien portée par \mathbb{N} . \clubsuit

Voyons maintenant les formules sur l'espérance et la variance de la loi géométrique. Je ne vous demande pas d'être capables de remonter directement la formule concernant l'espérance, mais éventuellement avec une petite indication. (Pour la variance, ce serait le cas échéant l'objet d'un exercice guidé).

Théorème (JB). *Pour tout $p \in]0, 1]$,*

$$\mathbb{E}(\text{Géom}^{\text{que}}(p)) = (1 - p) / p \quad (\text{JC})$$

$$\text{Var}(\text{Géom}^{\text{que}}(p)) = (1 - p) / p^2. \quad (\text{JD})$$

◇

Démonstration. Commençons par la preuve concernant l'espérance. Puisque l'expérience réussit avec probabilité p et échoue avec probabilité $1 - p$, il y a en moyenne $(1 - p)/p$ échecs pour un succès ; or la loi géométrique compte bien le nombre d'échecs pour un succès pris au hasard, lorsqu'on associe chaque échec au succès qui le suit directement [¶]. D'où le résultat.

Pour la variance, on écrit que si X suit une loi géométrique, par la formule de transfert :

$$\mathbb{E}(X^2) = \sum_{n=0}^{\infty} n^2 (1 - p)^n p. \quad (\text{JE})$$

Appelons M_2 cette espérance, et notons qu'en effectuant le changement de variable $n \leftarrow n + 1$, on a aussi

$$M_2 = \sum_{n=1}^{\infty} (n - 1)^2 (1 - p)^{n-1} p. \quad (\text{JF})$$

Par différence terme à terme (sachant que toutes nos sommes sont absolument convergentes),

$$\begin{aligned} pM_2 &= M_2 - (1 - p)M_2 = \underbrace{(n^2(1 - p)^n p)}_{=0} \Big|_{n=0} + \sum_{n=1}^{\infty} (n^2 - (n - 1)^2)(1 - p)^n p \\ &= \sum_{n=1}^{\infty} (2n - 1)(1 - p)^n p = \sum_{n=0}^{\infty} (2n - 1)(1 - p)^n p - (-p) = p + \mathbb{E}(2X - 1), \end{aligned}$$

la dernière égalité provenant de la formule de transfert. Or $\mathbb{E}(2X - 1) = 2\mathbb{E}(X) - 1 = (2 - 3p)/p$ en vertu de la formule précédemment calculée, d'où $pM_2 = (2 - 3p + p^2)/p$ et donc $\mathbb{E}(X^2) = (2 - 3p + p^2)/p^2$. Finalement, on conclut que

$$\text{Var}(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2 = \frac{2 - 3p + p^2}{p^2} - \left(\frac{1 - p}{p}\right)^2 = \frac{1 - p}{p^2}. \quad (\text{JG})$$

[¶]. Il convient cependant d'être très précautionneux avec ce genre d'argument, qui peut facilement conduire à des erreurs (qualifiées de « biais ») si on n'est pas assez attentif : par exemple, on pourrait être tenté de suivre l'argument suivant : « lorsqu'une loi géométrique vaut n , cela signifie qu'on a eu 1 succès pour n échecs, soit un taux de succès moyen pour cette tentative de $1/(n + 1)$. Or nous savons qu'en moyenne, le taux de succès vaut p ; lorsqu'on prend l'espérance de $1/(n + 1)$ pour la loi géométrique, on s'attend donc à trouver p ». Or en réalité $\mathbb{E}(1/(\text{Géom}^{\text{que}}(p) + 1))$ ne vaut pas p , mais $p|\ln p| / (1 - p) > p$... Le problème vient de ce que, lorsqu'on parle du taux global d'échecs, les expériences où la loi géométrique a pris une grande valeur (et donc où $1/(n + 1)$ a pris une petite valeur) ont un rôle plus grand que celles où elle a pris une petite valeur, puisqu'il y a eu besoin de plus d'essais dans le premier cas... Ainsi, faire la moyenne sur les essais et faire la moyenne sur les expériences n'est pas la même chose ! En ce qui concerne l'argument présenté dans la démonstration, heureusement, ce biais ne se manifeste pas : en effet, faire la moyenne par succès et faire la moyenne par expérience revient bien au même, puisqu'il y a toujours exactement un succès par expérience ☺



6.3 Les lois normales

Les lois normales sont tellement importantes, notamment en modélisation statistique, qu'elles méritent une section à elles seules ! C'est donc d'elles que nous allons parler dans ce qui suit.

Lois normales unidimensionnelles

Notation (JH). Attention : Dans ce cours, conformément à ce que vous avez vu au lycée et au semestre 5, le second paramètre d'une loi normale unidimensionnelle désignera sa *variance*. Cependant, de nombreux logiciels, par exemple *R*, utilisent pour leur part un paramétrage par *écart-type* ! ♥

! **Définition (JI)** (Loi normale). Pour $\mu \in \mathbb{R}$ et $\sigma \geq 0$, la *loi normale* d'espérance μ et de variance σ^2 ^{[[]]}, notée Normale(μ, σ^2), est la loi suivie par une variable aléatoire vérifiant les propriétés suivantes :

- La variable aléatoire s'écrit comme la somme d'un très grand nombre de « sous-variables » indépendantes ;
- Chaque sous-variable est concentrée dans une très petite plage de valeurs (ce qui est plus ou moins équivalent à dire que chaque sous-variable a un très petit écart-type) ;
- Les sous-variables ne sont pas trop « sauvages », au sens où il est extrêmement rare qu'elles prennent un comportement très différent de leur comportement typique ^[**],

lorsque cette variable aléatoire a pour espérance μ et pour variance σ^2 ^[††]. On a donc *par définition* $\mathbb{E}(\text{Normale}(\mu, \sigma^2)) = \mu$ et $\text{Var}(\text{Normale}(\mu, \sigma^2)) = \sigma^2$. ♥

Remarque (JJ). La loi normale peut être vue comme l'analogie, dans le cas continu, de la loi binomiale dans le cas discret ^[‡‡] ; en particulier, le théorème-limite central montre que pour tout $p \in]0, 1[$, lorsque n tend vers l'infini, la loi Binom^{le}(n, p) devient asymptotiquement « équivalente » à la loi Normale($np, np(1 - p)$). ♣

! **Proposition (JK).** Si X suit la loi Normale(μ, σ^2), alors pour $a, b \in \mathbb{R}$, la variable aléatoire $aX + b$ suit la loi Normale($a\mu + b, a^2\sigma^2$). ◇

Remarque (JL). En fait, il y a juste besoin de retenir que $aX + b$ suit une loi normale : l'espérance et l'écart-type de celle-ci se déduisent resp. de la linéarité de l'espérance et du caractère quadratique de la variance. ♣

[[]]. Il est d'usage, lorsque cela n'alourdit pas trop la notation, d'écrire le second paramètre d'une loi normale unidimensionnelle sous la forme d'un carré. Le paramètre σ représente alors l'*écart-type* de la loi normale, qui, contrairement à la variance, a le bon goût d'avoir la même homogénéité physique que la variable aléatoire elle-même, et donc d'être plus facile à interpréter intuitivement.

[**]. Je ne préciserai pas ce point ici : pour les curieux, *confer* l'énoncé du théorème-limite central dans le chapitre 5.

[††]. σ correspond donc à l'écart-type, sous réserve qu'il ait été pris positif évidemment.

[‡‡]. Par contre, le paramétrage des deux lois est différent : il n'y a pas d'homologie entre n et p d'une part, et μ et σ d'autre part.

Démonstration. Si nous notons $(X_i)_{i \in I}$ les sous-variables dont X est la somme, on peut voir $aX + b$ comme la somme des sous-variables $aX_i + b$. Il est alors facile de s'assurer que ces sous-variables vérifient elles aussi les hypothèses de la définition (JI), de sorte que $aX + b$ est bien aussi une loi normale. ♣

Remarque (JM). Attention, concernant les écarts-types, on a $\text{Var}^{1/2}(aX + b) = (a^2\sigma^2)^{1/2} = |a|\sigma$, avec apparition d'une valeur absolue! Pour éviter ce piège, bien garder à l'esprit qu'un écart-type est nécessairement positif. ♣

Remarque (JN). On observera en particulier que si X suit une loi Normale(μ, σ^2), alors la variable symétrique de X par rapport à μ , c.-à-d. $\mu - (X - \mu) = 2\mu - X$, suit la même loi que X . Autrement dit, la loi Normale(μ, σ^2) est *symétrique* par rapport à l'abscisse μ (cela se voit bien notamment sur le tracé de la densité), ce qui est une propriété qu'il est important de garder à l'esprit pour visualiser ce qui se passe. ♣

Corolaire (JO). Pour $\mu \in \mathbb{R}$ et $\sigma \in \mathbb{R}_+^*$, si X suit la loi Normale($0, 1$), alors $\sigma X + \mu$ suit la loi Normale(μ, σ^2). “Réciproquement”, si Y suit la loi Normale(μ, σ^2), alors $(Y - \mu) / \sigma$ suit la loi Normale($0, 1$). ♣ !

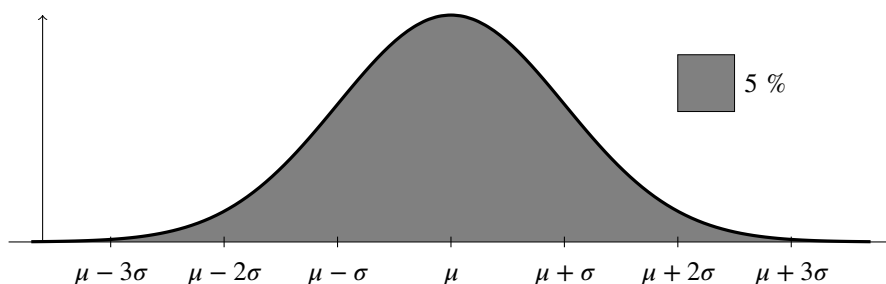
Remarque (JP). On voit donc que toutes les lois normales (à l'exception du cas où l'écart-type est nul, qui est dégénéré) se déduisent les unes des autres par des transformations affines : elles ont donc toutes la même *forme*. En outre, en termes de typologie des paramètres (confer § 6.1 *supra*), on peut observer que l'espérance est un paramètre de position, tandis que la variance est un paramètre d'échelle au carré (et l'écart-type un paramètre d'échelle “tout court” [*]). ♣

Définition (JQ) (Loi normale standard). La loi Normale($0, 1$) est appelée *loi normale standard*. ♣ !

Remarque (JR). En vertu de la Proposition (JK), il suffit de comprendre le comportement de la loi normale standard pour en déduire le comportement de toutes les autres lois normales. Par exemple, le fait qu'une loi normale standard ait 95 % de chances de tomber dans l'intervalle $[-1,96, 1,96]$ se traduit par le fait qu'une loi Normale(μ, σ^2) a 95 % de chances de tomber dans l'intervalle $[\mu - 1,96\sigma, \mu + 1,96\sigma]$. Cela était particulièrement important à l'époque où les propriétés (fonctions de répartition et quantiles) d'une loi de probabilité se retrouvaient à l'aide de *tables* : autant il aurait été impossible de tabuler toutes les lois normales envisageables, autant il suffisait de simplement tabuler la loi normale standard pour pouvoir faire tous les calculs!

Aujourd'hui, la plupart des logiciels qui permettent de faire des calculs sur la loi normale implémentent en interne la façon de passer d'une loi normale standard à une loi normale quelconque (c'est notamment le cas de R), de sorte qu'il est moins critique de savoir faire cette transformation ; cela dit, je ne saurais que trop vous recommander d'avoir quand même bien en tête la façon dont une loi normale quelconque correspond à une loi normale standard, parce que cela est très utile pour visualiser ce qui se passe! ♣

[*]. Attention toutefois : de manière générale, lorsqu'on applique la transformation affine $x \mapsto ax + b$ à une variable aléatoire gaussienne d'écart-type σ , l'image de cette v.a. aura pour écart-type $|a|\sigma$, ce qui n'est pas la même chose que $a\sigma$ lorsque $a < 0$; mais il n'en demeure pas moins que σ est bien un paramètre d'échelle lorsqu'on s'intéresse aux transformations affines *croissantes*. !

FIGURE 6.1 – Allure de la loi Normale(μ, σ^2)

Proposition (JS). Si X et Y sont deux variables aléatoires indépendantes de lois resp. Normale(μ, v) et Normale(v, w), alors $X + Y$ suit la loi Normale($\mu + v, v + w$). \diamond

Remarque (JT). En fait, il y a juste besoin de retenir que $X + Y$ suit une loi normale : l'espérance et la variance de celle-ci se déduisent resp. de la linéarité de l'espérance et de l'additivité de la variance (attention : de la *variance*, pas de l'écart-type!) pour des v.a. indépendantes. \clubsuit

Démonstration. On peut fabriquer la loi de $X + Y$ en prenant les sous-variables de X et celles de Y indépendantes ; et on se retrouve alors avec un seul grand ensemble de sous-variables qui satisfait toutes les conditions pour que sa somme, c.-à-d. $X + Y$, suive une loi normale. \spadesuit

Théorème (JU). Pour $\mu \in \mathbb{R}, \sigma \in \mathbb{R}_+^*$, la loi Normale(μ, σ^2) est à densité sur \mathbb{R} , avec (notant dx un voisinage infinitésimal de x) :

$$\mathbb{P}(\text{Normale}(\mu, \sigma^2) \in dx) = \frac{1}{\sqrt{2\pi}} \sigma^{-1} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \text{vol}_1(dx). \quad (\text{JV})$$

 \diamond

Remarque (JW). En revanche, il n'y a pas de formule fermée concernant la fonction de répartition de la loi normale, c.-à-d. pour calculer $\mathbb{P}(\text{Normale}(\mu, \sigma^2) \leq x)$. De nombreux logiciels implémentent tout de même le calcul de cette fonction à l'aide d'une fonction « spéciale » : soit la fonction de répartition de la loi normale standard, généralement notée Φ , soit la « fonction d'erreur », notée erf, définie par $\text{erf}(x) := 2\Phi(\sqrt{2}x) - 1$. \clubsuit

La figure 6.1 vous donne l'allure des lois normales. Notez en particulier que, bien que la loi normale ne soit pas bornée *stricto sensu*, la « courbe en cloche » formée par cette loi s'atténue néanmoins extrêmement vite aux extrêmes !

Il peut être utile de retenir, au moins approximativement, les valeurs suivantes :

Proposition (JX). On a les valeurs suivantes (arrondies par excès) :

$$\mathbb{P}(|\text{Normale}(0, 1)| > 1) = 32 \% ; \quad (\text{JY})$$

$$\mathbb{P}(|\text{Normale}(0, 1)| > 2) = 4,6 \% ; \quad (\text{JZ})$$

$$\mathbb{P}(|\text{Normale}(0, 1)| > 3) = 2,7 \text{‰}. \quad (\text{KA})$$

En pratique, on pourra retenir qu'il y a environ une chance sur 3 qu'une loi normale tombe à plus d'un écart-type de son espérance, seulement une chance sur 20 qu'elle tombe à plus de deux écarts-types, et à peine une chance sur 400 qu'elle tombe à plus de trois écarts-types ! \diamond

L'intervalle de fluctuation au risque 5 % pour les lois normales est tellement utilisé en pratique que, même si je ne vous demande pas formellement de l'apprendre par cœur, on peut le considérer comme relevant des connaissances standard exigibles de tout ingénieur :

Proposition (KB). *L'intervalle de fluctuation au risque 5 % pour la loi normale standard est $[\pm 1,96]$. Corolairement, l'intervalle de fluctuation au risque 5 % pour la loi Normale(μ, σ^2) est $[\mu \pm 1,96 \sigma]$.* \diamond

Lois normales multidimensionnelles

Définition (KC) (Loi normale multidimensionnelle). Soit $d \in \mathbb{N}$, $\vec{\mu} \in \mathbb{R}^d$ et $\mathbf{C} \in \mathbb{R}^{d \times d}$ une matrice symétrique positive (au sens des formes quadratiques). La loi normale d -dimensionnelle de paramètres $\vec{\mu}$ et \mathbf{C} , notée Normale($\vec{\mu}, \mathbf{C}$)^[†], est la loi suivie par une variable aléatoire vérifiant les propriétés suivantes :

- La variable aléatoire s'écrit comme la somme d'un très grand nombre de « sous-variables » (à valeurs dans \mathbb{R}^d) indépendantes ;
- Chaque sous-variable est concentrée dans une très petite plage de valeurs ;
- Les sous-variables ne sont pas trop « sauvages »,

lorsque cette variable aléatoire a pour espérance (au sens vectoriel) $\vec{\mu}$ et pour matrice de covariance \mathbf{C} . On a donc par définition $\mathbb{E}(\text{Normale}(\vec{\mu}, \mathbf{C})) = \vec{\mu}$ et $\text{Var}(\text{Normale}(\vec{\mu}, \mathbf{C})) = \mathbf{C}$. \heartsuit

Remarque (KD). Notez que je note la « matrice de covariance » comme une variance, car une matrice de covariance n'est jamais que l'analogue multidimensionnel d'une variance : par exemple, si \vec{X} est à valeurs vectorielles, on a, en termes matriciels, $\text{Var}(\vec{X}) = \mathbb{E}(\vec{X}\vec{X}^\top) - \mathbb{E}(\vec{X})\mathbb{E}(\vec{X})^\top$. \clubsuit

Théorème (KE). *Pour $\mathbf{A} \in \mathbb{R}^{k \times d}$ une matrice et $\vec{b} \in \mathbb{R}^k$, la mesure-image de la loi normale ci-dessus par l'application $\vec{x} \mapsto \mathbf{A}\vec{x} + \vec{b}$ est la loi normale de dimension k suivante :*

$$\mathbf{A} \times \text{Normale}(\vec{\mu}, \mathbf{C}) + \vec{b} = \text{Normale}(\mathbf{A}\vec{\mu} + \vec{b}, \mathbf{A}\mathbf{C}\mathbf{A}^\top). \quad (\text{KF})$$

\diamond

Remarque (KG). À nouveau, le seul point vraiment important dans ce théorème est le fait que la mesure-image soit également normale : les valeurs de ses espérance et variance respectives sont en fait des conséquences de la linéarité de l'espérance et de la bilinéarité de la variance, mais écrites ici dans un cadre multidimensionnel. \clubsuit

Théorème (KH). *Si la matrice \mathbf{C} est définie positive, alors la loi normale multidimensionnelle est à densité par rapport à la mesure de Lebesgue ; et pour dx un voisinage infinitésimal de $\vec{x} \in \mathbb{R}^d$, on a*

$$\mathbb{P}(\text{Normale}(\vec{\mu}, \mathbf{C}) \in dx) = ((2\pi)^d \det \mathbf{C})^{-1/2} \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu})^\top \mathbf{C}^{-1}(\vec{x} - \vec{\mu})\right) \text{vol}_d(dx). \quad (\text{KI})$$

[†]. Il n'est pas nécessaire de préciser la valeur de d , celle-ci correspondant nécessairement à la dimension commune de $\vec{\mu}$ et de \mathbf{C} .

En particulier, les lignes de niveau de la densité sont des ellipsoïdes, centrées sur $\vec{\mu}$, dont les axes sont orientés selon les vecteurs propres de la matrice \mathbf{C} et ont des dimensions proportionnelles aux racines carrées des valeurs propres respectives associées. \diamond

Conditionnement gaussien

Il est important de savoir que les vecteurs gaussiens se comportent très bien par conditionnement :

Théorème (KJ). Soit $(\vec{X}; \vec{Y})$ un vecteur gaussien de dimension $n + m$, où \vec{X} est le sous-vecteur constitué par les n premières coordonnées et \vec{Y} le sous-vecteur constitué par les m dernières coordonnées. Alors il existe un vecteur $\vec{\mu}' \in \mathbb{R}^n$ et des matrices $\mathbf{L} \in \mathbb{R}^{n \times m}$ et $\mathbf{V}' \in \mathbb{R}^{n \times n}$ tels que, pour tout $\vec{y} \in \mathbb{R}^m$:

$$\text{Loi}(\vec{X} \mid \vec{Y} = \vec{y}) = \text{Normale}(\mathbf{L}\vec{y} + \vec{\mu}', \mathbf{V}'). \quad (\text{KK})$$

En outre, si la loi de $(\vec{X}; \vec{Y})$ s'écrit par blocs :

$$\text{Loi} \begin{pmatrix} \vec{X} \\ \vec{Y} \end{pmatrix} = \text{Normale} \left(\begin{pmatrix} \vec{\mu} \\ \vec{v} \end{pmatrix}, \begin{pmatrix} \mathbf{V} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{W} \end{pmatrix} \right) \quad (\text{KL})$$

et que \mathbf{W} est de rang plein^[‡], on a $\mathbf{L} = \mathbf{C}\mathbf{W}^{-1}$, $\vec{\mu}' = \vec{\mu} - \mathbf{C}\mathbf{W}^{-1}\vec{v}$ et $\mathbf{V}' = \mathbf{V} - \mathbf{C}\mathbf{W}^{-1}\mathbf{C}^\top$. \diamond

Remarque (KM). En fait, la seconde partie du théorème donnant les formules pour \mathbf{L} , $\vec{\mu}'$ et \mathbf{V}' peut se retrouver à partir des formules de (bi)linéarité de l'espérance et de la covariance : en effet, si la première formule du théorème est vraie, cela signifie que $(\vec{X}; \vec{Y})$ a la même loi que $(\mathbf{L}\vec{Y} + \vec{\mu}' + \vec{\Xi}; \vec{Y})$ pour $\vec{\Xi}$ un vecteur indépendant de \vec{Y} de loi Normale($\vec{0}_n, \mathbf{V}'$) ; or on calcule par (bi)linéarité que

$$\text{Loi} \begin{pmatrix} \mathbf{L}\vec{Y} + \vec{\mu}' + \vec{\Xi} \\ \vec{Y} \end{pmatrix} = \text{Normale} \left(\begin{pmatrix} \mathbf{L}\vec{v} + \vec{\mu}' \\ \vec{v} \end{pmatrix}, \begin{pmatrix} \mathbf{L}\mathbf{W}\mathbf{L}^\top + \mathbf{V}' & \mathbf{L}\mathbf{C} \\ \mathbf{C}\mathbf{L}^\top & \mathbf{W} \end{pmatrix} \right); \quad (\text{KN})$$

et puisque cela doit correspondre à la loi de $(\vec{X}; \vec{Y})$, on en déduit les formules caractérisant \mathbf{L} , \mathbf{V}' et $\vec{\mu}'$ par identification, formules qu'il n'y a ensuite plus qu'à résoudre. \clubsuit

6.4 Autres lois continues remarquables

Lois uniformes

! **Définition (KO)** (Loi uniforme sur un intervalle). Pour $a, b \in \mathbb{R}$ avec $a < b$, la loi uniforme sur $]a, b[$, notée $\text{Unif}^{\text{me}}(a, b)$ ^[§], est la distribution de probabilité portée par $]a, b[$ qui attribue à tout sous-intervalle de $]a, b[$ une masse proportionnelle à sa longueur. \heartsuit

Remarque (KP). La loi uniforme sur un segment est clairement l'homologue, dans le cas continu, de la loi uniforme sur un intervalle d'entiers dans le cas discret : on peut même démontrer que, dans un certain sens, la loi uniforme sur $\llbracket a, b \rrbracket$ tend vers la loi uniforme sur $]a, b[$ lorsque $b - a$ tend vers l'infini. \clubsuit

! **Proposition (KQ).** Pour $a < b$, la loi $\text{Unif}^{\text{me}}(a, b)$ est à densité sur \mathbb{R} , avec (notant dx un voisinage infinitésimal de x) :

$$\mathbb{P}(\text{Unif}^{\text{me}}(a, b) \in dx) = \begin{cases} \frac{\text{vol}_1(dx)}{b-a} & \text{pour } x \in]a, b[; \\ 0 & \text{pour } x \notin]a, b[. \end{cases} \quad (\text{KR})$$

On calcule également facilement la fonction de répartition de la loi :

$$\mathbb{P}(\text{Unif}^{\text{me}}(a, b) \leq x) = \begin{cases} 0 & \text{pour } x \leq a; \\ \frac{x-a}{b-a} & \text{pour } x \in [a, b]; \\ 1 & \text{pour } x \geq b. \end{cases} \quad (\text{KS})$$

◇

Remarque (KT). Dans la formule de densité, ne vous cassez surtout pas la tête à savoir quelle est la densité pour les cas limites où x vaut a ou b : en fait, toutes les conventions sont également valables, puisque ajouter ou retrancher deux masses infinitésimales en a et b n'a aucun effet macroscopique sur la mesure ! D'ailleurs, comme vous l'avez vu au semestre précédent, la définition rigoureuse du concept de densité ne détermine la fonction de densité qu'à égalité presque-partout près, ce qui signifie en particulier que si on change la formule pour la densité en un nombre fini (ou même dénombrable) de points, la nouvelle densité obtenue est une formule tout aussi valable que l'ancienne. (Même si ici, cela serait évidemment assez stupide de vouloir changer la valeur de la densité ailleurs qu'en a ou b). ♣

Remarque (KU). Concernant la formule pour la fonction de répartition, lorsque x vaut a ou b , on se situe dans deux cas de la formule simultanément, mais ces deux cas donnent alors le même résultat : ce qui est logique puisque nos formules doivent décrire une fonction continue (la loi étant à densité, donc diffuse). ♣

Démonstration. On sait que pour $I \subseteq]a, b[$, $\mathbb{P}(\text{Unif}^{\text{me}}(a, b) \in I) = \alpha \text{vol}_1(I)$ (où $\text{vol}_1(I)$ n'est autre que la longueur de I) pour une certaine constante de proportionnalité α ; en outre, puisque la distribution de probabilité est portée par $]a, b[$, on doit avoir $\mathbb{P}(\text{Unif}^{\text{me}}(a, b) \in]a, b[) = 1$, d'où $\alpha = (b - a)^{-1}$. Dès lors, les deux formules de l'énoncé découlent directement de la définition de la loi uniforme, en observant, dans le second cas, que la probabilité donnée à $]-\infty, x]$ n'est autre que la probabilité donnée à $]a, x]$, vu que la mesure est portée par $]a, b[$. ◇

L'espérance de la loi uniforme intervalle s'obtient immédiatement par symétrie :

Proposition (KV). Pour $a \leq b$,

$$\mathbb{E}(\text{Unif}^{\text{me}}(a, b)) = \frac{a + b}{2}. \quad (\text{KW})$$

◇

Pour la variance, c'est un peu plus compliqué :

[‡]. Ce à quoi on peut toujours se ramener, en considérant si nécessaire un sous-vecteur \bar{Y}' de \bar{Y} de taille $\text{rg } \mathbf{W}$ tel que \bar{Y} soit presque-surement une fonction affine de \bar{Y}' .

[§]. À ne pas confondre avec la loi $\text{Unif}^{\text{me}}(\{a, b\})$, qui est la loi uniforme sur la paire $\{a, b\}$!

Théorème (KX). Pour $a \leq b$,

$$\text{Var}(\text{Unif}^{\text{me}}(a, b)) = \frac{1}{12}(b - a)^2. \tag{KY}$$

◇

Démonstration. La façon de faire la plus directe (mais pas la plus élégante [¶]) est de calculer en utilisant la formule de la densité de la loi et la formule de transfert que

$$\mathbb{E}(\text{Unif}^{\text{me}}(a, b)^2) = \int_{x=a}^b x^2 \frac{dx}{b-a} = (b-a)^{-1} \left[\frac{1}{3}x^3 \right]_a^b = \frac{1}{3} \frac{b^3 - a^3}{(b-a)} = \frac{1}{3}(a^2 + ab + b^2), \tag{KZ}$$

et d'en déduire au vu de la formule pour l'espérance que

$$\begin{aligned} \text{Var}(\text{Unif}^{\text{me}}(a, b)^2) &= \frac{1}{3}(a^2 + ab + b^2) - \left(\frac{b+a}{2} \right)^2 \\ &= \left(\frac{1}{3}a^2 + \frac{1}{3}ab + \frac{1}{3}b^2 \right) - \left(\frac{1}{4}a^2 + \frac{1}{2}ab + \frac{1}{4}b^2 \right) = \frac{1}{12}a^2 - \frac{1}{6}ab + \frac{1}{12}a^2 = \frac{1}{12}(b-a)^2. \end{aligned}$$

◇

! **Proposition (LA).** Pour X une variable suivant la loi $\text{Unif}^{\text{me}}(\alpha, \beta)$, pour $a, b \in \mathbb{R}$, la variable aléatoire $aX + b$ suit la loi :

$$\begin{cases} \text{Unif}^{\text{me}}(a\alpha + b, a\beta + b) & \text{si } a > 0; \\ \text{Unif}^{\text{me}}(a\beta + b, a\alpha + b) & \text{si } a < 0; \\ \delta_b & \text{si } a = 0. \end{cases} \tag{LB}$$

◇

Démonstration. Prouver rigoureusement cette proposition demande quelques calculs à l'aide de la formule de changement de variables géométrique, mais le résultat est tellement évident intuitivement que je pense qu'on peut vous demander de le retenir par cœur ☺

◇

Lois exponentielles

! **Définition (LC)** (Loi exponentielle). Pour $\lambda > 0$, la loi exponentielle de taux λ , notée $\text{Expon}^{\text{le}}(\lambda)$, est la loi de la durée de vie d'un individu qui, à tout instant, indépendamment du passé, a une probabilité λ par unité de temps de mourir ; autrement dit, sachant que l'individu est encore en vie au temps t , pour dt un accroissement de temps infinitésimal, l'individu a une probabilité λdt de mourir avant le temps $(t + dt)$. ♡

Remarque (LD). Attention, sauf erreur de ma part, ce n'est pas le même paramétrage que celui que vous avez vu au semestre 5 ! Une convention alternative fréquemment répandue étant en effet de paramétrer les lois exponentielles par leur *espérance de vie*, qui, avec les notations ci-dessus, correspond à $1/\lambda$. ♣

Remarque (LE). La loi exponentielle est l'homologue continu de la loi géométrique dans le cas discret, le rôle de la mort de l'individu remplaçant celui de la réussite de l'expérience. On peut même démontrer que, dans un certain sens, lorsque p tend vers 0, la loi $\text{Géom}^{\text{que}}(p)$ devient asymptotiquement égale à la loi $\text{Expon}^{\text{le}}(p)$. ♣

[¶]. Il est plus habile de commencer par observer que le résultat est nécessairement de la forme $\alpha(b - a)^2$, puis de calculer α sur un cas particulièrement simple, p. ex. pour $]a, b[=]\pm 1[$.

Les formules concernant la fonction de répartition et la densité de la loi exponentielle ne sont pas à connaître par cœur, mais doivent pouvoir être redémontrées au besoin :

Proposition (LF). *La loi $\text{Expon}^{\text{le}}(\lambda)$ est portée par \mathbb{R}_+ , avec, pour tout $x \in \mathbb{R}_+$,*

$$\mathbb{P}(\text{Expon}^{\text{le}}(\lambda) \geq x) = e^{-\lambda x}. \quad (\text{LG})$$

◇

Démonstration. La première chose qu'il faut comprendre, c'est que, puisque l'individu a en permanence la même probabilité de mourir par unité de temps, sa probabilité d'être encore en vie au temps x va décroître exponentiellement. Même si cela se saisit intuitivement, formalisons le raisonnement. Au vu de la définition de la façon dont notre individu meurt, il est clair que, sachant qu'il est encore en vie au temps t , son futur sera exactement le même que s'il venait de naître ! En particulier, sachant qu'il est en vie au temps t , la probabilité qu'il vive encore pendant une durée s (et donc que sa durée de vie totale atteigne $t + s$) est égale à la probabilité, si on le considère à la naissance, qu'il vive pendant une durée s ; ce qu'on peut écrire, en termes mathématiques : pour $X \sim \text{Expon}^{\text{le}}(\lambda)$,

$$\mathbb{P}(X \geq t + s \mid X \geq t) = \mathbb{P}(X \geq s), \quad (\text{LH})$$

ce qui donne « $\mathbb{P}(A_{t+s} \mid A_t) = \mathbb{P}(A_s)$ » en posant $\{X \geq u\} =: A_u$ pour alléger les notations. Utilisant la définition de l'espérance conditionnelle, on a donc

$$\frac{\mathbb{P}(A_t \text{ et } A_{t+s})}{\mathbb{P}(A_t)} = \mathbb{P}(A_s); \quad (\text{LI})$$

où nous notons que la définition des A_u fait que $A_{t+s} \implies A_t$, de sorte qu'on a en fait $\{A_t \text{ et } A_{t+s}\} = A_{t+s}$, et donc finalement

$$\mathbb{P}(A_{t+s}) = \mathbb{P}(A_t) \mathbb{P}(A_s). \quad (\text{LJ})$$

Ce raisonnement peut se comprendre très simplement en disant que pour vivre pendant un temps $t + s$, il faut d'abord vivre pendant un temps t , puis à nouveau ne pas mourir pendant un temps s entre t et $t + s$, et que ces deux événements sont indépendants [III].

Dès lors, l'application $u \mapsto \mathbb{P}(A_u)$ est un morphisme de semigroupes de $(\mathbb{R}_+, +)$ dans $([0, 1], \times)$, et comme ce morphisme est continu [**] puisqu'il n'y a aucun moment où la particule a une probabilité non nulle de mourir, c'est une application de la forme $u \mapsto e^{\tau u}$: c'est pour ainsi dire la définition même de l'exponentielle ! Pour finir, il reste à identifier τ : mais, puisque la probabilité (non conditionnée) de mourir par unité de temps est la dérivée de la probabilité d'être encore en vie, on a que la probabilité par unité de temps de mourir au temps 0 vaut $-(u \mapsto e^{\tau u})'(0) = \tau$ (là encore, c'est essentiellement la définition même de l'exponentielle), et puisque par ailleurs cette probabilité vaut λ , on arrive au résultat annoncé. ◇

[III]. À ceci près que « ne pas mourir entre t et $t + s$ » n'a pas vraiment de sens si on est déjà mort au temps t ; mais on comprend l'idée...

[**]. En fait, on n'a même pas besoin de la continuité : le fait que l'application soit décroissante — ce qui est trivial — suffit.

Corolaire (LK). La loi $\text{Expon}^{\text{le}}(\lambda)$ prend presque-surement des valeurs strictement positives, et est à densité sur \mathbb{R}_+^* , avec, notant dx un voisinage infinitésimal de x :

$$\mathbb{P}(\text{Expon}^{\text{le}}(\lambda) \in dx) = \lambda e^{-\lambda x} \text{vol}_1(dx) \quad (\text{LL})$$

◇

Démonstration. Le fait que la loi exponentielle ne prenne presque-jamais la valeur nulle vient de ce que la probabilité de mourir possède par définition une certaine densité par rapport au temps : dès lors, la probabilité de mourir à un instant bien précis est nulle.

Concernant la formule pour la densité, il y a deux façons de la prouver. La première est de prendre, pour se fixer les idées, $dx = [x, x + h[$ (avec h infinitésimalement petit) et de dire que, pour mourir entre les temps x et $x + h$, il faut commencer par survivre jusqu'au temps x (ce qui, en vertu de la proposition (LF), se produit avec une probabilité $e^{-\lambda x}$), puis, sachant cela, mourir pendant l'intervalle de temps infinitésimal de durée h suivant, ce qui, par définition de la loi exponentielle, se produit avec probabilité λh (on s'est placé ici dans le cas $x \geq 0$, le cas $x < 0$ étant trivial). En multipliant ces deux expressions, on obtient bien le résultat annoncé.

Une autre façon de prouver ce résultat est un raisonnement de probabilité très classique que vous devez absolument savoir maîtriser. On fixe à nouveau dx de la forme $[x, x + h[$, et on observe que, pour X une variable aléatoire réelle,

$$\{X \in [x, x + h[\} = \{X \geq x\} \text{ et } \neg\{X \geq x + h\}. \quad (\text{LM})$$

Or on a $\{X \geq x + h\} \implies \{X \geq x\}$, donc la probabilité de l'évènement de gauche s'obtient par différence :

$$\begin{aligned} \mathbb{P}(\text{Expon}^{\text{le}}(\lambda) \in [x, x + h[) \\ = \mathbb{P}(\text{Expon}^{\text{le}}(\lambda) \geq x) - \mathbb{P}(\text{Expon}^{\text{le}}(\lambda) \geq x + h) = e^{-\lambda x} - e^{-\lambda(x+h)} \end{aligned}$$

en vertu de la proposition (LF). On voit que cette expression est de la forme $f(x + h) - f(x)$ avec $f(u) = -e^{-\lambda u}$; par développement au premier ordre, elle vaut donc $f'(x)h$, ce qui est bien, à nouveau, le résultat annoncé. ♡

Théorème (LN). Pour $\lambda > 0$,

$$\mathbb{E}(\text{Expon}^{\text{le}}(\lambda)) = 1/\lambda ; \quad (\text{LO})$$

$$\text{Var}(\text{Expon}^{\text{le}}(\lambda)) = 1/\lambda^2. \quad (\text{LP})$$

◇

Remarque (LQ). Ces résultats sont analogues des formules pour la loi géométrique; et on aurait ainsi pu retrouver la formule pour l'espérance par le même genre de démonstration que pour la loi géométrique; cependant il m'a semblé que c'était nettement plus difficile à visualiser dans le cas continu, et j'ai donc choisi de ne pas vous imposer de savoir retrouver cette formule. ♡

Remarque (LR). En termes de typologie des paramètres, il n'est pas difficile de vérifier que le paramètre de taux d'une loi exponentielle correspond à un paramètre d'échelle inverse. (Les transformations affines à considérer pour cette famille de lois étant les transformations linéaires strictement croissantes). ♡

Deuxième partie

Statistique bayésienne

Chapitre 7

Concept d'inférence statistique

7.1 Motivation : Deux exemples

Afin de motiver la notion d'inférence statistique, nous allons ci-dessous présenter deux exemples. Ces deux exemples, après formalisation, nous accompagneront tout au long du polycopié.

Exemple (LS) (Le (bon ?) chasseur). Un candidat souhaite rejoindre le club des chasseurs du Bouchonnois. Le club ayant une politique de haute exigence, il souhaite savoir avec quelle précision le candidat sait se servir de son arme en vue d'un éventuel recrutement.

La « précision » d'un tireur est évaluée via un indicateur que nous appellerons *fiabilité*, qui est défini comme la proportion de plateaux^[*] que ce tireur est capable de toucher au balltrap dans certaines conditions de référence. Tout tireur possède une certaine fiabilité, qui est par construction un réel de $[0, 1]$; et en l'occurrence, nous appellerons $\theta_{\mathcal{J}}$ la fiabilité de notre candidat. Ici il est essentiel de comprendre que, bien que le nombre $\theta_{\mathcal{J}}$ soit défini sans ambiguïté, on ne peut pas déterminer sa valeur directement : en effet, quand on parle de la « proportion de plateaux que le tireur est capable de toucher », cela est à entendre dans le sens d'un nombre de plateaux *tendant vers l'infini*; mais dans les situations pratiques, il faudrait plutôt dire que le tireur a une *probabilité* $\theta_{\mathcal{J}}$ de toucher chaque plateau!

Au niveau de la procédure de recrutement, pour se faire une idée de la fiabilité du candidat, le club va lui demander de tirer sur une série de $25 =: n$ plateaux. Le club regardera le nombre de plateaux qui ont été touchés au cours de ce test, nombre que nous noterons X ; avec l'espoir qu'on pourra, grâce à la connaissance de X , se faire une « idée » assez fiable de ce que vaut $\theta_{\mathcal{J}}$.

Par ailleurs, le club envisage, s'il recrute le candidat, de le faire participer à une compétition dans laquelle il faudra tirer sur $75 =: m$ plateaux (dans les mêmes conditions que pour le test) : ainsi, mieux cerner la fiabilité du candidat permettra également au club d'avoir une idée de score que le candidat réalisera à la compétition si on l'y inscrit! ♣

Remarque (LT). Il importe de bien insister sur le fait que X ne permet pas de reconstituer $\theta_{\mathcal{J}}$ à coup sûr, du fait qu'il s'agit d'une *variable aléatoire* : plus précisément, en supposant que le fait que chaque plateau de l'épreuve soit ou non touché est indépendant, X suit la loi Binom^{le} $(n, \theta_{\mathcal{J}})$; or tout ce que le club voit à

[*]. Au balltrap, un *plateau* est une cible volante en argile, qui explose lorsqu'elle est touchée par une balle.

l'issue du test, c'est *une* réalisation de X ! Par exemple, mettons que le club voie que le candidat a touché 8 =: x_{\checkmark} plateaux. Vu que, pour toute valeur de θ entre 0 et 1 (exclus), il existe une probabilité non nulle que la loi Binom^{le}(n, θ) prenne la valeur 8, cette observation ne nous permet d'exclure *formellement* aucune valeur pour θ_{\checkmark} : peut-être que θ_{\checkmark} vaut 20 % et que le candidat a été chanceux lors du test ; peut-être que θ_{\checkmark} vaut 30 % et que le candidat a fait un score tout à fait conforme à ce à quoi il pouvait s'attendre ; peut-être que θ_{\checkmark} vaut 75 % et que le candidat a été très malchanceux... Mais bien sûr, certaines de ces possibilités semblent plus *raisonnables* que d'autres ; et c'est bien en ce sens qu'on pourra se faire une idée assez fiable de θ_{\checkmark} malgré tout ! ♣

Exemple (LU) (Le pédagogue). Un enseignant est chargé de donner un cours de thermodynamique avancée à un groupe d'élèves ingénieurs. La méthode d'évaluation des élèves en fin de cours est standardisée selon un protocole précis, et donne pour chaque élève un « score » figuré par un nombre réel. Traditionnellement, ce cours était enseigné selon une certaine méthode pédagogique ; et la très longue expérience sur cette méthode avait montré que le score d'un élève pris au hasard formé selon la méthode traditionnelle suit une loi Normale($\mu_{\text{réf}}, \sigma_{\text{réf}}^2$) avec $\mu_{\text{réf}} = 75$ et $\sigma_{\text{réf}} = 11$, les scores des différents élèves étant indépendants^[†].

Notre enseignant souhaite expérimenter une nouvelle méthode pédagogique, dont il espère évidemment qu'elle donnera de meilleurs résultats que l'ancienne. Diverses raisons l'amènent à penser que, avec la nouvelle méthode, les scores des élèves seront toujours indépendants et identiquement distribués selon une loi de la forme Normale($\mu_{\checkmark}, \sigma_{\checkmark}^2$), mais que μ_{\checkmark} et σ_{\checkmark} seront (potentiellement) *différentes* de $\mu_{\text{réf}}$ et $\sigma_{\text{réf}}$!

De même que dans le cas du chasseur, l'enseignant ne peut pas observer directement μ_{\checkmark} et σ_{\checkmark} . Tout ce qu'il peut faire, c'est tester sa méthode sur la promotion d'élèves qu'il a cette année, que nous supposons ici composée de 22 =: n élèves, et regarder les résultats : les scores obtenus seront des v.a. X_0, \dots, X_{n-1} i.i.d. Normale($\mu_{\checkmark}, \sigma_{\checkmark}^2$) — ou encore, de manière équivalente, le n -uplet (X_0, \dots, X_{n-1}) suivra la loi Normale($\mu_{\checkmark}, \sigma_{\checkmark}^2$)^{⊗ n} sur \mathbb{R}^n . Et à partir de ces X_i , l'enseignant tentera de “reconstituer” les valeurs de μ_{\checkmark} et σ_{\checkmark} : là encore, il ne pourra pas le faire avec une certitude absolue, mais certaines possibilités seront plus plausibles que d'autres. ♣

7.2 Modèle d'inférence statistique

Cas d'une inférence explicative

On donne ici les définitions mathématiques de la notion de modèle d'inférence statistique. On distinguera deux situations : celle de l'inférence *explicative*, et celle de l'inférence *prédictive*.

Définition générale

!! **Définition (LV)** (Modèle d'inférence statistique explicative). Un *modèle d'inférence statistique explicative* est la donnée d'un espace Θ , appelé *espace du paramètre*

[†]. Attention : Lorsque nous disons que le score d'un élève est « aléatoire », en l'occurrence cela ne veut pas dire que l'examen donne un résultat au hasard, mais que c'est l'*identité* de l'élève (est-on tombé sur une élève douée ou un élève plus en difficulté) qui est la (principale) source d'aléa.

caché, d'un espace \mathcal{X} , appelé *espace de l'observation*, et d'une application de Θ dans $\mathcal{M}_1(\mathcal{X})$ ^[‡] qui, à chaque $\theta \in \Theta$, associe une loi de probabilité P_θ sur \mathcal{X} , appelée *loi de l'observation sachant que le paramètre caché vaut θ* . ♡

Point (LW) (Variables aléatoires associées à un modèle statistique). Chaque fois que nous introduirons un modèle statistique, nous considérerons automatiquement un espace probabilisé sur lequel nous définirons une variable aléatoire θ à valeurs dans Θ , appelée *paramètre caché* (ou « paramètre caché en tant que variable aléatoire » quand il y aura lieu de préciser), et une variable aléatoire X à valeurs dans \mathcal{X} , appelée *observation* (ou « observation en tant que variable aléatoire »), de façon que pour tout $\theta \in \Theta$, on ait

$$\text{Loi}(X \mid \theta = \theta) = P_\theta. \quad (\text{LX})$$

♣

Principe (LY) (Interprétation d'un modèle d'inférence explicative).

- Dans un modèle d'inférence statistique (*explicative*), l'observation représente ce qu'on est (ou sera) en mesure d'observer au moment de réaliser notre analyse. On disposera donc, au moment de l'analyse, de la réalisation de X , que nous noterons x_\vee , et que nous qualifierons de « valeur effective de l'observation ».
- La réalisation du paramètre caché θ , que nous noterons θ_\vee et que nous appellerons « vraie valeur du paramètre caché », représente les informations qu'il pourrait être utile de connaître (dans le cadre du modèle) sur le comportement de l'observation, mais qu'on ne peut pas observer directement. Dans la mesure où θ_\vee ne peut pas être observé directement, nous sommes obligés d'envisager qu'il puisse prendre toutes sortes de valeurs, d'où l'intérêt d'avoir une vision « variable aléatoire » du paramètre caché pour formaliser cette multitude de possibilités. Dans ce contexte, l'espace du paramètre caché Θ est l'ensemble des valeurs qu'il est envisageable (dans le cadre de notre modèle) que θ puisse valoir.
- La loi de l'observation sachant le paramètre caché représente l'aléa « intrinsèque » sur l'observation, celui qui demeurerait fondamentalement même si on connaissait la vraie valeur du paramètre caché.

◇

Remarque (LZ). Bien comprendre que d'un point de vue « expérimental », quand on considère une situation donnée modélisée par un modèle d'inférence, la vraie valeur θ_\vee du paramètre caché existe « quelque part », avant même que nous observions la valeur de X , et qu'il n'existe aucun aléa ni aucun contrôle sur ce qu'elle vaut ! (en tout cas, pas au moment de commencer l'expérience). Pour autant, nous n'avons aucun moyen d'y accéder directement : tout ce que nous pourrions observer, c'est la valeur effective x_\vee de l'observation. ♣

Exemple (MA) (Formalisation du modèle du chasseur). Le modèle du chasseur, dans sa version explicative, a pour paramètre caché la variable aléatoire θ représentant la *fiabilité* du chasseur ; l'espace du paramètre caché étant $\Theta =]0, 1[$ ^[§].

[‡]. Rappelons que « $\mathcal{M}_1(\mathcal{X})$ » désigne l'ensemble des lois de probabilité sur \mathcal{X} .

[§]. Remarque : J'ai choisi d'exclure les cas-limites $\theta = 0$ et $\theta = 1$, afin d'éviter certains soucis techniques.

L'observation, notée X , représente le nombre de plateaux touchés lors du test. L'espace de l'observation est l'intervalle d'entiers $\llbracket 0, n \rrbracket =: \mathcal{X}$, où n est une valeur connue représentant le nombre d'essais lors du test, que nous prendrons égale à 25 à moins qu'il ne soit explicitement précisé autrement. La loi de l'observation sachant le paramètre caché est

$$\text{Loi}(X \mid \theta = \theta) := \text{Binom}^{\text{le}}(n, \theta). \quad (\text{MB})$$

En ce qui considère le nombre x_{\checkmark} de plateaux touchés par notre chasseur, nous le prendrons égal à 8 pour les applications numériques, à moins qu'il ne soit explicitement précisé autrement. ♣

Exemple (MC) (Formalisation du modèle du pédagogue). Le *modèle du pédagogue*, dans sa version explicative, a pour paramètre caché le couple de variables aléatoires (μ, σ) , où μ et σ représentent respectivement l'*espérance* et l'*écart-type* de la distribution (sur le long terme) des scores des élèves; l'espace du paramètre caché étant $\Theta =: \mathbb{R} \times \mathbb{R}_+^*$ (autrement dit, μ est à valeurs dans \mathbb{R} et σ à valeurs dans \mathbb{R}_+^*). L'observation représente la liste des scores de l'ensemble des élèves (rangés, disons, dans l'ordre alphabétique) : c'est donc un n -uplet de valeurs, de sorte que nous la noterons $\vec{X}_{\llbracket 0, n \rrbracket}$ (ou parfois simplement \vec{X} pour alléger), ou, de façon équivalente, (X_0, \dots, X_{n-1}) (X_0 étant alors le score de la première élève de la liste, X_1 celui du second de la liste, etc.). La valeur n , représentant le nombre d'élèves dans la promotion expérimentant la nouvelle méthode, est connue : on la prendra égale à 22, sauf mention explicite du contraire. Les scores sont supposés être des nombres réels, de sorte que l'espace de l'observation, noté \mathcal{X} , correspond à \mathbb{R}^n . La loi de l'observation sachant le paramètre caché est

$$\text{Loi}(\vec{X} \mid \mu = \mu \text{ et } \sigma = \sigma) := \text{Normale}(\mu, \sigma^2)^{\otimes n} : \quad (\text{MD})$$

autrement dit, si le paramètre caché vaut (μ, σ) , les scores des élèves sont i.i.d. de loi Normale(μ, σ^2). Pour les applications numériques de l'analyse, la réalisation effective de $\vec{X}_{\llbracket 0, n \rrbracket}$, notée $(x_{0\checkmark}, \dots, x_{(n-1)\checkmark})$ ou $\vec{x}_{\llbracket 0, n \rrbracket\checkmark}$, sera prise égale à

$$(91,7, 69,6, 81,5, 45,0, 69,7, 77,3, 42,1, 70,7, 66,9, 93,6, 75,7, \\ 93,7, 59,8, 75,2, 49,8, 88,2, 63,7, 90,0, 54,2, 57,2, 72,8, 58,3). \quad (\text{ME})$$

♣

Cas d'une inférence prédictive

!! **Définition (MF)** (Modèle d'inférence statistique prédictive). Un modèle d'inférence statistique prédictive est la donnée d'un espace du paramètre caché Θ , de deux espaces \mathcal{X} et \mathcal{Y} , appelés resp. *espace de l'observation* (ou « de l'observation passée ») et *espace de l'observation future*, et d'une application de Θ dans $\mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ qui, à chaque $\theta \in \Theta$, associe une loi de probabilité P_θ , appelée *loi de l'observation complétée* sachant que le paramètre caché vaut θ . ♡

! **Définition (MG)** (Variables aléatoires associées à un modèle prédictif). De même que dans le paradigme explicatif, chaque fois que nous introduisons un modèle statistique, nous considérerons automatiquement un espace probabilisé sur lequel nous définirons une variable aléatoire « paramètre caché » θ à valeurs dans Θ ,

et deux variables aléatoires X et Y à valeurs dans resp. \mathcal{X} et \mathcal{Y} , appelées resp. *observation (passée)* et *observation future*, de façon que pour tout $\theta \in \Theta$, la loi *jointe* de (X, Y) conditionnellement à $\{\theta = \theta\}$ soit P_θ :

$$\text{Loi}(X, Y \mid \theta = \theta) := P_\theta. \quad (\text{MH})$$

♡

Principe (MI) (Interprétation d'un modèle prédictif). *En termes d'interprétation, X représente l'observation à laquelle on pourra avoir accès au moment de faire l'analyse (et dont on connaîtra donc la réalisation x_\blacktriangledown), tandis que Y représente une observation à laquelle on pourra éventuellement avoir accès dans un certain futur après l'analyse.* ♠

!!

Exemple (MJ). Dans la version prédictive du modèle du chasseur, outre l'observation (passée) mentionnée en exemple, on introduit aussi une observation future, notée Y , qui représente le nombre de plateaux touchés par le candidat lors de la compétition à venir (s'il y participe). L'espace de l'observation future est l'intervalle d'entiers $\llbracket 0, m \rrbracket =: \mathcal{Y}$, où m représente le nombre de plateaux pour la compétition, qui est connu : pour les applications numériques, nous prendrons $m := 75$, sauf mention explicite du contraire. En version prédictive, la loi de l'observation complétée sachant le paramètre caché est :

$$\text{Loi}(X, Y \mid \theta = \theta) := \text{Binom}^{\text{le}}(n, \theta) \otimes \text{Binom}^{\text{le}}(m, \theta) : \quad (\text{MK})$$

cela signifie concrètement que, si le paramètre caché vaut θ , X et Y sont (conditionnellement) indépendantes et suivent resp. les lois $\text{Binom}^{\text{le}}(n, \theta)$ et $\text{Binom}^{\text{le}}(m, \theta)$; en particulier, pour la vraie valeur du paramètre caché, X et Y sont indépendantes et de lois respectives $\text{Binom}^{\text{le}}(n, \theta_\blacktriangledown)$ et $\text{Binom}^{\text{le}}(m, \theta_\blacktriangledown)$. ♠

Exemple (ML) (Modèle du pédagogue, version prédictive). Dans la version prédictive du modèle du pédagogue, outre l'observation (passée) mentionnée en exemple, on introduit aussi une observation future, notée $\vec{Y}_{\llbracket 0, m \rrbracket}$, qui représente les scores obtenus par la seconde promotion à tester la méthode, à valeurs dans $\mathbb{R}^m =: \mathcal{Y}$, où m représente le nombre d'élèves de la seconde promotion, qui est connu : pour les applications numériques, nous prendrons $m := 16$, sauf mention explicite du contraire. Dans ce cas, la loi de l'observation complétée sachant le paramètre caché est :

$$\text{Loi}(\vec{X}, \vec{Y} \mid \mu = \mu \text{ et } \sigma = \sigma) := \text{Normale}(\mu, \sigma^2)^{\otimes n} \otimes \text{Normale}(\mu, \sigma^2)^{\otimes m}, \quad (\text{MM})$$

ce qui signifie concrètement que, sachant la vraie valeur du paramètre caché, tous les scores des élèves (qu'ils appartiennent à la première ou à la seconde promotion) sont (conditionnellement) i.i.d. $\text{Normale}(\mu_\blacktriangledown, \sigma_\blacktriangledown^2)^{\otimes n}$. ♠

Remarque (MN). Attention ! Dans nos deux exemples favoris, dans le cadre prédictif, l'observation (passée) et l'observation future sont, sachant le paramètre caché, indépendantes : mais cela n'a aucune raison d'être le cas en général ! ♠

Remarque (MO). Attention ! Nonobstant la nomenclature utilisée, l'« observation future » n'est pas forcément quelque chose qui aura lieu plus tard dans le temps, ni même forcément quelque chose qu'on peut effectivement observer ! Ce qui compte, c'est la *structure mathématique* du modèle : l'« observation future », c'est le nom qu'on donne à une quantité sur laquelle subsiste un aléa, même sachant le paramètre caché (et l'observation passée) ! ♠

7.3 Notion d'inférence

Maintenant que nous avons formalisé la notion de modèle statistique, nous sommes prêts à expliquer plus précisément ce qu'on entend par « inférence statistique ».

Notion de quantité d'intérêt

Au préalable, il nous faut introduire le concept de *quantité d'intérêt* :

! **Définition (MP)** (Quantité d'intérêt). Une quantité d'intérêt correspondra à une *fonction* des quantités introduites dans le modèle, impossible à observer directement au moment de l'analyse. Plus précisément, on introduira deux types de quantités d'intérêt :

- Une quantité d'intérêt *explicative* est une fonction du paramètre caché, autrement dit, une v.a. de la forme $\gamma(\theta)$, pour $\gamma(\bullet)$ une fonction de Θ dans un certain espace \mathcal{E} .
- Une quantité d'intérêt *prédictive* est une fonction de l'observation future, autrement dit, une v.a. de la forme $g(Y)$, pour $g(\bullet)$ une fonction de \mathcal{Y} dans un certain espace \mathcal{E} .

L'espace dans lequel une certaine quantité d'intérêt est à valeurs sera appelé « espace de la quantité d'intérêt (en question) ». ♣

Remarque (MQ). L'appellation « quantité d'intérêt » provient de ce que, en pratique, le statisticien s'intéresse surtout au paramètre caché ou l'observation future au travers d'une certaine quantité d'intérêt. Cependant nous parlerons plus généralement de « quantités d'intérêt » dès lors que nous voulons parler d'une quantité de la forme $\gamma(\theta)$ ou $g(Y)$, quand bien même cette quantité ne constituerait qu'une étape intermédiaire de notre raisonnement et que nous ne sommes pas intrinsèquement « intéressés » par elle. ♣

Exemple (MR). Dans le modèle du chasseur, la quantité qui intéresse le plus le club est manifestement la fiabilité du candidat : autrement dit, θ elle-même. Il s'agit bien d'une quantité d'intérêt explicative, puisque c'est l'image du paramètre caché par l'application « identité ». ♣

Exemple (MS). Toujours dans le modèle du chasseur, une autre quantité d'intérêt explicative pourrait être $2\theta - \theta^2$, qui représente la probabilité le toucher un plateau lorsqu'on a droit à deux essais par plateau. ♣

Exemple (MT). Dans le modèle du pédagogue, une quantité particulièrement intéressante est μ , qui décrit l'efficacité globale de la nouvelle méthode pédagogique : c'est bien une quantité d'intérêt, puisque c'est l'image du paramètre caché par l'application « première composante ». ♣

Exemple (MU). Dans le modèle du pédagogue, une quantité d'intérêt *prédictive* intéressante pourrait être la pire score parmi les élèves de la seconde promotion (si on continue de leur appliquer la nouvelle méthode), à savoir, $\min\{Y_j \mid j \in \llbracket 0, m \rrbracket\}$. ♣

Inférence statistique

!! **Définition (MV)** (Inférence statistique). Une inférence statistique est un protocole qui, étant donné un modèle d'inférence statistique pour lequel on a pris connaissance

de la valeur effective de l'observation, vise à dire des choses aussi intéressantes que possibles sur la quantité d'intérêt ; et ce, de façon *mathématiquement rationnelle* et *numériquement quantifiable*. ♡

Définition (MW). On parle d'inférence *explicative* lorsque notre modèle statistique est explicatif et que le but est de “reconstituer” une certaine quantité d'intérêt explicative, resp. d'inférence *prédictive* lorsque notre modèle statistique est prédictif et que le but est de “prédire” une certaine quantité d'intérêt prédictive. ♡ !

Dans la définition (MV), nous avons parlé de « protocole ». Ceci est étroitement lié à la notion de « statistique » :

Point (MX). Dans la définition (MV), parler de « protocole » signifie qu'on a une procédure d'analyse qui est en capacité de fournir une réponse pour toute valeur possible de l'observation ; même si, *en pratique*, on ne donnera cette réponse que pour sa valeur *effective* x_{\checkmark} ... Une inférence statistique peut donc être vue comme consistant à calculer une certaine *fonction* de l'observation X . ♡

Cette idée de « fonction de l'observation » correspond, dans le jargon de la statistique, au concept de « statistique » :

Définition (MY) (Statistique). On appelle *statistique* une variable aléatoire qui dépend uniquement de l'observation. ♡ !!

Principe (MZ). *Le résultat final d'une analyse statistique doit nécessairement être (la réalisation de) une (ou des) statistique(s).* ♡ !!

Remarque (NA). On dit parfois que l'inférence statistique explicative est la “démarche inverse” du calcul des probabilités : dans les deux cas en effet, un paramètre θ influe sur la loi d'une observation aléatoire X ; mais, alors qu'en probabilités, on connaît la vraie valeur θ_{\checkmark} de θ et qu'on cherche à en déduire des choses sur X , en inférence statistique, on connaît la réalisation x_{\checkmark} de X et on cherche à en déduire des choses sur θ !

Dans le cas de l'inférence prédictive, on mélange *deux* niveaux d'incertitude :

- D'une part, l'incertitude qui provient du fait que θ_{\checkmark} est inconnu, qui est de nature “statistique” ;
- Et d'autre part le fait que, même si on connaissait la vraie valeur θ_{\checkmark} du paramètre caché, il resterait un aléa sur Y : cette incertitude-là étant de nature “probabiliste”.

♡

7.4 Paramètres des modèles

Notion de paramètre du modèle

Lorsque j'ai formalisé le modèle du chasseur, j'ai introduit des notations littérales ‘ n ’ et ‘ m ’ pour désigner les nombres de plateaux à tirer resp. lors du test et lors de la compétition. J'ai expliqué que ces nombres étaient connus, et qu'on les prendrait égaux à resp. 25 et 75 en l'absence de mention explicite du contraire. (Et de même, pour le modèle du pédagogue, on a introduit des notations ‘ n ’ et ‘ m ’ similaires). Cependant, on pourrait choisir n'importe quelles valeurs entières naturelles pour n et m sans que la définition ne cesse d'être une description correcte de modèle

statistique... Du coup, si nous n'avions pas spécifié les valeurs de n et m dans notre définition, nous aurions en fait défini une *famille* de modèles statistiques, indexée par $\mathbb{N} \times \mathbb{N}$.

Du point de vue de la nomenclature, on a envie de considérer qu'une telle famille de modèles constitue en fait « un » modèle “flou” où les notations ‘ n ’ et ‘ m ’ sont autorisées à prendre différentes valeurs. Dans ce cas, ‘ n ’ et ‘ m ’ seront qualifiées de « paramètres du modèle » :

! **Définition (NB)** (Paramètres du modèle). Un *paramètre du modèle*, dans un modèle statistique, est une quantité nommée dans la description du modèle, qu'on autorise à prendre différentes valeurs (dans un espace qu'on précisera), mais dont la valeur sera toujours connue à l'avance pour les cas qu'on sera amené à étudier effectivement. Autrement dit, un « modèle présentant des paramètres du modèle » est en fait une *famille* de modèles au sens de la définition (LV) (ou de la définition (MF)), paramétrée par les noms des paramètres du modèle. ♡

Remarque (NC). Lorsque le modèle étudié comporte un ou plusieurs paramètres de modèle, il convient que l'analyse statistique soit menée d'une façon pouvant s'appliquer à toutes les valeurs possibles de ce(s) paramètre(s); autrement dit, les statistiques constituant nos réponses devront faire appel aux paramètres du modèle sous leur forme *littérale*.

Dans une telle situation, quand bien même notre « modèle avec paramètre(s) » est en fait une famille de modèle au sens des définitions (LV) ou (MF), on continuera à utiliser le singulier pour dire qu'on fait « une » inférence sur « ce » modèle. ♣

Remarque (ND). Dans le cadre de nos exercices, nous conviendrons de l'usage suivant : on considère qu'une grandeur constitue un paramètre du modèle dès lors qu'elle reçoit une appellation littérale (et qui aurait été susceptible de prendre d'autres valeurs possibles sans nuire à la cohérence de la description) ou que la description du modèle prévoit explicitement qu'il puisse y avoir d'autres valeurs pour ce paramètre. ♣

Un cas particulier de modèles avec paramètre, qu'il est utile de connaître, est celui des *modèles d'échantillonnage* :

! **Définition (NE)** (Modèle d'échantillonnage). Nous nous plaçons ici dans le cadre des modèles explicatifs. Lorsque le modèle considéré contient un paramètre du modèle n tel que (les autres paramètres éventuels du modèle et) le paramètre caché étant fixé, l'observation est constituée par un n -uplet (X_0, \dots, X_{n-1}) de v.a.i.i.d. dont la loi ne dépend pas de n , alors on qualifie notre modèle le *modèle d'échantillonnage*, pour lequel le paramètre n constitue la *taille d'échantillon*. ♡

Exemple (NF). Le modèle du pédagogue (dans sa version explicative) est un modèle d'échantillonnage, le paramètre de taille d'échantillon correspondant à n . ♣

Exemple (NG). En revanche, le modèle du chasseur n'est pas, techniquement parlant, un modèle d'échantillonnage par rapport à n , puisque l'observation X (le nombre de tirs réussis) ne se décompose pas, dans le formalisme que nous avons choisi, en sous-observations indépendantes. ♣

Notion de régime asymptotique

Outre le fait qu'ils permettent de traiter d'un seul coup plusieurs variantes, les paramètres des modèles permettent d'introduire des notions de *régimes asymptotiques*.

tiques, par rapport auxquels on pourra obtenir des résultats valables *à la limite*. Cela se formalise ainsi :

Définition (NH) (Régime asymptotique). Dans un modèle statistique présentant des paramètres, on peut définir un *régime asymptotique* en considérant un certain passage à la limite d'un (ou plusieurs) paramètres du modèle. ♡

Exemple (NI). Dans le cadre du modèle du pédagogue (en version prédictive), on peut considérer l'asymptotique suivante : « la promotion actuelle et la seconde promotion ont la même taille, et cette taille tend vers l'infini ». Il est alors possible de démontrer, par exemple, qu'à partir de l'observation (passée), on peut construire une approximation de la quantité d'intérêt « pire score obtenu par la seconde promotion » dont l'erreur tendra vers zéro dans l'asymptotique considérée. ♣

7.5 Statistiques bayésienne et fréquentiste

J'ai écrit plus haut que, lorsqu'on définit un modèle de statistique (mettons ici, pour alléger les notations, qu'il s'agit d'un modèle explicatif), on introduit automatiquement un espace probabilisé où on a des variables aléatoires θ et X , où les lois conditionnelles $\text{Loi}(X \mid \theta = \theta)$ sont imposées par le modèle considéré.

Néanmoins, la donnée des différentes $\text{Loi}(X \mid \theta = \theta)$ ne permet pas de connaître la loi de X “tout court”, et encore moins celle de θ ! La description que nous avons donnée de nos variables aléatoires θ et X est donc “incomplète”.

Pour rendre cette description complète, il faudrait que nous précisions aussi la loi de θ . Lorsqu'on le fait, on se retrouve dans ce qu'on appelle un modèle de *statistique bayésienne*. À l'inverse, si on se refuse à le faire, on dit qu'on fait de la *statistique fréquentiste* :

Définition (NJ) (Statistiques bayésienne et fréquentiste).

- Un modèle de statistique *bayésienne* est un modèle statistique (explicatif ou prédictif) dans lequel, outre le modèle lui-même, on donne aussi une loi pour le paramètre caché. Cette loi est qualifiée de *loi à priori* du paramètre caché ; on l'appelle aussi plus simplement la *prior* [¶].
- On dit qu'on fait de la statistique *fréquentiste* lorsqu'on se refuse à introduire une prior. Dans ce cas, on ne devra manipuler θ et X , et faire nos analyses, que d'une façon qui soit indépendante de la loi de θ , et ne dépende donc que du modèle statistique à proprement parler !

♡

Principe (NK) (Interprétation de la prior). *La loi à priori de θ représente, dans le cadre de notre modèle, à quel point nous nous attendions à ce que θ_{\checkmark} vaille tant ou tant avant de commencer à regarder l'observation. Certes, au moment où on commence à collecter les données, la valeur θ_{\checkmark} est déjà fixée définitivement ; mais comme nous n'y avons pas accès, nous la traitons comme s'il s'agissait d'une quantité aléatoire !*

Une autre façon de formuler les choses est la suivante : la prior exprime la connaissance que nous avons, avant l'expérience, du paramètre caché, et des valeurs que, à nos yeux, il est plus ou moins susceptible de prendre. ◇

[¶]. En anglais : *prior*.

Exemple (NL) (Prior par défaut pour le modèle du chasseur). On peut traiter le modèle du chasseur dans un cadre bayésien en utilisant la prior suivante : de par l'expérience qu'il a déjà eue avec ses candidats précédents, sans rien savoir sur le candidat actuel, le club considère qu'il est plus probable que la fiabilité du candidat soit proche de 0 ou de 1, plutôt qu'elle ne soit proche de 1/2 : plus précisément, il considère que le paramètre caché θ est distribué à priori selon la loi suivante, dite *loi de l'arcsinus* :

$$\mathbb{P}_{\text{à priori}}(\theta \in d\theta) = \frac{1}{\pi} \theta^{-1/2} (1 - \theta)^{-1/2} \text{vol}_1(d\theta). \quad (\text{NM})$$

Dans la suite de ce polycopié, lorsque le modèle du chasseur sera traité dans le paradigme bayésien, la prior considérée sera celle de l'arcsinus, à moins qu'il ne soit explicitement précisé autrement. ♣

Exemple (NN) (Prior par défaut pour le modèle du pédagogue). Quand nous traiterons le modèle du pédagogue dans le cadre bayésien, la prior que nous utiliserons sera la suivante (sauf mention explicite du contraire) : μ et σ seront prises indépendantes à priori, avec μ suivant (à priori) la loi Normale($\mu_{\text{réf}} - \sigma_{\text{réf}}, \sigma_{\text{réf}}^2$) et σ suivant la loi Expon^{le}($\sigma_{\text{réf}}^{-2}/2$)^{-1/2} (càd. la mesure image, par l'application $x \mapsto x^{-1/2}$, de la loi Expon^{le}($\sigma_{\text{réf}}^{-2}$)). On peut en l'occurrence calculer explicitement la densité correspondante : pour $\mu \in \mathbb{R}, \sigma \in \mathbb{R}_+^*$, on a ainsi

$$\mathbb{P}_{\text{pr}}(\mu \in d\mu \text{ et } \sigma \in d\sigma) = \frac{\sigma_{\text{réf}}}{\sqrt{2\pi}} \frac{1}{\sigma^3} \exp\left(-\frac{(\mu - \mu_{\text{réf}} + \sigma_{\text{réf}})^2}{2\sigma_{\text{réf}}^2} - \frac{\sigma_{\text{réf}}^2}{2\sigma^2}\right) \text{vol}_1(d\mu) \text{vol}_1(d\sigma). \quad (\text{NO})$$

Une visualisation graphique de cette prior est donnée en figure 7.1. ♣

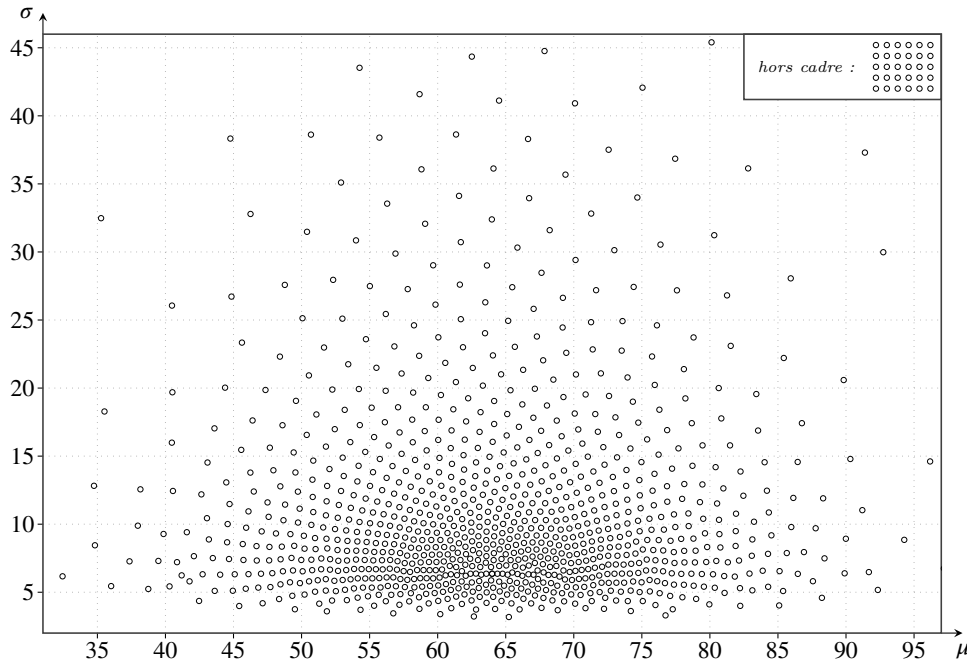


FIGURE 7.1 – Notre prior par défaut pour le modèle du pédagogue. Pour représenter cette loi sur \mathbb{R}^2 , on a fait figurer un petit cercle par zone de probabilité 1 ‰ : ainsi, plus les cercles sont rapprochés les uns des autres, plus la densité est élevée.

Remarque (NP). Vu que la priore représente à quel point nous nous attendons à telle ou telle valeur pour le paramètre caché, elle comporte en général une part de *subjectivité* : nous verrons d'ailleurs dans la § 15 en quoi la détermination de la priore s'apparente parfois à un “art” plutôt qu'à l'application rigoureuse d'une procédure algorithmique. ♣

Remarque (NQ). Nous verrons dans le chapitre 15 comment l'idée générale de « exprimer notre connaissance du paramètre caché avant l'expérience » peut se décliner en pratique pour en arriver effectivement au choix d'une loi à priori dans une situation donnée. ♣

Dans le cas de la statistique bayésienne, deux notations très importantes interviennent : celles concernant les contextes probabilistes « à priori » et « à postériori » :

Définition (NR). Dans le cadre bayésien, on parle de « contexte probabiliste à priori » pour signifier que, non seulement on traite le paramètre caché comme une “authentique” variable aléatoire (i.e., on ne conditionne pas par une valeur donnée de θ), mais surtout, qu'on n'a pas conditionné par l'information sur la valeur de l'observation. Le contexte probabiliste à priori sera noté « $\mathbb{P}_{\text{pr}}(\bullet)$ » : bien que cette notation soit *techniquement* synonyme de $\mathbb{P}(\bullet)$, il est fortement préférable de souligner explicitement cette absence de conditionnement !

Par contraste, on parle de « contexte probabiliste à postériori » lorsque, toujours dans le cas où θ est une authentique variable aléatoire, on raisonne conditionnellement à la valeur effectivement prise par l'observation. Le contexte à postériori sera noté « $\mathbb{P}_{\text{post}}(\bullet)$ » : cette notation est techniquement synonyme de $\mathbb{P}(\bullet \mid X = x_{\checkmark})$. ♡

7.6 Description d'un modèle statistique en langage ordinaire

Point (NS). Dans les exercices, on décrira souvent les modèles de statistique à l'aide de phrases du langage ordinaire (i.e. non mathématique). Typiquement les énoncés prendront plus ou moins la forme suivante (où les passages entre crochets en italique désignent la forme générale de ce qui s'insèrera dans l'énoncé, et où les passages entre parenthèses pourront éventuellement être absents) :

[Tel objet du monde réel ou abstrait] (que nous noterons [notation du type 'X' ou du type ('X', 'Y')]) suit la loi [description d'une loi pour l'objet dépendant de certains paramètres, les paramètres en question étant du type ' θ_{\checkmark} ' (et ' λ ')], où ([un ou plusieurs des paramètres, ceux du type λ] vaut [telle valeur connue pour λ] et [un ou plusieurs des paramètres, ceux du type θ_{\checkmark}] est inconnu, à valeurs dans [un certain ensemble pour θ_{\checkmark}]. (On suppose à priori que [θ_{\checkmark}] est la réalisation d'une variable aléatoire (notée [une notation de type θ]) suivant la loi [une certaine loi pour θ]).

Dans la situation à laquelle nous sommes confrontés en pratique, on a trouvé que [(la partie en X de) l'objet] a pris la valeur [telle valeur pour X] (que nous noterons [notation du type ' x_{\checkmark} ']).

Notre but est de mieux savoir ce que vaut [quantité de la forme $\gamma(\theta_{\checkmark})$ ou de la forme $g(Y)$], en en déterminant [non d'une méthode statistique, p. ex. « un estimateur » ou « un intervalle de prédiction »].

Un tel texte s'interprète de la façon suivante : [tel objet du monde réel ou abstrait] se réfère à l'observation complétée (s'il y a lieu de compléter, s'entend). [La partie X de l'objet] est l'observation. Lorsqu'il n'y a pas de [partie Y de l'objet], c'est qu'on est dans un modèle d'inférence explicative ; et lorsque l'observation complétée est du type (X, Y) [la partie Y de l'objet] constitue l'observation future, et on est alors dans un modèle de prédiction statistique. [Les paramètres du type $\theta_{\mathcal{J}}$] sont la réalisation du paramètre caché, [les paramètres du type λ] étant les paramètres du modèle. L'espace du paramètre caché est [un certain ensemble pour $\theta_{\mathcal{J}}$] Un point délicat est l'interprétation de la description de [une loi pour l'objet] : en effet, l'énoncé décrit une *unique* loi dans laquelle intervient un $\theta_{\mathcal{J}}$ *inconnu* ; mais du point de vue du modèle statistique, il faut comprendre qu'il s'agit implicitement de la description d'une *famille* de lois indexée par [un certain ensemble pour $\theta_{\mathcal{J}}$], une pour chaque valeur θ dans cet ensemble. Lorsque la phrase commençant par « on suppose à priori que » est présente, cela signifie qu'on introduit un contexte bayésien, avec pour priore [une certaine loi pour θ]. [Telle valeur pour X] est la valeur effective de l'observation (passée). [La quantité de la forme $\gamma(\theta_{\mathcal{J}})$ ou $g(Y)$] désigne la quantité d'intérêt qui fait l'objet de notre étude : si celle-ci est de la forme $\gamma(\theta_{\mathcal{J}})$, c'est qu'on s'intéresse à une question d'inférence ; si elle est de la forme $g(Y)$, c'est qu'on s'intéresse à une question de prédiction. ♣

Remarque (NT). Dans cette formulation, l'énoncé donne seulement la loi de l'observation (complétée) sous le *véritable* contexte probabiliste, autrement dit, sachant la vraie valeur $\theta_{\mathcal{J}}$ du paramètre caché ; mais il le fait sans rien préciser sur $\theta_{\mathcal{J}}$: c'est pourquoi on peut considérer que cette description fournit en fait une *famille* de lois indexée par Θ . Cette simplification rédactionnelle sera régulièrement utilisée dans le polycopié ainsi que dans les exercices, parfois même implicitement. Elle justifie ainsi l'introduction de deux nouvelles notations : ♣

Notation (NU). Pour $\theta \in \Theta$ une valeur possible du paramètre caché, le contexte probabiliste $\mathbb{P}(\bullet \mid \theta = \theta)$ — autrement dit, le contexte sachant que le paramètre caché vaut θ — pourra également être noté « $\mathbb{P}_{\theta}(\bullet)$ ».

Parmi ces contextes, on appelle « véritable contexte probabiliste », et on note « $\mathbb{P}_{\mathcal{J}}(\bullet)$ », le contexte probabiliste correspondant à la vraie valeur du paramètre caché : autrement dit, $\mathbb{P}_{\mathcal{J}}(\bullet) := \mathbb{P}_{\theta_{\mathcal{J}}}$. ♡

Revenons maintenant au point (NS), en illustrant celui-ci par un exemple concret.

Exemple (NV). Voici un avatar du modèle du pédagogue qui pourrait constituer le début d'un énoncé d'exercice :

Dans le cadre d'un nouveau cours, une enseignante s'attend à ce que les notes de ses élèves soient i.i.d. suivant une loi Normale $(\mu_{\mathcal{J}}, \sigma_{\mathcal{J}}^2)$, pour des valeurs de $\mu_{\mathcal{J}}$ et $\sigma_{\mathcal{J}}$ qui sont encore complètement inconnues pour lui (puisque le cours est nouveau). L'enseignante donne le cours deux années de suite, d'abord à une promotion de n élèves, puis à une promotion de p élèves. Au moment que nous considérons, l'enseignante a gardé en mémoire les notes de la première promotion, mais la correction de la seconde promotion n'ayant pas encore eu lieu, les notes de celle-ci sont encore inconnues. Notre enseignante se demande, avant de commencer à corriger les copies, à quelles valeurs elle peut raisonnablement s'attendre pour la moyenne de la seconde promotion. Pour les

applications numériques, nous prendrons $n = 16$, $p = 21$, et les notes de la première promotion (prises dans un ordre arbitraire, p. ex. l'ordre alphabétique des noms) sont : 13,2, 9,1, 1,6, 8,6, 10,7, 12,9, 11,9, 8,9, 15,5, 13,1, 13,5, 11,6, 11,7, 11,4, 10,0, 10,9.

Dans ce cas, l'interprétation en termes de modèle est la suivante :

Observation complétée : Notes de tous les élèves des deux promotions. Cette observation complétée se divise en une observation passée, correspondant aux notes de la première promotion, et une observation future, correspondant aux notes de la seconde promotion. Si on est amené à introduire soi-même des notations pour les observations, un choix logique sera de noter l'observation (passée) comme le n -uplet (X_0, \dots, X_{n-1}) , les X_i correspondant aux notes respectives des élèves, et l'observation future comme le p -uplet (Y_0, \dots, Y_{p-1}) . Concernant l'espace de l'observation, les notes doivent être considérées comme réelles puisqu'elles suivent des lois normales : l'espace de l'observation passée est donc \mathbb{R}^n , et celui de l'observation future est \mathbb{R}^p . (Comme nous le verrons, n et p sont ici des paramètres du modèle, de sorte qu'ils ont parfaitement le droit d'intervenir dans la description des espaces des observations!).

Observation effective : Les 16 notes données en fin d'énoncé correspondent à la réalisation de l'observation (passée) : $x_{0\checkmark} := 13,2$, $x_{1\checkmark} := 9,1$, etc.

Paramètre caché : (La réalisation de) le paramètre caché est constitué du couple $(\mu_{\checkmark}, \sigma_{\checkmark})$, puisque ce sont les quantités inconnues qui déterminent la loi de l'observation (complétée). Les variables aléatoires associées à ces réalisations pourront être notées resp. μ et σ . L'espace du paramètre caché n'ayant pas été précisé ici, nous allons le prendre aussi large que mathématiquement possible : μ sera donc à valeurs dans \mathbb{R} et σ à valeurs dans \mathbb{R}_+ .

Loi de l'observation (complétée) sachant le paramètre caché : C'est le point le plus important dans la description d'un modèle! L'affirmation que les notes des différents élèves sont i.i.d. de loi Normale($\mu_{\checkmark}, \sigma_{\checkmark}^2$) correspond à la description de la loi de l'observation complétée sachant la vraie valeur du paramètre caché. En termes plus formels, cela peut se traduire par le fait que $\text{Loi}_{\checkmark}(X, Y) = \text{Normale}(\mu_{\checkmark}, \sigma_{\checkmark}^2)^{\otimes(n+p)}$. Comme on a écrit cela sans aucune connaissance de la valeur $(\mu_{\checkmark}, \sigma_{\checkmark})$, cela signifie en fait que, pour tout $(\mu, \sigma) \in \mathbb{R} \times \mathbb{R}_+^*$, on a

$$\text{Loi}(X, Y \mid \mu = \mu \text{ et } \sigma = \sigma) = \text{Normale}(\mu, \sigma^2)^{\otimes(n+p)}. \quad (\text{NW})$$

Paramètres du modèle : Deux autres paramètres interviennent dans la description du modèle, qui sont n et p : ce ne sont pas des paramètres cachés dans la mesure où la taille des promotions est parfaitement connue pour notre analyse statistique (en l'occurrence, l'énoncé nous précise même que $n := 16$ et $p := 21$) ; ce sont donc des paramètres du modèle!

Priore : Aucune priore n'étant mentionnée pour le modèle, nous serons forcément dans un cadre fréquentiste.

Quantité d'intérêt : L'énoncé dit explicitement que la quantité d'intérêt constituant notre objectif sera la moyenne des Y_j , soit $(\sum_{j=1}^p Y_j) / p$: c'est bien une fonction de l'observation future, comme il fallait s'y attendre pour un modèle de prédiction.

♣

7.7 Récapitulatif des principales notations

Variables aléatoires, valeurs possibles et réalisations

Notation (NX). Dans tout ce polycopié, je suivrai autant que possible l'usage de désigner les *variables aléatoires* par des lettres majuscules (X, Y, M, S, \dots), et les *valeurs* de ces variables aléatoires par les lettres minuscules correspondantes (x, y, m, s, \dots).

Pour les variables aléatoires désignées par des lettres *grecques*, le contraste « majuscule / minuscule » sera remplacé par le contraste « gras / maigre » : les variables aléatoires seront donc désignées par $\theta, \alpha, \beta, \lambda, \dots$, et leurs valeurs par $\theta, \alpha, \beta, \lambda, \dots$. Si vous trouvez pénible de rendre le gras en écriture manuscrite, vous pourrez émuler le gras par du souligné ($\underline{\theta}, \underline{\alpha}, \underline{\beta}, \underline{\lambda}, \dots$) ou par du double-barres ($\mathbb{\theta}, \mathbb{\alpha}, \mathbb{\beta}, \mathbb{\lambda}, \dots$). ♡

Notation (NY). Pour distinguer la valeur *effectivement prise* par une variable aléatoire (autrement dit, sa réalisation), on utilisera le symbole '✓' (« coche ») en indice : x_{\checkmark} sera ainsi la valeur effectivement prise par la variable aléatoire X . Pour désigner d'autres valeurs possibles remarquables pour X , on utilisera des notations comme $x_0, x_1, x_{\text{réf}}, x_*$, ... On utilisera enfin le simple 'x' pour désigner une valeur possible complètement générique, ou pour une variable muette. ♡

Remarque (NZ). Dans le cadre de ce cours, afin de gagner du temps, l'association entre lettres majuscules et majuscules (ou grasses et maigres), ainsi que l'usage de la coche, sera souvent implicite : si une variable aléatoire s'appelle W , on pourra manipuler directement la notation ' w_{\checkmark} ' sans préciser au préalable que cela désigne la valeur effectivement prise par W ...! ♣

Remarque (OA). Dans vos exercices et examens, une tolérance vous sera accordée concernant l'omission du symbole '✓' (sauf si celui-ci figure explicitement dans l'énoncé) ou du passage en gras, à moins que cela ne crée un souci d'intelligibilité [||]. ♣

Il existe une convention concernant l'usage des lettres latines et grecques en statistique :

Convention (OB).

- Dans le cadre de ce polycopié, on utilisera des lettres grecques pour désigner le paramètre caché, ainsi que pour les quantités d'intérêt explicatives et pour leurs estimateurs (sur ce dernier point, confer chapitre 11).
- À l'inverse, on utilisera des lettres latines pour désigner l'observation et l'observation future, ainsi que pour les statistiques de test (sur ce point, confer chapitre 12) et les quantités d'intérêt prédictives.
- Les paramètres du modèle, quant à eux, pourront aussi bien être désignés par des lettres latines que par des lettres grecques. ♡

Remarque (OC). Attention : rien ne garantit en revanche que ces conventions soient respectées dans vos exercices ! En effet, dans une situation industrielle où vous aurez à proposer vous-mêmes un modèle statistique, il va de soi que personne n'aura pris soin de respecter les conventions en amont à votre place... ♣

[||]. Par exemple, il faudra impérativement respecter le gras lorsqu'on écrira « $\mathbb{P}_{\text{post}}(\theta = \theta)$ » : sinon, on obtiendrait « $\mathbb{P}_{\text{post}}(\theta = \theta)$ », où il y aurait la même chose des deux côtés du signe '='...!

Notations raccourcies

Dans les pages précédentes, nous avons introduit différentes notations condensées pour certains contextes probabilistes que nous serons amenés à rencontrer particulièrement souvent en statistique. Récapitulons-les ci-dessous :

Notation (OD).

- Pour θ une valeur possible pour le paramètre caché, le contexte probabiliste sachant que le paramètre caché vaut θ (autrement dit, le contexte $\mathbb{P}(\bullet \mid \theta = \theta)$) pourra être noté “ $\mathbb{P}_\theta(\bullet)$ ”.
- Le véritable contexte probabiliste, celui connaissant la vraie valeur du paramètre caché (autrement dit, le contexte $\mathbb{P}_{\theta_\vee}(\bullet)$), pourra être noté “ $\mathbb{P}_\vee(\bullet)$ ”.
- En statistique bayésienne, lorsqu'on ne conditionne pas par rapport à la valeur de l'observation (ni par quoi que ce soit d'autre), on note le contexte probabiliste correspondant “ $\mathbb{P}_{\text{pr}}(\bullet)$ ”, histoire de bien clarifier cette absence de conditionnement. Techniquement parlant, c'est la même chose que le contexte $\mathbb{P}(\bullet)$; mais on préfère utiliser des notations qui rendent toujours explicite le contexte considéré.
- Toujours en statistique bayésienne, le contexte probabiliste à postériori (autrement dit, le contexte $\mathbb{P}(\bullet \mid X = x_\vee)$, où ‘ X ’ se réfère à l'observation et ‘ x_\vee ’ à sa réalisation effective) pourra être noté “ $\mathbb{P}_{\text{post}}(\bullet)$ ”.

♡

Un autre type de raccourci fréquemment utilisé sera le suivant :

Notation (OE). Dans l'expression d'une probabilité de la forme « $\mathbb{P}(X = x)$ », on pourra sous-entendre le “ $X =$ ” et écrire simplement « $\mathbb{P}(x)$ ». De même, on pourra simplifier « $\mathbb{P}(X \in dx)$ » en « $\mathbb{P}(dx)$ ». Similairement, à droite d'une barre de conditionnement, une expression comme « $\text{Loi}(Y \mid X = x)$ » pourra être raccourcie en « $\text{Loi}(Y \mid x)$ ».

♡

Notations génériques

Dans le cadre de ce cours, lorsque nous voudrions introduire de nouvelles méthodes d'analyse statistique, il serait un peu pénible de devoir re-préciser à chaque fois toutes les notations permettant de se référer aux différents aspects d'un modèle statistique... Pour faciliter les choses, nous pourrions donc utiliser un certain nombre de notations par défaut, que j'appellerai « notations génériques » dans la suite de ce polycopié :

Définition (OF) (Notations génériques). La phrase « *Soit un modèle statistique explicatif avec les notations génériques* » signifie : « Soit un modèle statistique dans lequel :

- L'espace du paramètre caché est noté ‘ Θ ’;
- Le paramètre caché est noté ‘ θ ’;
- L'espace de l'observation est noté ‘ \mathcal{X} ’;
- L'observation est notée ‘ X ’;
- La notation ‘ γ ’ se réfère à une quantité d'intérêt explicative, qui peut aussi s'écrire “ $\gamma(\theta)$ ”;
- L'espace dans lequel la quantité d'intérêt γ prend ses valeurs est noté ‘ \mathcal{G} ’.

La phrase « *Soit un modèle statistique prédictif avec les notations génériques* » signifie la même chose que ci-dessus, aux nuances suivantes près :

- Cette fois-ci, il y a également une observation future, notée ' Y ' ;
- L'espace dans lequel Y prend ses valeurs est noté ' \mathcal{Y} ' ;
- La notation ' γ ' n'est plus censée se référer à quelque chose en particulier, pas plus que la notation " $\gamma(\bullet)$ " ;
- En revanche, la notation ' G ' se réfère maintenant à une quantité d'intérêt prédictive, qui peut aussi s'écrire " $g(\theta)$ " ;
- La notation ' \mathcal{G} ' correspond à présent à l'espace dans lequel la quantité d'intérêt G prend ses valeurs.

En outre, en cas d'utilisation des notations génériques, si nous souhaitons considérer une certaine asymptotique pour un paramètre du modèle, nous appellerons ce paramètre λ et désignerons l'asymptotique considérée par « $\lambda \rightarrow \infty$ ». \heartsuit

Chapitre 8

Le théorème de Bayes

8.1 Notion de postérieure

Contextes à postérieur et à priori

(Ce qui suit est écrit en notations génériques). En statistique bayésienne, le fait de disposer d'une loi de probabilité pour θ (la « prior ») permet d'en déduire, via les lois conditionnelles $\text{Loi}(X, Y \mid \theta = \theta)$, la loi jointe du triplet (θ, X, Y) . Par conséquent, on peut aussi de mettre à parler de la loi de θ et de Y *conditionnellement à telle ou telle valeur de X* , et donc en particulier, de leur loi sachant la valeur qui a été *effectivement observée* pour X ! On parle alors de *probabilités à postérieur* :

Définition (OG) (À postérieur). En statistique bayésienne, le *contexte probabiliste à postérieur* désigne le contexte conditionné par rapport à $\{X = x_{\checkmark}\}$ (en notations génériques). Ainsi, pour A un évènement, la probabilité à postérieur de A désigne la probabilité $\mathbb{P}(A \mid X = x_{\checkmark})$; et on alors en déduire les notions de loi à postérieur, espérance à postérieur, variance à postérieur, etc. !!

Dans le cadre de ce cours, rappelons que, plutôt que d'écrire explicitement « $\mathbb{P}(\bullet \mid X = x_{\checkmark})$ », on pourra désigner une probabilité à postérieur par la notation $\mathbb{P}_{\text{post}}(\bullet)$ (ou $\mathbb{P}_{\text{po}}(\bullet)$), et de même une loi à postérieur par $\text{Loi}_{\text{post}}(\bullet)$, etc. ♡

Principe (OH). *Un modèle de statistique bayésienne donnant toute l'information suffisante au calcul de la loi jointe de (θ, X, Y) , on peut toujours y calculer les probabilités (resp. lois, $\mathcal{E}c.$) à postérieur ! À l'inverse, en statistique fréquentiste, la notion de probabilité à postérieur, et tout ce qui s'y rapporte, n'a pas de sens, car le modèle ne donne pas suffisamment d'information pour qu'on puisse les définir... !* !!

Définition (OI). La loi à postérieur du paramètre caché est appelée simplement « la postérieure ». ♡

À l'inverse, dans une situation bayésienne, lorsqu'on souhaite insister sur le fait qu'on n'a pas conditionné par la valeur de l'observation, on parle de contexte probabiliste à priori :

Définition (OJ). On utilise la locution « à priori » pour souligner qu'on se place dans un contexte probabiliste où on n'a pas conditionné par la valeur de l'observation. Le contexte probabiliste à priori sera noté $\mathbb{P}_{\text{pr}}(\bullet)$. !

Bien que “ $\mathbb{P}_{\text{pr}}(\bullet)$ ” soit formellement synonyme de “ $\mathbb{P}(\bullet)$ ”, on s’efforcera de préciser systématiquement l’aspect « à priori » dès lors que l’objet dont on parle dépend du choix de la priore. \heartsuit

De la loi du paramètre caché à celle des quantités d’intérêt

Dans quelques pages, la section 8.2 nous dira comment, dans un modèle de statistique bayésienne, on peut déterminer la postérieure, autrement dit la loi à postériori du paramètre caché θ . Dans la présente sous-section, nous allons voir l’argument qui justifie toute l’importance de la postérieure dans nos analyses : à savoir que, à partir de la postérieure, on peut déterminer les lois à postériori de nos quantités d’intérêt !

Ci-dessous, j’ai indiqué divers points comme étant à retenir : en l’occurrence, c’est la *méthode de raisonnement* qui doit être retenue, pas le résultat lui-même !

! *Point (OK)* (Loi à postériori d’une quantité d’intérêt explicative). Commençons par le cas d’une quantité d’intérêt explicative $\gamma(\theta) =: \boldsymbol{\gamma}$, où $\gamma(\bullet)$ est une certaine fonction déterministe de Θ dans un certain espace \mathcal{E} . Il est clair que, si nous connaissons la loi à postériori de θ , nous pouvons en déduire la loi de $\boldsymbol{\gamma}$ dans ce même contexte : en effet, pour $A \subseteq \mathcal{E}$, on a que

$$\mathbb{P}_{\text{post}}(\boldsymbol{\gamma} \in A) = \mathbb{P}_{\text{post}}(\gamma(\theta) \in A) = \mathbb{P}_{\text{post}}(\theta \in \gamma^{-1}(A)) = \mathbb{P}(\text{Loi}_{\text{post}}(\theta) \in \gamma^{-1}(A)). \quad (\text{OL})$$

Par conséquent, si nous connaissons la postérieure, nous sommes en mesure de déterminer la loi à postériori de la quantité d’intérêt $\boldsymbol{\varphi}$, et donc par ricochet ses espérance à postériori, quantiles à postériori, etc. \clubsuit

Remarque (OM). En fait, le raisonnement ci-dessus revient à utiliser que la loi à postériori de $\gamma(\theta)$ est la mesure-image, par l’application $\gamma(\bullet)$, de la loi à postériori de θ . \clubsuit

Remarque (ON). En fait, le raisonnement du point (OK) n’a rien de spécifique au contexte probabiliste à postériori : on pourrait donc, exactement de la même manière, déterminer la loi à priori de $\boldsymbol{\gamma}$ à partir de la priore ! Le raisonnement n’a rien non plus de spécifique au fait qu’on s’intéresse à une fonction de θ : on pourrait similairement, sous le contexte probabiliste sachant la valeur du paramètre caché, déterminer la loi $\text{Loi}_{\theta}(t(X))$ d’une statistique ou celle $\text{Loi}_{\theta}(g(Y))$ d’une quantité d’intérêt à partir de resp. $\text{Loi}_{\theta}(X)$ et $\text{Loi}_{\theta}(Y)$! \clubsuit

! *Point (OO)* (Loi à priori de l’observation). À présent, regardons comment déterminer la loi à priori de l’observation. (On utilise ci-après les notations génériques). « Déterminer la loi de X », cela revient à calculer la probabilité d’avoir $\{X \in dx\}$ pour dx un voisinage infinitésimal d’un point arbitraire $x \in \mathcal{X}$: nous considérons donc dans la suite un tel voisinage infinitésimal.

Le modèle statistique nous donnant la loi de l’observation *sachant telle ou telle valeur du paramètre caché*, nous allons utiliser la formule des probabilités totales pour faire apparaître les différentes valeurs possibles du paramètre caché :

$$\mathbb{P}_{\text{pr}}(X \in dx) = \int_{\theta \in \Theta} \mathbb{P}_{\text{pr}}(\theta \in d\theta \text{ et } X \in dx). \quad (\text{OP})$$

De là, on fait intervenir la loi sachant la valeur de θ en utilisant la règle de chaîne :

$$\mathbb{P}_{\text{pr}}(\theta \in d\theta \text{ et } X \in dx) = \mathbb{P}_{\text{pr}}(\theta \in d\theta) \mathbb{P}_{\text{pr}}(X \in dx \mid \theta \in d\theta). \quad (\text{OQ})$$

Le membre de droite ci-dessus se simplifie à trois niveaux. Premièrement, conditionner par un évènement de précision infinitésimale comme $\{\theta \in d\theta\}$, cela signifie en fait la même chose que conditionner par $\{\theta = \theta\}$. Deuxièmement, quand on considère le contexte $\mathbb{P}_{\text{pr}}(\bullet \mid \theta = \theta)$ (qui, formellement, est synonyme de $\mathbb{P}(\bullet \mid \theta = \theta)$), il n'y a en fait pas lieu de préciser « à priori » : en effet, ce contexte-ci ayant été fixé par le modèle indépendamment du choix de la priore, il s'agit d'un concept fréquentiste ! Troisièmement, on se rappelle que le contexte $\mathbb{P}(\bullet \mid \theta = \theta)$ peut aussi être noté $\mathbb{P}_{\theta}(\bullet)$: au final, « $\mathbb{P}_{\text{pr}}(X \in dx \mid \theta \in d\theta)$ » peut donc se ré-écrire simplement « $\mathbb{P}_{\theta}(X \in dx)$ » !

En mettant bout à bout tout ce qui précède, on a finalement obtenu la loi à priori de X :

$$\mathbb{P}_{\text{pr}}(X \in dx) = \int_{\theta \in \Theta} \mathbb{P}_{\theta}(X \in dx) \mathbb{P}_{\text{pr}}(\theta \in d\theta). \quad (\text{OR})$$

♣

Remarque (OS). Le même raisonnement, *mutatis mutandis*, permettrait également de déterminer la loi à priori de Y . ♣

Exemple (OT). Appliquons la méthodologie ci-dessus en calculant, dans le modèle du chasseur, la probabilité à priori (en utilisant la priore de l'arcsinus pour θ) que le candidat rate tous ses plateaux : autrement dit, il s'agit de déterminer $\mathbb{P}_{\text{pr}}(X = 0)$.

Par la formule des probabilités totales, on décompose selon la valeur de θ :

$$\mathbb{P}_{\text{pr}}(X = 0) = \int_{\theta \in]0,1[} \mathbb{P}(\theta \in d\theta \text{ et } X = 0). \quad (\text{OU})$$

La règle de chaîne, nous donne alors

$$\mathbb{P}_{\text{pr}}(X = 0) = \int_{\theta \in]0,1[} \mathbb{P}_{\text{pr}}(\theta \in d\theta) \mathbb{P}(X = 0 \mid \theta = \theta). \quad (\text{OV})$$

Dans cette expression, nous avons alors clairement distingué le facteur $\mathbb{P}_{\text{pr}}(\theta \in d\theta)$, qui dépend uniquement de la priore et vaut $\frac{1}{\pi} \theta^{-1/2} (1 - \theta)^{-1/2} \text{vol}_1(d\theta)$, et le facteur fréquentiste $\mathbb{P}(X = 0 \mid \theta = \theta)$, qui peut être déterminé indépendamment de la nature de la priore : dans la mesure où $\text{Loi}_{\theta}(X) = \text{Binom}^{\text{le}}(n, \theta)$, on a $\mathbb{P}_{\theta}(X = 0) = \mathbb{P}(\text{Binom}^{\text{le}}(n, \theta) = 0 \mid =)(1 - \theta)^n$. Finalement, la probabilité à priori que le candidat ne touche aucun plateau vaut

$$\frac{1}{\pi} \int_{\theta \in]0,1[} \theta^{-1/2} (1 - \theta)^{n-1/2} \text{vol}_1(d\theta) = \frac{1}{\pi} \int_{\theta=0}^1 \theta^{-1/2} (1 - \theta)^{n-1/2} d\theta, \quad (\text{OW})$$

dont on calcule numériquement que cela vaut 11,2 %. ♣

Pour déterminer la loi à *posteriori* de l'observation future, la méthode est la même, mais il y a une subtilité supplémentaire à cause du conditionnement par $\{X = x_{\checkmark}\}$:

Point (OX) (Loi à *posteriori* de l'observation future). Avec les notations génériques, pour $y \in \mathcal{Y}$ et dy un voisinage infinitésimal de y , on commence comme précédemment par découper selon la valeur de θ , puis par utiliser la règle de chaîne :

$$\begin{aligned} \mathbb{P}_{\text{post}}(Y \in dy) &= \int_{\theta \in \Theta} \mathbb{P}_{\text{post}}(\theta \in d\theta \text{ et } Y \in dy) \\ &= \int_{\theta \in \Theta} \mathbb{P}_{\text{post}}(Y \in dy \mid \theta \in d\theta) \mathbb{P}_{\text{post}}(\theta = \theta). \quad (\text{OY}) \end{aligned}$$

!

La dernière étape consiste à ré-écrire $\mathbb{P}_{\text{post}}(Y \in dy \mid \theta = \theta)$; mais cette fois-ci c'est un peu plus délicat à comprendre. Nous nous rappelons que le contexte probabiliste " $\mathbb{P}_{\text{post}}(\bullet)$ " signifie en fait " $\mathbb{P}(\bullet \mid X = x_{\checkmark})$ " sauf que là, nous avons un conditionnement supplémentaire par $\{\theta = \theta\}$... Un tel "double conditionnement", en fait, signifie qu'on a conditionné par la *conjonction* des deux événements considérés ! Par conséquent, $\mathbb{P}_{\text{post}}(Y \in dy \mid \theta = \theta)$ se ré-interprète simplement comme $\mathbb{P}(Y \in dy \mid X = x_{\checkmark} \text{ et } \theta = \theta)$, ou encore comme $\mathbb{P}(Y \in dy \mid \theta = \theta \text{ et } X = x_{\checkmark})$. Mais ici, nous pouvons décider d'interpréter le conditionnement par $\{\theta = \theta\}$ comme relevant du *contexte* probabiliste : et finalement, $\mathbb{P}_{\text{post}}(Y \in dy \mid \theta = \theta)$ est la même chose que $\mathbb{P}_{\theta}(Y \in dy \mid X = x_{\checkmark})$.

Une fois arrivé à cette écriture, on comprend que $\mathbb{P}_{\theta}(Y \in dy \mid X = x_{\checkmark})$ est en fait un objet fréquentiste, dans la mesure où la loi jointe $\text{Loi}_{\theta}(X, Y)$ est entièrement déterminée par le modèle lui-même, sans tenir compte de la priore. On a ainsi obtenu une expression de la loi à postériori de Y à partir de la postérieure et de quantités fréquentistes :

$$\mathbb{P}_{\text{post}}(Y \in dy) = \int_{\theta \in \Theta} \mathbb{P}_{\theta}(Y \in dy \mid X = x_{\checkmark}) \mathbb{P}_{\text{post}}(\theta = \theta). \quad (\text{OZ})$$

♣

L'ultime point serait de savoir comment déterminer la loi à priori ou à postériori d'une *quantité d'intérêt prédictive* à partir de resp. la priore et la postérieure. Mais en fait, pour cela, il suffit de combiner les points précédents avec l'idée de mesure-image abordée dans le point (OK) : on détermine d'abord la loi à postériori de Y ; puis on en déduit celle de $g(Y)$ par mesure-image ! 😊

8.2 La formule de Bayes

Dérivation de la formule

☛ *Le but de la présente sous-section est uniquement d'expliquer d'où provient la formule de Bayes que nous allons présenter dans la suite ; le lecteur peu curieux de comprendre le pourquoi des choses peut donc la sauter sans dommage.*

Supposons donné un modèle de statistique bayésienne (en notations génériques) ; nous aimerions calculer la postérieure à partir de la priore et de la description du modèle. « Calculer la postérieure » revient à déterminer $\mathbb{P}_{\text{post}}(\theta \in d\theta)$ pour $d\theta$ un voisinage infinitésimal d'un point arbitraire $\theta \in \Theta$; soit donc un tel $d\theta$.

Par du contexte à postériori, $\mathbb{P}_{\text{post}}(\theta \in d\theta)$ désigne la probabilité conditionnelle $\mathbb{P}(\theta \in d\theta \mid X = x_{\checkmark})$; et par définition du conditionnement, cette probabilité conditionnelle vaut

$$\frac{\mathbb{P}_{\text{pr}}(X \in dx_{\checkmark} \text{ et } \theta \in d\theta)}{\mathbb{P}_{\text{pr}}(X \in dx_{\checkmark})}, \quad (\text{PA})$$

pour dx_{\checkmark} un voisinage infinitésimal de x_{\checkmark} . [L'introduction de ce voisinage servant à éviter d'avoir un quotient dégénéré $0/0$ dans le cas où les $\text{Loi}_{\theta}(X)$ seraient diffuses en x_{\checkmark}].

On applique alors la règle de chaîne au numérateur (en utilisant à nouveau l'équivalence entre « conditionner par $d\theta$ » et « conditionner par θ ») :

$$\mathbb{P}_{\text{pr}}(X \in dx_{\checkmark} \text{ et } \theta \in d\theta) = \mathbb{P}_{\text{pr}}(\theta \in d\theta) \mathbb{P}_{\theta}(X \in dx_{\checkmark}) ; \quad (\text{PB})$$

quant au dénominateur, on va simplement utiliser que c'est une constante (vu qu'on raisonne à $d\mathcal{X}$ fixé). Finalement, on a obtenu que

$$\mathbb{P}_{\text{post}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}) \propto \mathbb{P}_{\boldsymbol{\theta}}(X \in dx_{\mathcal{X}}) \times \mathbb{P}_{\text{pr}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}), \quad (\text{PC})$$

où il n'y a pas besoin de préciser la constante de proportionnalité, celle-ci pouvant être retrouvée par la condition de normalisation d'une loi de probabilité : à savoir, en l'occurrence, que $\int_{\boldsymbol{\theta} \in \Theta} \mathbb{P}_{\text{post}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}) = 1$.

Dans cette formule, le facteur $\mathbb{P}_{\boldsymbol{\theta}}(X \in dx_{\mathcal{X}})$ est en général infinitésimalement petit, ce qui n'est pas commode. Néanmoins, dans un tel cas, on peut trouver (normalement) une mesure de référence $\text{vol}(\bullet)$ sur \mathcal{X} et une fonction $\ell_{x_{\mathcal{X}}}(\bullet)$ sur $\boldsymbol{\theta}$, telles que, indépendamment du choix du voisinage $dx_{\mathcal{X}}$, on ait toujours $\mathbb{P}_{\boldsymbol{\theta}}(X \in dx_{\mathcal{X}}) = \ell_{x_{\mathcal{X}}}(\boldsymbol{\theta}) \text{vol}(dx_{\mathcal{X}})$: dès lors, on peut faire rentrer le facteur $\text{vol}(d\mathcal{X})$ dans la constante de proportionnalité pour donner à la formule la formule suivante, sans infinitésimaux superflus :

$$\mathbb{P}_{\text{post}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}) \propto \ell_{x_{\mathcal{X}}}(\boldsymbol{\theta}) \times \mathbb{P}_{\text{pr}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}). \quad (\text{PD})$$

La quantité $\ell_{x_{\mathcal{X}}}(\boldsymbol{\theta})$ sera, dans la suite, appelée *vraisemblance que le paramètre caché vaille $\boldsymbol{\theta}$* (sachant que l'observation vaut $x_{\mathcal{X}}$); et la formule à laquelle nous sommes arrivés sera appelée *formule de Bayes*.

La formule

La formule de Bayes repose sur le concept de *vraisemblance*, qui sera développé en détail dans le chapitre 10. La définition générale de la vraisemblance (qui n'est pas à connaître, car en pratique on n'aura pas besoin d'y revenir pour faire nos calculs) est la suivante :

Définition (PE) (Vraisemblance). On utilise ici les notations génériques. La fonction de vraisemblance du paramètre caché au vu de l'observation effectivement réalisée, notée $\boldsymbol{\theta} \mapsto \mathcal{L}(\boldsymbol{\theta} = \boldsymbol{\theta} \mid X = x_{\mathcal{X}})$ — ou simplement $\mathcal{L}(\boldsymbol{\theta} = \boldsymbol{\theta})$, voire $\mathcal{L}(\boldsymbol{\theta})$, lorsqu'il n'y a pas d'ambiguïté —, est l'application qui, à $\boldsymbol{\theta} \in \Theta$, associe la probabilité $\mathbb{P}_{\boldsymbol{\theta}}(X \in dx_{\mathcal{X}})$ (où $dx_{\mathcal{X}}$ désigne un voisinage infinitésimal arbitraire de $x_{\mathcal{X}}$), *cette fonction devant néanmoins être considérée à constante multiplicative près*, où la constante (dont l'expression fera généralement intervenir le choix de $x_{\mathcal{X}}$) doit avoir un ordre de grandeur tel que la fonction de vraisemblance ne soit ni infinitésimalement petite, ni infinitésimalement grande. \heartsuit

Théorème (PF) (Formule de Bayes). *On utilise les notations génériques. Pour une observation effective $x_{\mathcal{X}}$, la postérieure sur $\boldsymbol{\theta}$ se déduit de la priorie en multipliant par la vraisemblance, à constante multiplicative près :*

$$\mathbb{P}_{\text{post}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}) \propto \mathcal{L}(\boldsymbol{\theta} = \boldsymbol{\theta} \mid X = x_{\mathcal{X}}) \times \mathbb{P}_{\text{pr}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}). \quad (\text{PG})$$

\diamond

Remarque (PH). Dans le cas où la priorie est une loi discrète, les infinitésimaux ne sont plus nécessaires : la formule devient alors simplement

$$\mathbb{P}_{\text{post}}(\boldsymbol{\theta} = \boldsymbol{\theta}) \propto \mathcal{L}(\boldsymbol{\theta} = \boldsymbol{\theta} \mid X = x_{\mathcal{X}}) \times \mathbb{P}_{\text{pr}}(\boldsymbol{\theta} = \boldsymbol{\theta}). \quad (\text{PI})$$

\clubsuit

Remarque (PJ). La formule de Bayes ne donne la postérieure qu'à constante multiplicative près; mais on peut toujours retrouver cette constante sans avoir à refaire les calculs depuis le début, en utilisant la condition de normalisation des lois de probabilité!

Concrètement, si on a obtenu une expression de la forme « $\mathbb{P}_{\text{post}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}) \propto f(\boldsymbol{\theta})\text{vol}(\text{d}\boldsymbol{\theta})$ », alors on va chercher la constante $Z \in \mathbb{R}_+^*$ telle que $\mathbb{P}_{\text{post}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}) = Z^{-1}f(\boldsymbol{\theta})\text{vol}(\text{d}\boldsymbol{\theta})$ en écrivant que

$$1 = \int_{\boldsymbol{\theta} \in \Theta} \mathbb{P}_{\text{post}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}) = \int_{\boldsymbol{\theta} \in \Theta} Z^{-1}f(\boldsymbol{\theta})\text{vol}(\text{d}\boldsymbol{\theta}) = \frac{\int_{\boldsymbol{\theta} \in \Theta} f(\boldsymbol{\theta})\text{vol}(\text{d}\boldsymbol{\theta})}{Z}, \quad (\text{PK})$$

d'où $Z = \int_{\boldsymbol{\theta} \in \Theta} f(\boldsymbol{\theta})\text{vol}(\text{d}\boldsymbol{\theta})$, et finalement

$$\mathbb{P}_{\text{post}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}) = \frac{f(\boldsymbol{\theta})}{\int_{\boldsymbol{\theta}' \in \Theta} f(\boldsymbol{\theta}')\text{vol}(\text{d}\boldsymbol{\theta}')} \text{vol}(\text{d}\boldsymbol{\theta}). \quad (\text{PL})$$

♣

8.3 Notion de priore impropre

En analyse bayésienne, on sera régulièrement confronté à une situation qu'on appelle la *priore impropre*.

! **Définition (PM).** Une priore impropre, c'est l'utilisation comme priore une loi de probabilité définie à constante multiplicative près... mais où la constante en question serait formellement infinie! ♥

Exemple (PN). Par exemple, la « loi de probabilité uniforme sur \mathbb{R} », définie formellement par

$$P(\text{d}x) \propto \text{vol}_1(\text{d}x), \quad (\text{PO})$$

autrement dit « $P(\text{d}x) = Z^{-1}\text{vol}_1(\text{d}x)$ » pour une constante de normalisation Z appropriée, n'est pas une véritable loi de probabilité, car il faudrait prendre $Z = \int_{x \in \mathbb{R}} \text{vol}_1(\text{d}x) = \infty$... Cependant cette expression n'est pas totalement dénuée de sens, car, par exemple, elle reflète l'idée selon laquelle il est deux fois plus probable de se trouver entre 3 et 5 qu'entre 8 et 9. ♣

Le point essentiel est que, dans la plupart des cas, le fait d'avoir pris une priore impropre n'empêche pas d'avoir une postérieure bien définie :

! **Proposition (PP).** *Lorsqu'on a un modèle avec une priore impropre, il faut faire tous les calculs en traitant la constante de normalisation dégénérée comme si c'était une véritable constante (mais sans jamais l'expliciter, évidemment) : cela n'empêche en rien de faire les calculs; et la plupart du temps, on arrivera à la fin à une postérieure qui sera une véritable loi de probabilité!* ◇

Exemple (PQ). Dans le cas du modèle du pédagogue, on pourrait envisager d'utiliser la « priore de Haar »^[*], qui consisterait en l'occurrence à prendre

$$\mathbb{P}(\boldsymbol{\mu} \in \text{d}\boldsymbol{\mu} \text{ et } \boldsymbol{\sigma} \in \text{d}\boldsymbol{\sigma}) \propto \boldsymbol{\sigma}^{-2} \text{vol}_1(\text{d}\boldsymbol{\mu}) \text{vol}_1(\text{d}\boldsymbol{\sigma}). \quad (\text{PR})$$

[*]. Il s'agit d'un type de priore non informative, confer § 15.5.

Il s'agit là d'une priore impropre, car l'intégrale de $\sigma^{-2} \text{vol}_1(d\mu) \text{vol}_1(d\sigma)$ sur Θ est infinie... Cependant, comme nous le verrons dans la § 8.4 ci-dessous, l'application de la formule de Bayes conduit alors à trouver pour postérieure

$$\mathbb{P}(\mu \in d\mu \text{ et } \sigma \in d\sigma \mid X = x_{\mathcal{J}}) \propto \sigma^{-24} \exp\left(-\frac{11 \times (\mu - 70,3)^2 + 2,54 \times 10^3}{\sigma^2}\right) \text{vol}_2(d\mu \times d\sigma), \quad (\text{PS})$$

qui cette fois-ci décrit bien une loi propre, vu que l'intégrale de $\sigma^{-24} \exp(-((11(\mu - 70,3)^2 + 2540)\sigma^{-2}))$ est convergente sur $\mathbb{R} \times \mathbb{R}_+^*$! \clubsuit

8.4 Exemples de calculs

Calcul de postérieure pour le chasseur

Quelle est la fonction de vraisemblance pour le modèle du chasseur ? \mathcal{X} étant discret, nous pouvons appliquer le théorème (TY) (confer chapitre 10) :

$$\begin{aligned} \mathcal{L}(\theta = \theta \mid X = x_{\mathcal{J}}) &= \mathbb{P}(X = x_{\mathcal{J}} \mid \theta = \theta) = \mathbb{P}(\text{Binom}^{\text{le}}(n, \theta) = x_{\mathcal{J}}) \\ &= \binom{n}{x_{\mathcal{J}}} \theta^{x_{\mathcal{J}}} (1 - \theta)^{n - x_{\mathcal{J}}}. \end{aligned}$$

Cependant, nous rappelant que la vraisemblance n'est définie, en tant que fonction de θ , qu'à constante multiplicative près, il s'avère judicieux de simplifier le facteur constant $\binom{n}{x_{\mathcal{J}}}$ pour alléger les calculs, et de prendre ainsi

$$\mathcal{L}(\theta = \theta \mid X = x_{\mathcal{J}}) = \theta^{x_{\mathcal{J}}} (1 - \theta)^{n - x_{\mathcal{J}}}. \quad (\text{PT})$$

Pour calculer la loi à postériori, nous avons également besoin de nous rappeler la priore sur θ : en l'occurrence, comme expliqué en page 82, il s'agit de la loi arcsinus, caractérisée par :

$$\mathbb{P}_{\text{pr}}(\theta \in d\theta) := Z_{\text{pr}}^{-1} \theta^{-1/2} (1 - \theta)^{-1/2} \text{vol}_1(d\theta), \quad (\text{PU})$$

où Z_{pr} est la constante de normalisation faisant de cette mesure une loi de probabilité. (Nous savons en fait que celle-ci vaut en l'occurrence π ; mais je souhaite insister sur le fait que la détermination cette valeur est sans importance pour notre raisonnement).

En vertu de la formule de Bayes, on trouve alors que, pour $d\theta$ un voisinage infinitésimal de θ :

$$\begin{aligned} \mathbb{P}(\theta \in d\theta \mid X = x_{\mathcal{J}}) &\propto \mathcal{L}(\theta = \theta \mid X = x_{\mathcal{J}}) \times \mathbb{P}_{\text{pr}}(\theta \in d\theta) \\ &= \theta^{x_{\mathcal{J}}} (1 - \theta)^{n - x_{\mathcal{J}}} \times Z_{\text{pr}}^{-1} \theta^{-1/2} (1 - \theta)^{-1/2} \text{vol}_1(d\theta) \\ &\propto \theta^{x_{\mathcal{J}} - 1/2} (1 - \theta)^{n - x_{\mathcal{J}} - 1/2} \text{vol}_1(d\theta) = \theta^{7/2} (1 - \theta)^{16/2} \text{vol}_1(d\theta). \end{aligned}$$

Nous avons ainsi déterminé la postérieure pour θ ... mais seulement à constante multiplicative près. Cependant, si le but est uniquement d'*identifier* la postérieure, il est inutile d'aller plus loin : en effet, cette constante est forcément égale à l'unique valeur qui donne à la postérieure une masse totale égale à 1 ! Or, en l'occurrence, un

statisticien expert connaît justement une loi classique, appelée la loi Bêta($8\frac{1}{2}$, $17\frac{1}{2}$), qui est une loi portée par $]0, 1[$ et telle que

$$\mathbb{P}(\text{Bêta}(8\frac{1}{2}, 17\frac{1}{2}) \in dx) \propto x^{7\frac{1}{2}}(1-x)^{16\frac{1}{2}} \text{vol}_1(dx) : \quad (\text{PV})$$

la postérieure coïncidant avec la loi Bêta($8\frac{1}{2}$, $17\frac{1}{2}$) à constante multiplicative près, elle y est nécessairement *égale*, étant donné que postérieure comme loi bêta sont des lois de probabilité, et qu'elles doivent donc avoir la même masse totale (à savoir, 1).

Noter par ailleurs que, si on avait quand même besoin de déterminer la constante de proportionnalité pour la postérieure (cela arrive parfois), le calcul se ferait simplement en écrivant qu'on soit satisfaire la condition de normalisation d'une mesure de probabilité : on a donc

$$\mathbb{P}_{\text{post}}(\theta \in d\theta) = Z_{\text{post}}^{-1} \theta^{7\frac{1}{2}}(1-\theta)^{16\frac{1}{2}} \text{vol}_1(d\theta) \quad (\text{PW})$$

avec

$$Z_{\text{post}} = \int_{\theta \in \Theta} \theta^{7\frac{1}{2}}(1-\theta)^{16\frac{1}{2}} \text{vol}_1(d\theta) = \int_{\theta=0}^1 \theta^{7\frac{1}{2}}(1-\theta)^{16\frac{1}{2}} d\theta = 7,748 \times 10^{-8}. \quad (\text{PX})$$

Notez en particulier qu'on n'a pas besoin, pour faire ce calcul, de connaître la valeur Z_{pr} de la constante de normalisation de la priore! \curvearrowright

Nous avons donc déterminé la loi à postériori de la fiabilité du chasseur. Concrètement, cela a (notamment) l'interprétation suivante : si on définit un « mauvais chasseur » comme un chasseur dont la fiabilité est inférieure à $1/4$, alors, *au vu du résultat de son test*, le club estime que la probabilité que son candidat soit un mauvais chasseur vaut $\mathbb{P}_{\text{post}}(\theta < 1/4) = \mathbb{P}(\text{Bêta}(8\frac{1}{2}, 17\frac{1}{2}) < 1/4) = 20,6 \%$; alors que, si on avait demandé au club ce qu'il pensait du candidat *avant* le test, il lui aurait donné une probabilité d'être un mauvais chasseur égale à $\mathbb{P}_{\text{pr}}(\theta < 1/4) = \mathbb{P}(\text{LoiArcsin} < 1/4) = 33,3 \%$. Le test a donc un peu rassuré le club quant au risque que son candidat soit un mauvais chasseur, mais il demeure un doute substantiel!

Loi à postériori d'une quantité d'intérêt explicative pour le chasseur

Si la club ne s'intéresse pas à θ elle-même, mais à la quantité d'intérêt explicative $\tau := 2\theta - \theta^2 =: \varphi(\theta)$ (laquelle quantité correspond à la probabilité de toucher un plateau quand on a droit à *deux* essais), il va tout simplement dire que la loi (à postériori) de τ est la *mesure-image* de celle de θ par l'application $\varphi(\bullet)$: en effet, pour $\tau \in]0, 1[$ et $d\tau$ un voisinage infinitésimal de τ , dire que τ (qui est $\varphi(\theta)$) est dans $d\tau$, c'est la même chose que dire que θ est dans l'image réciproque $\varphi^{-1}(d\tau)$!

En l'occurrence, l'application $\varphi(\bullet)$ étant bijective de $]0, 1[$ dans $]0, 1[$, le calcul va correspondre à un simple changement de variable. Plus précisément, par bijectivité, l'image réciproque de $d\tau$ par $\varphi(\bullet)$ est aussi son image directe par la bijection réciproque $\varphi^{-1}(\bullet)$, dont on détermine facilement la formule :

$$\varphi^{-1}(\tau) = 1 - (1 - \tau)^{1/2}. \quad (\text{PY})$$

Puisque $d\tau$ est un voisinage infinitésimal de τ , $\varphi^{-1}(d\tau)$ est un voisinage infinitésimal de $d\tau$. En outre, en vertu de la formule de changement de variable, la longueur de $\varphi^{-1}(d\tau)$ est $|[\varphi^{-1}]'(\tau)|$ fois plus grande que celle de $d\tau$: d'où, après calcul de dérivée,

$$\text{vol}_1(\varphi^{-1}(d\tau)) = \frac{1}{2(1-\tau)^{1/2}} \text{vol}_1(d\tau)^{[\ddagger]}. \quad (\text{PZ})$$

L'un dans l'autre, on a donc

$$\begin{aligned}\mathbb{P}_{\text{post.}}(\boldsymbol{\tau} \in d\boldsymbol{\tau}) &= \mathbb{P}_{\text{post.}}(\boldsymbol{\theta} \in \boldsymbol{\varphi}^{-1}(d\boldsymbol{\tau})) = Z_{\text{post.}}^{-1} \boldsymbol{\varphi}^{-1}(\boldsymbol{\tau})^{7/2} (1 - \boldsymbol{\varphi}^{-1}(\boldsymbol{\tau}))^{16/2} \text{vol}_1(\boldsymbol{\varphi}^{-1}(d\boldsymbol{\tau})) \\ &= \frac{1}{2Z_{\text{post.}}} (1 - \boldsymbol{\tau})^{7/4} (1 - (1 - \boldsymbol{\tau})^{1/2})^{7/2} \text{vol}_1(d\boldsymbol{\tau}) : \quad (\text{QA})\end{aligned}$$

on a donc déterminé l'expression explicite (à constante de normalisation près) de la loi à postérieure de $\boldsymbol{\tau}$!

Loi à postérieure de la quantité d'intérêt prédictive pour le chasseur

Maintenant, voyons comment on peut passer de la postérieure (concernant le paramètre caché lui-même) à la loi à postérieure de la quantité d'intérêt prédictive. Nous utilisons le raisonnement présenté dans le point (OX) : pour $y \in \llbracket 0, m \rrbracket$, on a par formule des probabilités totales sur $\boldsymbol{\theta}$, puis par règle de chaîne, que

$$\begin{aligned}\mathbb{P}_{\text{post.}}(Y = y) &= \int_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \mathbb{P}_{\text{post.}}(\boldsymbol{\theta} \in d\boldsymbol{\theta} \text{ et } Y = y) \\ &= \int_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \mathbb{P}_{\text{post.}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}) \mathbb{P}_{\text{post.}}(Y = y \mid \boldsymbol{\theta} = \boldsymbol{\theta}) ; \quad (\text{QB})\end{aligned}$$

puis, dans l'expression « $\mathbb{P}_{\text{post.}}(Y = y \mid \boldsymbol{\theta} = \boldsymbol{\theta})$ », qui fait intervenir à la fois le conditionnement par $\{X = x_{\checkmark}\}$ (en tant que contexte probabiliste) et le conditionnement par $\{\boldsymbol{\theta} = \boldsymbol{\theta}\}$ (en tant que contexte explicite), on échange les rôles que ces deux conditionnement jouent, obtenant ainsi que cette expression peut aussi s'écrire $\mathbb{P}_{\boldsymbol{\theta}}(Y = y \mid X = x_{\checkmark})$: ce qui, pour ce modèle spécifique où X et Y sont indépendantes sous tous les contextes $\mathbb{P}_{\boldsymbol{\theta}}(\bullet)$, se simplifie encore en $\mathbb{P}_{\boldsymbol{\theta}}(Y = y)$.

On utilise alors les formules pour resp. la postérieure et la loi de l'observation sachant le paramètre caché pour obtenir que

$$\begin{aligned}\mathbb{P}_{\text{post.}}(Y = y) &= \int_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \mathbb{P}_{\boldsymbol{\theta}}(Y = y) \mathbb{P}_{\text{post.}}(\boldsymbol{\theta} \in d\boldsymbol{\theta}) \\ &\propto \int_{\boldsymbol{\theta} \in]0, 1[} \mathbb{P}(\text{Binom}^{\text{le}}(m, \boldsymbol{\theta}) = y) \boldsymbol{\theta}^{7/2} (1 - \boldsymbol{\theta})^{16/2} \text{vol}_1(d\boldsymbol{\theta}), \quad (\text{QC})\end{aligned}$$

et on utilise enfin que

$$\mathbb{P}(\text{Binom}^{\text{le}}(m, \boldsymbol{\theta}) = y) = \binom{m}{y} \boldsymbol{\theta}^y (1 - \boldsymbol{\theta})^{m-y} = \frac{75!}{y!(75-y)!} \boldsymbol{\theta}^y (1 - \boldsymbol{\theta})^{75-y} \quad (\text{QD})$$

pour arriver à l'expression de la loi à postérieure pour y , à constante de normalisation près :

$$\begin{aligned}\mathbb{P}_{\text{post.}}(Y = y) &\propto \frac{75!}{y!(75-y)!} \int_{\boldsymbol{\theta}=0}^1 \boldsymbol{\theta}^{y+7/2} (1 - \boldsymbol{\theta})^{91/2-y} \text{vol}_1(d\boldsymbol{\theta}) \\ &= \frac{75!}{y!(75-y)!} \frac{(y+7/2)!(91/2-y)!}{100!} \propto \frac{(y+7/2)!(91/2-y)!}{y!(75-y)!}, \quad (\text{QE})\end{aligned}$$

[†]. En l'occurrence la valeur absolue n'apparaît pas, car la dérivée est toujours positive.

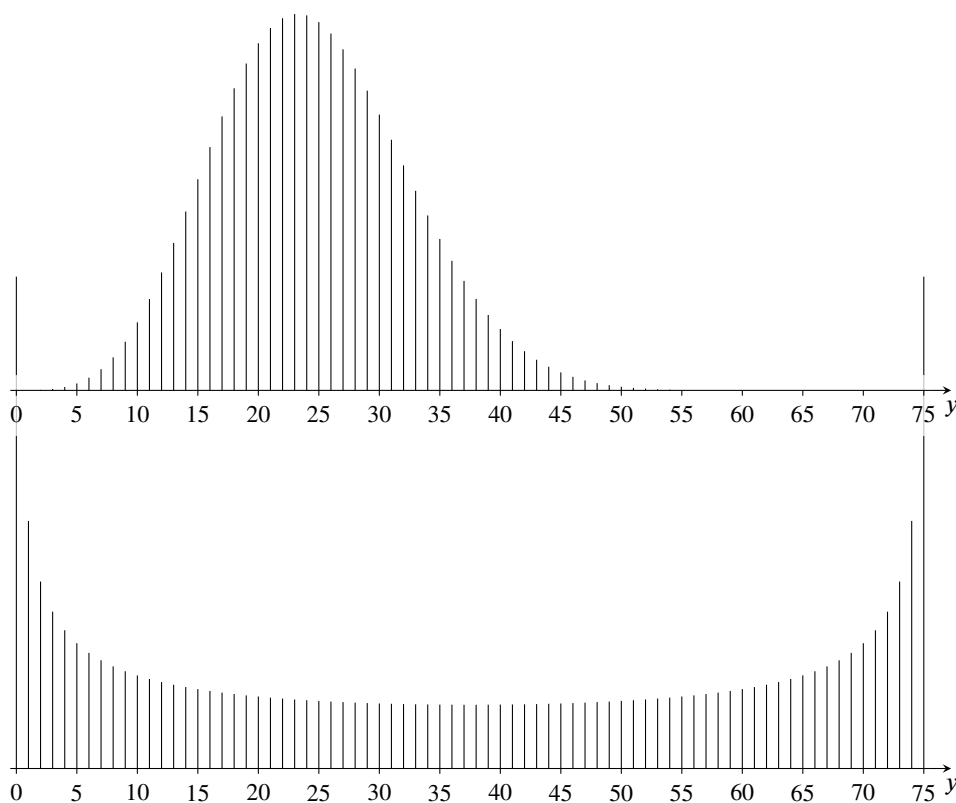


FIGURE 8.1 – En haut, la loi à postérieure de Y dans le modèle du chasseur. En bas, à titre de comparaison, la loi à priori de Y , tracée à la même échelle.

où on a utilisé la formule de l'« intégrale bêta » (que vous n'êtes pas supposés connaître) pour l'égalité centrale ^[‡].

On peut enfin récupérer la constante de proportionnalité de cette loi en écrivant la condition de normalisation : ce qui donne finalement que

$$\mathbb{P}_{\text{post}}(Y = y) = Z_Y^{-1} \frac{(y + 7\frac{1}{2})! (91\frac{1}{2} - y)!}{y! (75 - y)!}, \quad (\text{QF})$$

avec

$$Z_Y = \sum_{y=0}^{75} \frac{(y + 7\frac{1}{2})! (91\frac{1}{2} - y)!}{y! (75 - y)!}. \quad (\text{QG})$$

Cette loi peut se déterminer numériquement : vous en trouverez le tracé en figure 8.1, accompagné du tracé de la loi à priori de Y à titre de comparaison. On voit que, alors qu'à priori le club n'avait pour ainsi dire aucune idée du score que son candidat ferait en compétition, grâce au résultat du test, le club considère à présent qu'il peut être quasi-certain (à 97,8 % de probabilité, plus précisément) que le candidat atteindra entre 8 et 42 cibles à la compétition !

Calcul de postérieure pour le pédagogue

Comme dans l'exemple précédent, nous commençons par calculer la vraisemblance pour l'observation effective $(x_{0\checkmark}, \dots, x_{(n-1)\checkmark})$. Comme l'observation $(X_{1\checkmark}, \dots,$

^[‡]. La formule fait intervenir des factorielles d'arguments non entiers : mais les ordinateurs savent bien calculer ces dernières.

$X_{n\checkmark}$) se décompose en n sous-observations indépendantes, nous pouvons appliquer le théorème (UD) :

$$\mathcal{L}(\boldsymbol{\theta} = (\mu, \sigma) \mid X = (x_{0\checkmark}, \dots, x_{(n-1)\checkmark})) = \prod_{i=0}^{n-1} \mathcal{L}(\boldsymbol{\theta} = (\mu, \sigma) \mid X_i = x_{i\checkmark}), \quad (\text{QH})$$

où X_i désigne le score obtenu par l'élève numéro i . Maintenant, quelle est la vraisemblance $\mathcal{L}(\boldsymbol{\theta} \mid X_i = x_{i\checkmark})$ du paramètre caché pour la "sous-expérience" consistant à regarder seulement le score de l'élève i ? On remarque que, sachant que $\boldsymbol{\theta}$ vaut (μ, σ) , la loi de X_i est la loi Normale (μ, σ^2) , caractérisée par

$$\mathbb{P}(X_i \in dx \mid \boldsymbol{\theta} = (\mu, \sigma)) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \text{vol}_1(dx) : \quad (\text{QI})$$

vu qu'il s'agit d'une loi à densité sur \mathbb{R} , on peut appliquer le théorème (UA) pour en déduire que

$$\mathcal{L}(\boldsymbol{\theta} = (\mu, \sigma) \mid X_i = x_{i\checkmark}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_{i\checkmark} - \mu)^2}{2\sigma^2}\right), \quad (\text{QJ})$$

expression dans laquelle nous nous empresserons évidemment de simplifier le facteur constant (par rapport à $\boldsymbol{\theta}$) $1/\sqrt{2\pi}$. En appliquant le théorème (UD), on obtient alors :

$$\mathcal{L}(\boldsymbol{\theta} = (\mu, \sigma) \mid X = (x_{0\checkmark}, \dots, x_{(n-1)\checkmark})) = \sigma^{-n} \exp\left(-\frac{\sum_{i=0}^{n-1} (x_i - \mu)^2}{2\sigma^2}\right). \quad (\text{QK})$$

Il s'avèrera commode de ré-organiser encore cette expression, en particulier la quantité intervenant dans l'exponentielle. Pour cela, nous introduisons la *moyenne empirique* des X_i , soit $M := (\sum_{i=0}^{n-1} X_i) / n$, ainsi que leur *écart-type empirique*, soit $S := (\sum_{i=0}^{n-1} (X_i - M)^2 / n)^{1/2}$ (statistiques dont, suivant nos conventions, les réalisations seront notées resp. m_{\checkmark} et s_{\checkmark}) : on vérifie aisément que

$$\sum_{i=0}^{n-1} (X_i - \mu)^2 = n((M - \mu)^2 + S^2), \quad (\text{QL})$$

d'où :

$$\mathcal{L}(\boldsymbol{\theta} = (\mu, \sigma) \mid X = (x_{0\checkmark}, \dots, x_{(n-1)\checkmark})) = \sigma^{-n} \exp\left(-\frac{1}{2}n\sigma^{-2}((m_{\checkmark} - \mu)^2 + s_{\checkmark}^2)\right). \quad (\text{QM})$$

Au vu de la formule pour notre prior (confer p. ??), la formule de Bayes nous donne finalement que

$$\begin{aligned} \mathbb{P}(\boldsymbol{\theta} \in d\boldsymbol{\mu} \times d\sigma \mid X = x_{\checkmark}) &\propto \\ &\sigma^{-(n+3)} \exp\left(-\frac{n(\mu - m_{\checkmark})^2 + ns_{\checkmark}^2 + \sigma_{\text{réf}}^2}{2\sigma^2} - \frac{(\mu - \mu_{\text{réf}} + \sigma_{\text{réf}})^2}{2\sigma_{\text{réf}}^2}\right) \text{vol}_1(d\boldsymbol{\mu}) \text{vol}_1(d\sigma) \\ &= \sigma^{-25} \exp\left(-\frac{11 \times (\mu - 70,3)^2 + 2,77 \times 10^3}{\sigma^2} - \frac{(\mu - 64)^2}{242}\right) \text{vol}_2(d\boldsymbol{\mu} \times d\sigma). \end{aligned}$$

On notera que l'expression de cette postérieure est assez compliquée ; en particulier, $\boldsymbol{\mu}$ et σ ne sont plus indépendantes pour la loi à postérieure.

Le figure 8.2 donne l'allure de cette postérieure, obtenue numériquement à partir de la formule ci-dessus. On remarque que l'observation des notes de la première promotion a permis à l'enseignant d'affiner considérablement sa connaissance du paramètre caché : à présent, il est en mesure de dire que μ_{\checkmark} se situe probablement entre 65 et 75, et que σ_{\checkmark} se situe probablement entre 12 et 20 !

Loi à postériori de μ pour le pédagogue

Si ce qui intéresse spécifiquement notre enseignant est de connaître la valeur de la quantité d'intérêt μ , il peut obtenir la loi à postériori de μ à partir de la loi à postériori de (μ, σ) par mesure-image : ce qui, en l'occurrence, revient simplement à intégrer selon les différentes possibilités pour σ :

$$\begin{aligned} \mathbb{P}_{\text{post}}(\mu \in d\mu) &= \\ \int_{\sigma \in \mathbb{R}_+^*} \mathbb{P}_{\text{post}}(\mu \in d\mu \text{ et } \sigma \in d\sigma) &\propto \exp\left(-\frac{(\mu - \mu_{\text{réf}} + \sigma_{\text{réf}})^2}{2\sigma_{\text{réf}}^2}\right) \text{vol}_1(d\mu) \times \\ &\int_{\sigma \in \mathbb{R}_+^*} \sigma^{-(n+3)} \exp\left(-\frac{n(\mu - m_{\mathcal{J}})^2 + ns_{\mathcal{J}}^2 + \sigma_{\text{réf}}^2}{2\sigma^2}\right) \text{vol}_1(d\sigma) \quad (\text{QN}) \end{aligned}$$

Ici on utilise que, pour $k > 1$, pour $a \in \mathbb{R}_+^*$, on a par le changement de variable $y \leftarrow a^{-1/2}x$ que

$$\int_{x=0}^{\infty} x^{-k} e^{-ax^2} dx = \int_{y=0}^{\infty} (a^{1/2}y)^{-k} e^{-1/y^2} a^{1/2} dy \propto a^{-(k-1)/2} \quad (\text{QO})$$

(la proportionnalité s'entendant à valeur de k fixée) pour en déduire que

$$\mathbb{P}_{\text{post}}(\mu \in d\mu) \propto \exp\left(-\frac{(\mu - \mu_{\text{réf}} + \sigma_{\text{réf}})^2}{2\sigma_{\text{réf}}^2}\right) \times \frac{1}{(n(\mu - m_{\mathcal{J}})^2 + ns_{\mathcal{J}}^2 + \sigma_{\text{réf}}^2)^{n/2+1}} \text{vol}_1(d\mu). \quad (\text{QP})$$

Variante avec la priore de Haar

Si, au lieu de prendre la priore par défaut pour notre modèle de pédagogue, nous avons pris la priore de Haar (cf. exemple (PQ)), tout se serait passé de la même manière (à condition de ne jamais chercher à expliciter les constantes de normalisation !), à ceci près que nous aurions dû remplacer le facteur « $\sigma^{-3} \exp(-(\mu - \mu_{\text{réf}} + \sigma_{\text{réf}})^2 / 2\sigma_{\text{réf}}^2 - \sigma_{\text{réf}}^2 / 2\sigma^2)$ » qui apparaît dans la priore par défaut par le facteur « σ^{-2} » qui en est l'analogie pour la priore de Haar : on serait alors arrivé à

$$\begin{aligned} \mathbb{P}(\theta \in d\mu \times d\sigma \mid X = x_{\mathcal{J}}) &\propto \sigma^{-(n+2)} \exp\left(-\frac{n(\mu - m_{\mathcal{J}})^2 + ns_{\mathcal{J}}^2}{2\sigma^2}\right) \text{vol}_1(d\mu) \text{vol}_1(d\sigma) \\ &= \sigma^{-24} \exp\left(-\frac{11 \times (\mu - 70,3)^2 + 2,54 \times 10^3}{\sigma^2}\right) \text{vol}_2(d\mu \times d\sigma); \quad (\text{QQ}) \end{aligned}$$

et si nous avons voulu nous intéresser plus spécifiquement à la quantité d'intérêt μ , le même argument dans le cas de la priore par défaut nous aurait permis d'obtenir que, cette fois-ci,

$$\mathbb{P}_{\text{post}}(\mu \in d\mu) \propto \frac{\text{vol}_1(d\mu)}{((\mu - m_{\mathcal{J}})^2 + s_{\mathcal{J}}^2)^{(n+1)/2}}. \quad (\text{QR})$$

On observe ainsi que le résultat de l'analyse dépend du choix de la priore ! Voir une illustration de ce phénomène en figure 8.3.

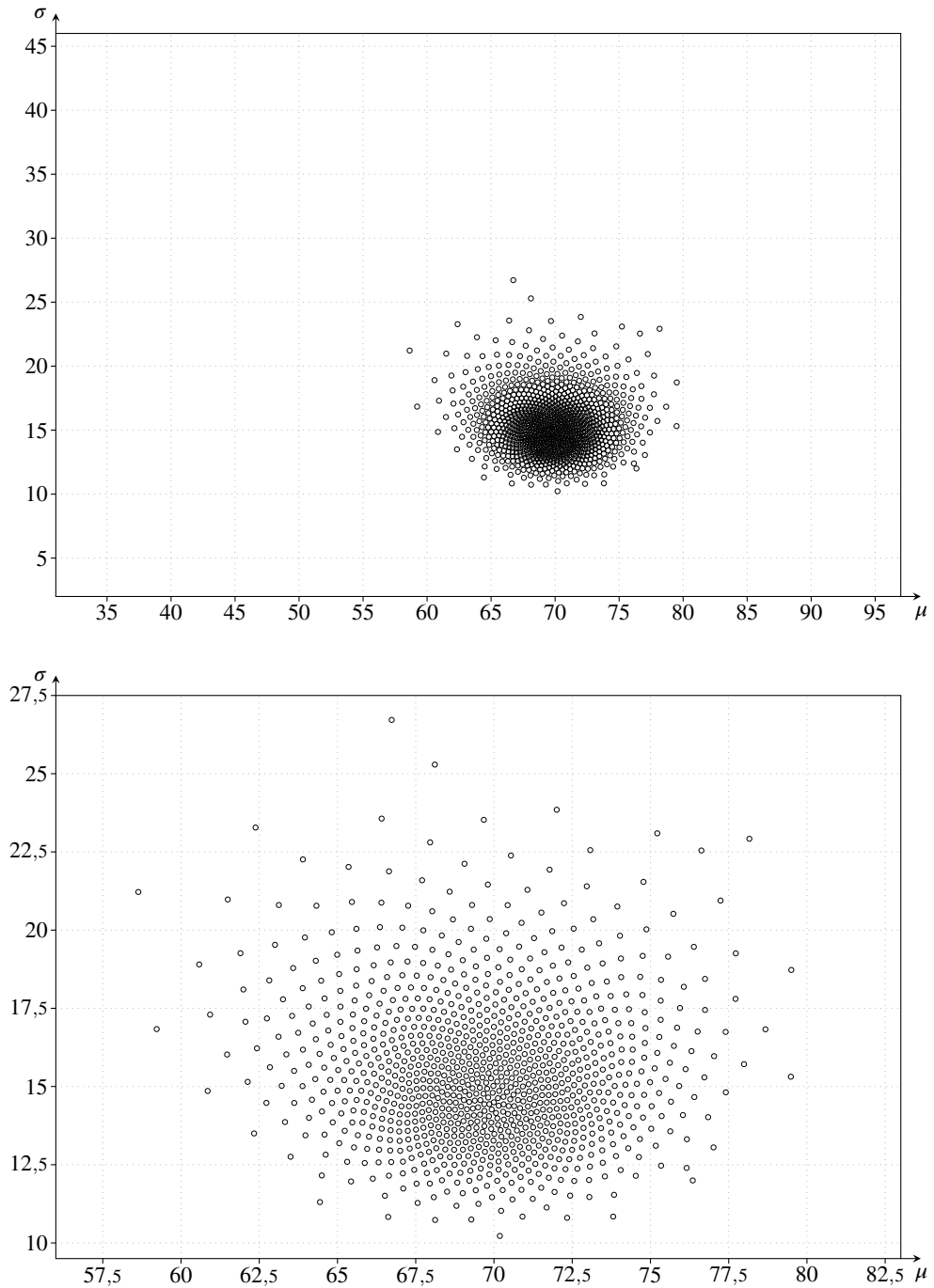


FIGURE 8.2 – Postérieure pour le modèle du pédagogue, avec un petit cercle par zone de probabilité 1 ‰ : en haut, avec la même échelle que pour la priore ; en bas, zoom sur la zone la plus chargée.

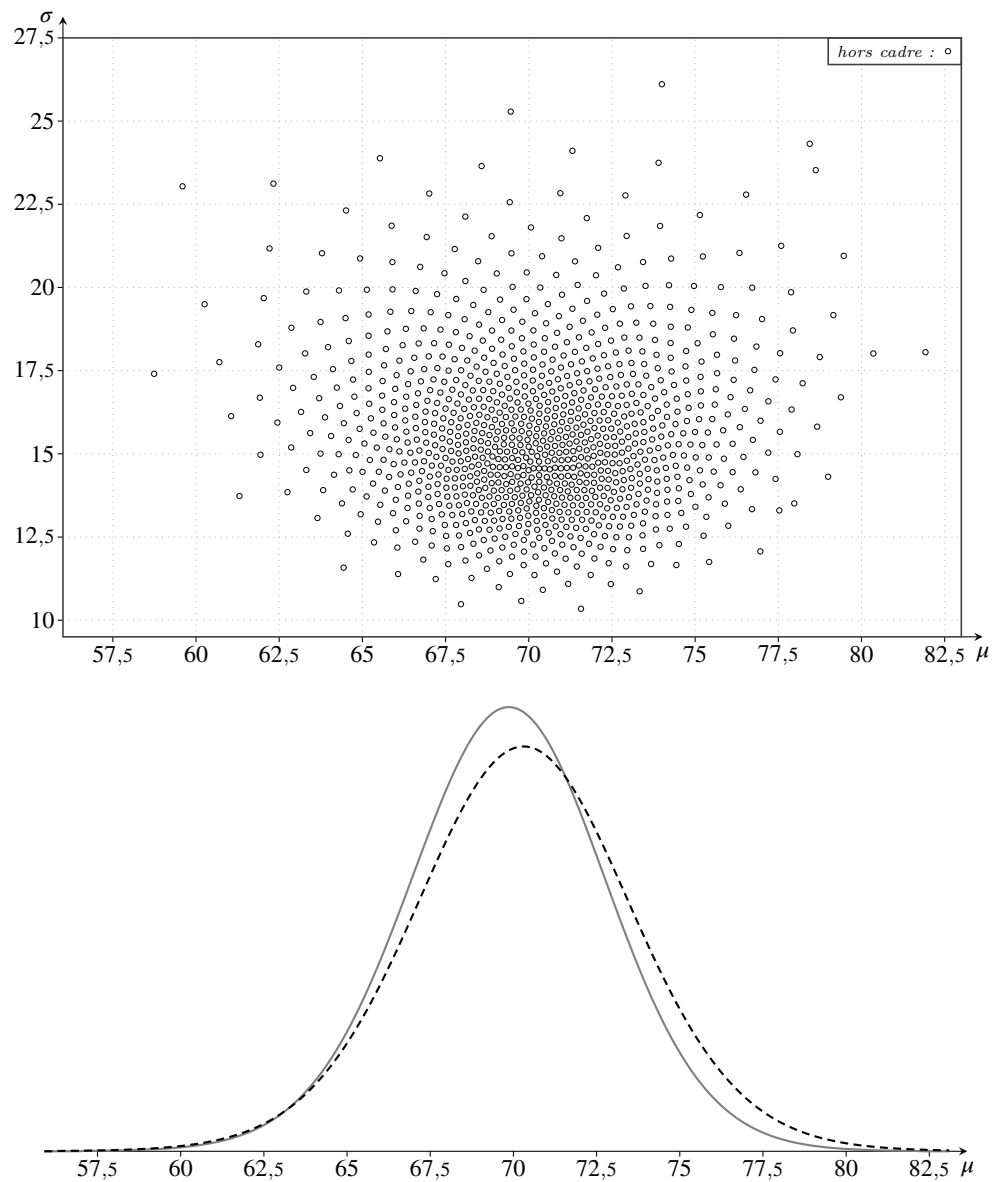


FIGURE 8.3 – En haut : Postérieure pour le modèle du pédagogue avec priore de Haar, à la même échelle qu'en figure 8.2. En bas : Densité de la loi à postérieure de μ dans les cas respectifs de la priore de référence (trait plein gris) et de la priore de Haar (trait pointillé noir).

Chapitre 9

Méthodes statistiques bayésiennes

9.1 Importance des lois à postériori

À l'issue du chapitre précédent, nous avons appris à exploiter le théorème de Bayes pour, une certaine priore étant donnée, déduire de notre observation la postérieure, puis, par ricochet, les *lois à postériori* de quantités d'intérêt $\gamma(\theta)$ (cas explicatif) ou $g(Y)$ (cas prédictif). Il est important d'observer que lois à postériori sont des *statistiques*, et par conséquent, constituent des réponses techniquement acceptables pour une analyse statistique, de même que toute quantité s'en déduisant :

Remarque (QS). La postérieure ne dépendant que^[*] de l'observation, et pas du paramètre caché, il s'agit de la réalisation d'une statistique. Du point de vue technique, la statistique en question peut s'écrire comme l'image de la v.a. X par la fonction déterministe^[†]

$$\begin{aligned} \mathcal{X} &\rightarrow \mathcal{M}_1(\Theta) \\ x &\mapsto \text{Loi}(\theta \mid X = x) : \end{aligned} \tag{QT}$$

Il s'agit ainsi d'une variable aléatoire (en tant que fonction de X) à valeurs dans les *lois de probabilité* sur Ω : qu'une variable aléatoire soit à valeurs « lois » est certes un peu paradoxal, avec un aspect “méta” ; mais cela ne pose aucune, objection du point de vue technique ! ☺

Corolairement, dans la mesure où nous avons vu que la loi à postériori d'une quantité d'intérêt $\gamma(\theta)$ ou $g(Y)$ pouvait se calculer directement à partir de la postérieure (et, éventuellement, de la valeur de l'observation passée), ces lois à postériori sont elles aussi des statistiques. ☺

En fait, comme nous le verrons dans la section suivante, non seulement la loi à postériori de $g(Y)$ (mettons ici qu'on est dans un contexte prédictif) est une statistique, mais elle est, dans un sens, “exhaustive”, en ce qu'elle décrit l'*intégralité* de l'information (et de l'incertitude) dont nous disposons sur $g(Y)$ à l'issue de notre analyse ! C'est donc un jalon absolument fondamental dans l'exploitation de l'observation dont nous disposons.

Cependant, du point de vue pratique, la loi à postériori d'une quantité d'intérêt comporte deux écueils :

[*]. Outre, s'entend, du modèle statistique et de la priore choisie.

[†]. Regardez attentivement la formule : même si le symbole ‘ X ’ intervient, c'est bien une fonction *déterministe* qui est décrite !

- Tout d’abord, elle ne dit pas directement *quoi faire* : car en fin de compte, notre but ne sera pas de déterminer une loi, mais de prendre une décision...!
- D’autre part, à supposer que notre but soit simplement de transmettre l’information obtenue sur la quantité d’intérêt au commanditaire de notre étude, l’objet « loi de probabilité à posteriori », qui est à valeurs dans l’espace de dimension (en général) infinie $\mathcal{M}_1(\mathcal{G})$, est particulièrement compliqué à décrire complètement, et sa description complète risque donc d’être particulièrement compliquée à appréhender par lesdits commanditaires... Il peut donc s’avérer utile de “résumer” cette distribution à posteriori par simplement une poignée de nombres réels.

La suite du présent chapitre vise à traiter ces deux problématiques. La § 9.2 s’attaquera à la question très générale de « quoi faire » ; tandis que les § ?? à ?? viseront plus modestement à résumer la loi à posteriori en quelques nombres-clés.

9.2 Décision optimale en contexte bayésien

Loi à posteriori et décision optimale

Dans cette section, nous nous intéresserons uniquement au cas d’une quantité d’intérêt *prédictive* $g(Y) =: G$: ce cas permet en effet de traiter aussi le cas explicatif, quitte à adjoindre à l’observation future une copie du paramètre caché.

Nous supposons choisie une certaine fonction de perte $\ell(\bullet, \bullet)$ pour notre problème, le premier argument de cette fonction de perte se référant à une quantité d’intérêt prédictive $g(Y) =: G$, et le second à une décision à prendre au sein d’un certain espace \mathcal{D} . Cette décision ne pouvant être prise qu’en fonction des informations dont nous disposons déjà au moment de l’analyse, cela devra être une fonction de l’observation, autrement dit, une *statistique* : nous la noterons donc $\hat{d}(X)$.

Notre but est alors de choisir une $\hat{d}(X) =: \hat{D}$ de façon à minimiser l’espérance de notre perte (confer § 3.5) :

$$\text{On cherche } \hat{d} : \mathcal{X} \rightarrow \mathcal{D} \text{ minimisant } \mathbb{E}(\ell(g(Y), \hat{d}(X))). \quad (\text{Opt}_{\mathbb{E}})$$

Le point essentiel est alors le suivant :

!! **Principe (QU)** (Décision optimale en contexte bayésien, principe général). *La réalisation $\hat{d}(x_{\mathcal{J}})$ de la décision optimale se trouve en résolvant le problème probabiliste où on doit prendre une décision déterministe dans un contexte où G suit la loi $\text{Loi}_{\text{post}}(G)$! Un énoncé plus formel est donné ci-après.* \diamond

Dit en termes plus formels :

Théorème (QV) (Décision optimale en contexte bayésien, version formelle). *Pour tout $x \in \mathcal{X}$, notons*

$$P_x := \text{Loi}(g(Y) \mid X = x), \quad (\text{QW})$$

et notons \mathcal{S}_x l’ensemble des solutions du problème d’optimisation « particularisé » suivant :

$$\text{On cherche } \hat{d} \in \mathcal{D} \text{ maximisant } \mathbb{E}^{G \sim P_x}(u(G, \hat{d})). \quad (\text{Opt}_x)$$

Alors^[†] une application $\hat{d} : \mathcal{X} \rightarrow \mathcal{D}$ constitue une solution optimale au problème $(\text{Opt}_{\mathbb{E}})$ si et seulement si on a $\hat{d}(x) \in \mathcal{S}_x$ pour (presque-)tout x . \diamond

! *Remarque (QX)*. Le principe (QU) montre donc que, dans le paradigme bayésien, seule la *réalisation* de la postérieure importe pour prendre une décision optimale : c'est extrêmement commode! ☺ En statistique fréquentiste, par contre, cette propriété ne sera plus valable... ☹

Remarque (QY). En général, dans la pratique, le but n'est pas tant de déterminer l'optimum exact (à supposer qu'il existe) que de s'approche très près de l'optimum (ce qui est toujours possible, quand bien même l'optimum lui-même n'est pas atteignable). Mais cela ne change rien au principe (QU) : il faut juste remplacer l'optimisation exacte par une optimisation *approchée*! ☺

Remarque (QZ). Nous comprenons donc pourquoi, dans la section précédente, nous avons dit que la loi à postériori de $g(Y)$ fournissait *l'intégralité* de l'information dont nous disposons sur la quantité d'intérêt au vu de l'observation : de fait, il suffit de connaître cette loi à postériori pour prendre optimalement n'importe quelle décision concernant $g(Y)$!

Notions de risques

La thématique de la prise de décision en contexte bayésien est l'occasion d'introduire le vocabulaire lié aux différentes notions de *risque*. En statistique, « risque » signifie simplement « espérance de la perte » :

Définition (RA) (Marco-définition du « risque »). En statistique, une quantité d'intérêt $g(Y)$ et une stratégie de décision $d(X)$ étant données, ainsi qu'une fonction de perte $\ell(\bullet, \bullet)$ relative à l'inadéquation entre telle valeur de la quantité d'intérêt et telle possibilité de décision, on parle de *risque* associé à cette stratégie de décision pour désigner, de manière générale, « l'espérance de la perte $\ell(g(Y), \hat{d}(X))$ ». ♥

Dans la mesure où, en statistique, la notion d'« espérance » se rapporte en fait à un certain *contexte probabiliste*, à chaque contexte probabiliste de trouve associé une notion de « risque » différente :

Définition (RB). Une stratégie de décision $\hat{d}(X) =: \hat{D}$ étant donnée, dans un contexte associé à une quantité d'intérêt $g(Y)$ et une fonction de perte $(g, \hat{d}) \mapsto \ell(g, \hat{d})$, on appelle *risque intégré* de cette stratégie, qu'on pourra noter par exemple $R_{\hat{D}}^{\text{int}}$, l'espérance *a priori* de la perte engendrée par cette stratégie :

$$R_{\hat{D}}^{\text{int}} := \mathbb{E}_{\text{pr}}(\ell(g(Y), \hat{d}(X))). \quad (\text{RC})$$

Définition (RD). Dans le même cadre, la *fonction de risque bayésien* ^[§] associée à la stratégie $\hat{d}(X)$, qu'on pourra noter par exemple $R_{\hat{D}}^{\text{bay}}(\bullet)$, est l'application qui, à une valeur donnée pour l'observation passée, associe la perte moyenne *conditionnellement à cette valeur de l'observation* :

$$R_{\hat{D}}^{\text{bay}}(x) := \mathbb{E}(\ell(g(Y), \hat{d}(x)) \mid X = x). \quad (\text{RE})$$

Définition (RF). Toujours dans le même cadre, le *risque à postériori*, qu'on pourra noter par exemple R_D^{post} , et l'espérance à *postériori* $\mathbb{E}_{\text{post}}(\ell(g(Y), \hat{d}(x_{\checkmark})))$ de la perte engendrée par la stratégie de décision ; en d'autres termes, c'est la valeur $R_D^{\text{bay}}(x_{\checkmark})$ de la fonction de risque bayésien en l'observation effective. ♡

Avec le vocabulaire ci-dessus, le début de cette section peut se reformuler en disant que pour minimiser le risque intégré, il faut chercher à minimiser la fonction de risque bayésien, ce qui peut se faire « *x par x* » : en particulier, trouver la décision optimale pour la valeur effective x_{\checkmark} de l'observation revient à minimiser le risque à postériori.

Dans le cadre de la statistique bayésienne, nous pourrions nous arrêter ici ; cependant, tant que nous y sommes, introduisons aussi le vocabulaire du risque fréquentiste, dont nous nous re-servirons dans le chapitre 11 :

Définition (RG). Toujours dans le même cadre que pour les définitions (RB) à (RF), la *fonction de risque (fréquentiste)*, qu'on pourra noter par exemple $R_D^{\text{fréq}}(\bullet)$, est l'application qui, à une valeur donnée pour le paramètre caché, associe la perte moyenne lorsque le paramètre caché prend cette valeur :

$$R_D^{\text{fréq}}(\theta) := \mathbb{E}_{\theta}(\ell(g(Y), \hat{d}(x))). \quad (\text{RH})$$

On parle parfois aussi de *risque véritable*, qu'on peut noter R_D^{\checkmark} , pour parler de la valeur $R_D^{\text{fréq}}(\theta_{\checkmark})$ prise par la fonction de risque fréquentiste en la vraie valeur du paramètre caché : il s'agit alors de l'espérance de la perte sous la loi véritable. ♡

Remarque (RI). Conformément à ce que son nom suggère, la fonction de risque fréquentiste est un concept fréquentiste : elle ne dépend que du modèle statistique lui-même (plus, bien sûr, de la fonction de perte choisie, et de la décision considérée), pas de la priore! ♣

Remarque (RJ). Noter que les fonctions de risques resp. bayésien et fréquentiste n'ont pas la même domaine : l'argument de la première vit dans l'espace de l'observation, tandis que l'argument de la seconde vit dans l'espace du paramètre caché! ♣

On peut enfin mentionner le lien suivant entre fonction de risque fréquentiste et risque intégré :

Théorème (RK). Toujours avec les mêmes notation que ci-dessus, on a la relation :

$$R_D^{\text{int}} = \int_{\theta \in \Theta} R_D^{\text{fréq}}(\theta) \mathbb{P}_{\text{pr}}(\theta \in d\theta). \quad (\text{RL})$$

◇

On déduit de ce théorème que « trouver la stratégie de décision qui minimise le risque intégré », c'est la même chose que « trouver la stratégie de décision qui minimise la fonction de risque fréquentiste pour tous les θ à la fois ».

[‡]. Sous réserve de considérations techniques d'intégrabilité que je passe sous silence ici.

[§]. Bien préciser « bayésien » !

9.3 Probabilités à postériori

Enjeu de résumer les lois à postériori

Dans les dernières sections de ce chapitre, nous ne nous intéressons pas à la façon de prendre une décision à partir d'une distribution de probabilité à postériori, mais plus simplement à la façon de *communiquer* efficacement ladite distribution de probabilité à postériori aux commanditaires de notre analyse statistique : ce sera à eux ensuite de prendre leurs décisions à partir de celle-ci, que ce soit en utilisant la méthode du chapitre précédent ou des techniques plus “artisanales” !

En effet, décrire une distribution de probabilité à postériori, c'est compliqué, dans la mesure où celle-ci vit généralement dans un espace de dimension infinie... En pratique, l'ingénieur désireux de transmettre le résultat de son analyse statistique optera donc en général pour une réponse synthétique résumant ce résultat en un ou deux nombres. Il existe trois principales façons de “synthétiser” ainsi un résultat, qui constitueront autant de sections de ce chapitre :

- La première idée est de dire que, pour résumer la loi à postériori de notre quantité d'intérêt, on peut tout simplement donner les probabilités à postériori que ladite quantité d'intérêt tombe dans telle ou telle région de l'espace dans lequel elle vit ! Cette idée sera explorée dans la suite de cette section-ci. À vrai dire, il s'agit là d'une notion essentiellement triviale : cependant, nous profiterons de l'étude des probabilités à postériori pour introduire le vocabulaire des *tests*, préparant ce faisant le terrain à l'étude des tests fréquentistes dans la partie III du cours.
- La seconde idée est de dire que, décrire une loi de probabilité, cela revient en première approche à dire dans quelle région de l'espace on doit s'attendre à tomber. Cela nous conduira à la notion d'*intervalle de confiance bayésien*, que nous étudierons dans la § 9.4.
- La troisième idée est d'essayer de compacter au maximum l'information sur la loi à postériori de la quantité d'intérêt en la ramenant à une valeur, censée être “la plus représentative”, pour guider simplement la prise de décision : cela conduira aux notions d'*estimateur* et de *prédicteur* (en l'occurrence, dans le cadre bayésien), qui feront l'objet de la § 9.5.

Les trois façons évoquées ci-dessus de synthétiser la loi à postériori de la quantité d'intérêt ont chacune son utilité : selon le contexte, il sera plus pertinent d'utiliser l'un ou l'autre de ces trois types de résultats. Dans la suite cette section, nous allons commencer par nous focaliser sur la première de ces façons.

Probabilité à postériori d'un évènement

La façon la plus évidente de résumer en quelques nombres réels les calculs de lois à postériori auxquels conduit une analyse bayésienne est tout simplement de dire la probabilité que telle ou telle chose se passe. Ainsi, si $\{\theta \in \Theta_1\}$ (resp. $\{Y \in \mathcal{Y}_1\}$) est une quantité d'intérêt explicative (resp. prédictive) à valeurs booléennes — autrement dit, que c'est un évènement ne dépendant que du paramètre caché (resp. de l'observation future) —, une information pertinente sur cette quantité d'intérêt booléenne est tout simplement sa probabilité à postériori : $\mathbb{P}_{\text{post}}(\theta \in \Theta_1) \stackrel{\text{déf}}{=} \mathbb{P}(\theta \in \Theta_1 \mid X = x_{\checkmark})$, resp. $\mathbb{P}_{\text{post}}(Y \in \mathcal{Y}_1) \stackrel{\text{déf}}{=} \mathbb{P}(Y \in \mathcal{Y}_1 \mid X = x_{\checkmark})$, où x_{\checkmark} est la valeur effective de l'observation.

L'exemple ci-dessous montre sur quel genre de quantités d'intérêt booléennes on peut être amené à se pencher :

Exemple (RM).

- Plaçons-nous dans le modèle du chasseur. Mettons que, aux compétitions de tir, un club de chasse fasse forte impression sur les autres clubs si lorsque son compétiteur atteint au moins 46 cibles, et qu'à l'inverse il se ridiculise lorsque son compétiteur n'atteint pas plus de 18 cibles. Si le club du Bouchonnois se demande « quelle est la probabilité, au vu de ce qu'il a montré lors de son test, que le candidat impressionne les autres clubs lors de la compétition ? (si nous le sélectionnons) », ou « quelle est la probabilité que le candidat nous ridiculise ? », cela revient à regarder les probabilités à postériori des évènements resp. $\{Y \geq 46\}$ et $\{Y \leq 18\}$, lesquels sont des quantités d'intérêt prédictives, puisque ne dépendant que de Y (et à valeurs booléennes, puisqu'étant des évènements).
- Plaçons-nous dans le modèle du pédagogue. Les évènements $\{\mu < \mu_{\text{réf}} - \sigma_{\text{réf}}\}$, $\{\mu \in [\mu_{\text{réf}} - \sigma_{\text{réf}}, \mu_{\text{réf}} - \sigma_{\text{réf}} / 3[]\}$, $\{|\mu - \mu_{\text{réf}}| \leq \sigma_{\text{réf}} / 3\}$, $\{\mu \in]\mu_{\text{réf}} + \sigma_{\text{réf}} / 2, \mu_{\text{réf}} + \sigma_{\text{réf}}]\}$, $\{\mu > \mu_{\text{réf}} + \sigma_{\text{réf}}\}$ sont des quantités d'intérêt (booléennes) explicatives qui expriment informellement les idées respectives que « la nouvelle méthode conduit à des résultats catastrophiques », « la nouvelle méthode conduit à des résultats nettement moins bons que l'ancienne », « ... à peu près comparables à l'ancienne », « ... nettement meilleurs que l'ancienne », ou « ... extraordinaires ». Et du point de vue pratique, il est clair que calculer les probabilités à postériori des ces cinq évènements nous donnera une information intéressante sur ce qu'on peut dire de la nouvelle méthode au vu des résultats de la première promotion. ♣

Les quantités d'intérêt booléennes de type explicatif ont reçu, dans le jargon statistique, un nom plus spécifique : on les appelle des « hypothèses ». Cette nomenclature sera fréquemment ré-utilisée dans la suite, notamment au chapitre 12. Retenons donc :

- ! **Définition (RN)** (Hypothèse). En statistique, une *hypothèse* est une quantité d'intérêt explicative à valeurs booléennes. Ce qu'on peut reformuler de multiples façons :
- En d'autres termes, une hypothèse est un évènement θ -mesurable ^[¶] ;
 - En d'autres termes, une hypothèse est un évènement de la forme $\{\theta \in \Theta_1\}$ pour un certain $\Theta_1 \subseteq \Theta$: pour cette raison, nous parlerons parfois par abus de langage de « l'hypothèse Θ_1 ».
 - En d'autres termes, une hypothèse est une proposition dont la valeur de vérité ne dépend que du paramètre caché θ .

On utilisera souvent la lettre ' \mathcal{H} ' pour désigner les hypothèses. ♥

Exemple (RO).

- Dans le modèle du chasseur, la proposition « le candidat est capable de toucher au moins un plateau sur quatre sur le long terme », autrement dit l'évènement $\{\theta \geq 1/4\}$, est une hypothèse.
- Dans le modèle du pédagogue, la proposition « la dispersion des notes est la même avec la nouvelle méthode qu'avec l'ancienne », autrement dit l'évè-

[¶]. J'entends par là que c'est un évènement mesurable par rapport à la tribu $\sigma(\theta)$ engendrée par θ .

nement $\{\sigma = \sigma_{\text{réf}}\}$ (qu'on pourrait aussi écrire « $\theta \in \mathbb{R} \times \{\sigma_{\text{réf}}\}$ »), est une hypothèse.

- Par contre, dans le modèle du chasseur, la proposition « le chasseur atteindra au moins 19 plateaux lors de la compétition » n'est pas une hypothèse, car elle se réfère à l'observation future, et non pas au paramètre caché ! En revanche, il s'agit bien d'une quantité d'intérêt booléenne (de nature prédictive), et déterminer sa probabilité à postériori serait également tout à fait intéressant en l'occurrence !

♣

Le calcul de la probabilité à postériori d'une hypothèse, ou plus généralement d'une quantité d'intérêt booléenne, revient alors à dire quelle masse la loi à postériori donne à cette quantité d'intérêt :

Exemple (RP). Dans le modèle (explicatif) du chasseur, avec les valeurs standards ($n = 25$, prioré arcsinus, $x_{\checkmark} = 8$), supposons que le club cherche à savoir si le niveau du candidat est suffisant pour toucher au moins un plateau sur quatre sur le long terme : autrement dit, on cherche à connaître la probabilité à postériori de l'hypothèse $\{\theta \geq 1/4\}$. Vu que la loi à postériori de θ est alors une distribution Bêta($8\frac{1}{2}$, $17\frac{1}{2}$) (cf. § 8.4), cette probabilité est tout simplement $\mathbb{P}(\text{Bêta}(8\frac{1}{2}, 17\frac{1}{2}) \geq 1/4)$, soit $1 - \text{répartBêta}(8\frac{1}{2}, 17\frac{1}{2}; (1/4)-)$, ce qu'on peut calculer à l'aide par exemple du logiciel *R* : [III]

```
> 1 - pbeta(1 / 4, 8 + 1 / 2, 17 + 1 / 2)
[1] 0.7944297
```

Ainsi, au vu de son examen, le club estime qu'il a 79,4 % de chances que le candidat ait le niveau pour toucher au moins un plateau sur quatre. (Notons qu'on aurait pu calculer que, à priori, la probabilité qu'il eût ce niveau avait été estimée à 66,7 %).

♣

Profitons au passage de ces considérations sur les probabilités à postériori pour attirer l'attention sur un écueil crucial de la statistique bayésienne :

Proposition (RQ). *Si la probabilité à priori d'un évènement d'intérêt donnée est nulle (resp. égale à 1), alors, quelles que soient les données, sa probabilité à postériori sera nulle (resp. égale à 1) aussi. Par conséquent, lorsqu'on souhaite étudier la probabilité à postériori d'un évènement d'intérêt, pour que notre analyse soit pertinente, il faut absolument qu'on ait donné à cet évènement une probabilité à priori non nulle !*

!!

◇

Démonstration. En fait, de manière générale, si A et B sont des évènements, on a $\mathbb{P}(A) = 0 \implies \mathbb{P}(A | B) = 0$: en effet, si A est de probabilité nulle, $\{A \text{ et } B\}$ est également de probabilité nulle en tant que sous-évènement de A ; et dès lors, $\mathbb{P}(A | B)$, qui par définition vaut $\mathbb{P}(A \text{ et } B) / \mathbb{P}(B)$, est nulle aussi. Par passage au complémentaire, on a de même que $\mathbb{P}(A) = 1 \implies \mathbb{P}(A | B) = 1$.

Quand on y réfléchit bien, cela est très intuitif : si on est *absolument sûr* dès le début d'une certaine hypothèse, rien de ce qui va se passer après ne peut ébranler une telle certitude absolue !

◇

[III]. Les lois bêta étant diffuses, le fait de prendre la fonction de répartition en version continue à gauche ou continue à droite est sans importance ici.

Remarque (RR). Dans la proposition (RQ), la formulation « il faut qu'on ait donné à cet événement une probabilité à priori non nulle » provient de ce que, dans certains cas, on dispose d'une certaine latéralité sur le choix de la priore (confer chap. 15) : dans ce cas, il faut absolument éviter de choisir une priore qui ne permettrait pas d'étudier de façon non triviale la question qui nous intéresse! ♣

Exemple (RS). Dans le cas du pédagogue, par exemple, la priore que nous avons prise en page 82 attribue une probabilité nulle à l'hypothèse $\{\mu = \mu_{\text{réf}}\}$ que la nouvelle méthode soit rigoureusement aussi efficace que l'ancienne. Par conséquent, quels que soient les observations, la probabilité à postériori de l'hypothèse $\{\mu = \mu_{\text{réf}}\}$ sera nulle... Si on souhaite envisager la possibilité que les deux méthodes pédagogiques soient rigoureusement équivalentes en moyenne, cette modélisation ne convient donc pas : il faudrait, par exemple, proposer comme priore un mélange de notre priore initiale et d'une priore réservée à l'hypothèse $\{\mu = \mu_{\text{réf}}\}$. ♣

Vision bayésienne des tests d'hypothèse

Avoir parlé de probabilités à postériori va nous servir de prétexte pour introduire un vocabulaire et des idées que nous retrouverons lorsque nous étudierons la notion de *tests* en statistique fréquentiste (cf. chap. 12). Il sera alors utile de comparer les notions fréquentistes que nous développerons avec leurs pendants bayésiens : c'est pourquoi nous allons maintenant parler des pendants bayésiens en question, quand bien même leur intérêt intrinsèque est assez limité!

Dans la suite de cette sous-section, nous allons plus spécifiquement nous intéresser à une alternative entre deux hypothèses complémentaires l'une à l'autre, auxquelles nous ferons jouer des rôles *différents* dans nos analyses, de sorte que nous leur donnerons aussi des appellations différentes. (La raison d'être de ce vocabulaire n'a essentiellement aucun intérêt dans le cadre bayésien, mais sera en revanche cruciale dans le chapitre 12).

! **Définition (RT).** Lorsqu'une hypothèse \mathcal{H}_0 est qualifiée de « nulle », l'hypothèse complémentaire (autrement dit, l'évènement $\neg\mathcal{H}_0$) est appelée *hypothèse alternative*. Réciproquement, lorsqu'une hypothèse \mathcal{H}_1 est qualifiée d'« alternative », l'hypothèse complémentaire est appelée *hypothèse nulle*.

On utilise conventionnellement les indices respectifs '0' et '1' pour étiqueter les hypothèses resp. nulle et alternative. ♡

Définition (RU) (“*p*-valeur bayésienne”). Deux hypothèses complémentaires \mathcal{H}_0 (nulle) et \mathcal{H}_1 (alternative) étant données, dans le cadre d'une analyse bayésienne, la “*p*-valeur bayésienne” est la probabilité à postériori que l'hypothèse nulle soit vraie, autrement dit que l'hypothèse alternative soit fausse. Il s'agit évidemment de la réalisation d'une certaine statistique bayésienne, statistique qu'on appellera également « la “*p*-valeur bayésienne” » quand on la verra en tant que variable aléatoire (le contexte permettant de trancher). ♡

Remarque (RV). Vous aurez noté que j'ai mis des guillemets autour de l'expression « *p*-valeur bayésienne » : c'est parce que cette expression n'est pas standard dans le vocabulaire de la statistique... En fait, je vous *déconseille* même de l'employer : parlez simplement de « probabilité à postériori que \mathcal{H}_1 soit fausse » le cas échéant!

Néanmoins, l'expression « *p*-valeur booléenne » demeure utile sur le plan *pédagogique*, car elle montre que l'idée de « probabilité à postériori que \mathcal{H}_1 soit fausse » est l'analogue, dans le cadre bayésien, de la notion fréquentiste de « *p*-valeur » (tout

court) que nous verrons au chapitre 13. Or, cette notion de p -valeur (fréquentiste) est assez subtile à saisir : et en pratique, les débutants ont tendance à la confondre avec la fameuse “ p -valeur bayésienne” de la définition (RU) ci-dessus... En donnant explicitement les deux définitions et en soulignant que l’une est le pendant bayésien de l’autre, j’espère attirer votre attention sur la nuance entre ces deux concepts !
 ☺

Voyons maintenant la notion de ce que j’appellerai, dans ce cours, un « test booléen »^[**]. Ce concept sera ré-utilisé dans le cadre fréquentiste ; vous pouvez donc d’ores et déjà le retenir :

Définition (RW) (Test booléen). Un *test (booléen)* est une statistique à valeurs booléennes, autrement dit une fonction (déterministe) de l’observation (passée) qui ne peut valoir que VRAI ou FAUX. !

Un test booléen T est qualifié de « test de l’hypothèse nulle \mathcal{H}_0 » (ou de l’hypothèse alternative \mathcal{H}_1) lorsqu’il “essaye de deviner” qui de \mathcal{H}_0 ou de \mathcal{H}_1 est vraie : dans ce cas, lorsque t_{\checkmark} vaut VRAI, on penche en faveur de l’hypothèse alternative tandis que lorsque t_{\checkmark} vaut FAUX, on penche du côté de l’hypothèse nulle. ♥

Remarque (RX). Notez qu’il faut toujours préciser (éventuellement en le sous-entendant par l’indice ‘0’ ou ‘1’) si l’hypothèse dont on est en train de parler joue le rôle d’hypothèse nulle ou d’hypothèse alternative ! Ainsi on n’écrira jamais « T est un test de l’hypothèse $\{\theta \geq 1/4\}$ » : on écrira soit « T est un test de l’hypothèse nulle $\{\theta \geq 1/4\}$ », soit « T est un test de l’hypothèse alternative $\{\theta \geq 1/4\}$ ». ☺

Un principe très important, qui prendra toute son importance dans le paradigme fréquentiste, est que la notion de « test » sous-entend une *dissymétrie* entre la façon de traiter l’hypothèse alternative et celle de traiter l’hypothèse nulle :

Principe (RY) (Dissymétrie de l’interprétation d’un test booléen). *Lorsqu’on parle de tests booléens, l’usage est de faire jouer des rôles différents à \mathcal{H}_0 et \mathcal{H}_1 ! Le principe étant que, lorsque le test vaudra VRAI, cela signifiera qu’on a un haut niveau de confiance en valeur de l’hypothèse alternative (autrement dit, qu’on la considère comme « très probablement vraie ») ; tandis que lorsque le test vaut FAUX, cela signifie simplement qu’on considère l’hypothèse \mathcal{H}_0 comme plausible, mais sans que cela implique nécessairement qu’elle soit particulièrement probable !* !!

Définition (RZ) (“Test bayésien” de niveau α). Pour $\alpha \in]0, 1[$ (sachant qu’on s’intéressera en général à des cas où $\alpha \leq 1/2$), un test $T = t(X)$ de l’hypothèse nulle \mathcal{H}_0 (ou de l’hypothèse alternative \mathcal{H}_1) est dit de “niveau de risque bayésien” α (ou « de “niveau de confiance” bayésien $1 - \alpha$ ») lorsque la véracité du test entraîne que la probabilité à postériori de l’hypothèse nulle est inférieure ou égale à α , ce qui est équivalent à dire que la probabilité à postériori de l’hypothèse alternative est supérieure ou égale à $1 - \alpha$:

$$\forall x \in \mathcal{X} \quad t(x) \implies \mathbb{P}(\mathcal{H}_0 \mid X = x) \leq \alpha. \quad (\text{SA})$$

♥

Remarque (SB). Notez l’usage des guillemets, confer remarque (RV) : j’ai introduit ce vocabulaire à des fins purement pédagogiques, et je vous déconseille de l’utiliser en pratique ! ☺

[**]. La plupart des références parlent simplement de « test », sans préciser « booléen » : dans ce cours néanmoins, je parlerai aussi de « tests par p -valeur », d’où cette précision.

Remarque (SC). Notez que le fait d'« être un test de “niveau bayésien” α de l'hypothèse alternative \mathcal{H}_1 » ne concerne que ce qui se passe lorsque le test vaut VRAI : à l'inverse, lorsque la réalisation d'un test booléen vaut FAUX, savoir que ce test est de “niveau bayésien” α ne nous apporte aucune information supplémentaire... Cela est cohérent avec le principe (RY), qui veut que de toutes façons, c'est seulement dans le cas où un test booléen vaudra VRAI qu'il fournira une conclusion réellement intéressante.

Cette remarque vaudra également dans le paradigme fréquentiste. \clubsuit

En vertu de la remarque (SC), un test de “niveau bayésien” α d'une hypothèse donnée sera d'autant plus intéressant qu'il répondra souvent VRAI. Il est donc naturel de se demander, parmi tous les tests d'une hypothèse donnée ayant un certain “niveau bayésien”, s'il en existe un qui répond VRAI plus souvent que les autres. C'est effectivement le cas : parmi tous les tests (booléens) de “niveau bayésien” α d'une hypothèse nulle \mathcal{H}_0 donnée, le test optimal^[††] est le test $t(X)$ défini par

$$t(x) := (\mathbb{P}(\mathcal{H}_0 \mid X = x) \leq \alpha). \quad (\text{SD})$$

De cela, on en déduit la connexion suivante entre les “ p -valeurs bayésiennes” et la façon de construire un test booléen de “niveau bayésien” donné. J'ai noté ces points comme importants, car ils continueront d'être valables dans le cadre fréquentiste :

!! **Point (SE).** Ainsi, le “bon” moyen de construire un test booléen de “niveau bayésien” α pour une hypothèse donnée est de répondre en disant si la “ p -valeur bayésienne” de cette hypothèse est, ou pas, inférieure ou égale à α . \clubsuit

! **Proposition (SF).** Réciproquement, si une statistique $p(X)$ est une^[††] “ p -valeur bayésienne” pour une hypothèse donnée, alors le test $\{p(X) \leq \alpha\}$ de cette hypothèse est de “niveau bayésien” α . \diamond

9.4 Intervalles de confiance et de prédiction bayésiens

! **Définition (SG)** (Intervalle de confiance bayésien). Dans le cas où φ (ou g , dans le cas d'une modèle d'inférence prédictive) est à valeurs réelles, un *intervalle de confiance bayésien* au niveau de risque α (ou « au niveau de confiance $1 - \alpha$ »)^[*] est un intervalle réel $I(X)$ ne dépendant que de l'observation (autrement dit, une statistique à valeurs « intervalles ») tel que, pour (presque-)tout $x \in \mathcal{X}$, on ait

$$\mathbb{P}(\varphi(\theta) \notin I(x) \mid X = x) \leq \alpha : \quad (\text{SH})$$

autrement dit, l'intervalle $I(X)$ comprend une proportion au moins $1 - \alpha$ de la loi de $\varphi(\theta)$ conditionnellement à l'observation de X . \heartsuit

[††]. Du point de vue technique, l'« optimalité » prend le sens suivant : pour tout autre test $t'(X)$ de niveau bayésien α de notre hypothèse, on aura que $t'(X) \implies t(X)$ presque-surement (c'est-à-dire que l'ensemble des x pour lesquels on a $t'(x)$ est inclus dans l'ensemble des x pour lesquels on a $t(x)$).

[††]. On dit que $p(X)$ est « une » “ p -valeur bayésienne” pour l'hypothèse nulle \mathcal{H}_0 lorsque (avec les notations génériques), pour tout $x \in \mathcal{X}$, on a $\mathbb{P}(\mathcal{H}_0 \mid X = x) \leq p(x)$. (La réalisation du cas d'égalité correspondant à « la » p -valeur).

[*]. En pratique, comme α est toujours choisi $\leq 1/2$, on ne précise pas si le niveau dont on parle est un niveau de risque ou de confiance : ainsi, un niveau de 5 % désigne implicitement un niveau de risque, et un niveau de 80 % désigne implicitement un niveau de confiance.

Définition (SI) (Intervalle de prédiction bayésien). Dans le cas où on a un modèle de prévision statistique, avec g à valeurs réelles, un *intervalle de prédiction bayésien* au risque α pour $g(Y)$ est une statistique $I(X)$ à valeurs intervalles telle que, pour (presque-)tout $x \in \mathcal{X}$, on ait

$$\mathbb{P}(g(Y) \notin I(x) \mid X = x) \leq \alpha : \quad (\text{SJ})$$

autrement dit, l'intervalle $I(X)$ comprend une proportion au moins $1 - \alpha$ de la loi de $g(Y)$ conditionnellement à l'observation de X . ♡

Remarque (SK). Il est absolument essentiel de bien préciser « bayésien » dans les locutions « intervalle de confiance bayésien » et « intervalle de prédiction » : si vous l'omettez en effet, vous obtenez la locution « intervalle de confiance » (resp. « intervalle de prédiction »), qui se réfère par défaut à un intervalle de confiance (resp. prédiction) *fréquentiste* ^[†] : or la définition des intervalles de confiance (resp. prédiction) fréquentistes est fondamentalement *différente* de celle de leurs contreparties bayésiennes ! La même remarque vaut pour les intervalles de prédiction. ♣

Remarque (SL). Attention, la nomenclature « intervalle de confiance bayésien » peut varier suivant les sources. En particulier, certains auteurs désignent par « intervalle de confiance bayésien » un type d'objet différent (à savoir, un intervalle *construit* selon une procédure bayésienne, mais dont les propriétés requises sont néanmoins *fréquentiste*) : dans ce cours néanmoins nous n'introduirons pas de tels objets, de sorte qu'il n'y aura pas ce risque d'ambiguïté. (Et dès lors, l'appellation « intervalle de confiance bayésien » m'a semblé préférable pour bien montrer le parallélisme entre bayésianisme et fréquentisme). Si vous craignez l'ambiguïté, vous pouvez préférer la locution « intervalle de croyance » pour désigner les intervalles de confiance bayésiens.

En anglais, les intervalles de confiance bayésiens sont souvent appelés « *credible intervals* ». En français, cela est parfois rendu par « intervalles de crédibilité »... Cette dernière traduction est néanmoins **EXTRÊMEMENT** malheureuse, car elle laisse entendre que l'intervalle de confiance bayésien serait l'ensemble des choses qu'il n'est « pas déraisonnable » de croire, alors que ce sont au contraire les intervalles de confiance *fréquentistes* qui correspondent à cette idée-là ! L'intervalle de confiance bayésien, lui, correspond à une description de l'état actuel de notre croyance, pas seulement en tant compte de la crédibilité vis-à-vis de l'observation, mais aussi de nos renseignements à priori reflétés par la priore : ce qui peut conduire certaines hypothèses « crédibles » à être rejetées (par que notre priore les pénalise), ou au contraire certaines hypothèses « peu crédibles » à être acceptées (parce que notre priore les favorise). Je vous interdis donc catégoriquement, dans le cadre de ce cours, d'employer la locution « intervalle de crédibilité » pour parler d'intervalles de confiance bayésien : le cas échéant, cela sera comptabilisé comme une faute ! ♣

Remarque (SM). En fait, rien dans les définitions précédentes n'oblige à ce que $I(X)$ soit un *intervalle* à proprement parler : on pourrait plus généralement définir une *zone* de confiance (resp. prédiction), ce qui aurait l'avantage de se généraliser immédiatement au cas où $\varphi(\bullet)$ est à valeurs dans un espace de dimension > 1 . Cependant les intervalles présentent l'avantage pratique de pouvoir se résumer à deux nombres : c'est pourquoi, dans les usages courants de la statistique, c'est quasiment toujours en termes d'intervalles qu'on exprimera de telles zones ; et c'est donc pourquoi nous nous limiterons à ce cas ici. ♣

Définition (SN) (Intervalle de confiance bayésien, bis). En pratique, on parle souvent de « l' »intervalle de confiance (resp. prédiction) bayésien au niveau α , qui correspond alors à la construction standard de l'intervalle de fluctuation pour la loi à postériori de φ : autrement dit, la réalisation de cet intervalle de confiance (resp.

[†]. Et le risque de confusion est d'autant plus grand que, contrairement au cas des tests (cf. remarque (RV)), la façon de rédiger les phrases de conclusion concernant les intervalles de confiance ou de prédiction est très similaire entre les cas bayésien et fréquentiste !

prédiction) $[\text{Qtile}(\text{Loi}_{\text{post}}(\boldsymbol{\varphi}); \alpha/2), \text{Qtile}(\text{Loi}_{\text{post}}(\boldsymbol{\varphi}); 1 - \alpha/2)]$, où $\text{Qtile}(P; \beta)$ désigne le quantile au niveau β d'une distribution de probabilité P ^[‡]. ♡

Remarque (SO). Il peut aussi arriver qu'on ait besoin d'un *minorant* (resp. *majorant*) bayésien au risque α , à savoir d'une statistique $m(X)$ telle que, pour (presque-)tout $x \in \mathcal{X}$, on ait $\mathbb{P}(\boldsymbol{\varphi}(\boldsymbol{\theta}) < m(x) \mid X = x) \leq \alpha$ (resp. $\mathbb{P}(\boldsymbol{\varphi}(\boldsymbol{\theta}) > m(x) \mid X = x) \leq \alpha$). Il s'agit en fait d'un cas particulier d'intervalle de confiance (resp. prédiction) bayésien où on impose à la borne de droite (resp. de gauche) d'être infinie. Dans ce cas, le minorant (resp. majorant) bayésien optimal au risque α est évidemment $\text{Qtile}(\text{Loi}_{\text{post}}(\boldsymbol{\varphi}); \alpha)$, resp. $\text{Qtile}(\boldsymbol{\varphi}); 1 - \alpha$. ♣

Remarque (SP). Certains ouvrages construisent « l' » intervalle de confiance (resp. prédiction) bayésien au niveau α en faisant en sorte qu'il soit aussi court que possible, ce qui revient à chercher un seuil sur la densité à postériori, aussi haut que possible, tel que la probabilité à postériori de tomber au-dessus de ce seuil soit au moins égale à $1 - \alpha$. Nous ne suivrons pas cette approche qui est plus compliquée à calculer (et qui, par ailleurs, est susceptible de conduire à des zones de confiance (resp. prédiction) qui ne sont *pas* des intervalles); mais il est utile que vous sachiez qu'une telle convention existe si jamais vous la rencontrez. Une troisième convention envisageable, mais plus rare, serait enfin de fixer un seuil sur la *vraisemblance* du paramètre caché (ce qui aurait l'avantage d'être invariant par reparamétrisation, mais ne donnerait pas forcément un intervalle aussi court—et pas toujours un intervalle non plus, d'ailleurs). ♣

Remarque (SQ). Lorsqu'on souhaite arrondir les valeurs numériques calculées pour un intervalle de confiance ou de prédiction (qu'ils soient bayésiens ou fréquentistes, d'ailleurs), il convient toujours de procéder à l'arrondi en *élargissant* l'intervalle, afin pour rester cohérent avec la définition d'intervalle de confiance : on arrondit donc les bornes inférieures par dessous et les bornes supérieures par dessus. ♣

Exemple (SR). Calculons l'intervalle de confiance bayésien au risque 5 % (c'est une valeur que l'usage a consacrée comme standard en l'absence de contexte particulier) pour le cas du chasseur, avec la priore arcsinus, pour $x_{\checkmark} = 8$. Il nous faut donc trouver les quantiles aux niveaux 2,5 % et 97,5 % de la loi à postériori $\text{Bêta}(8\frac{1}{2}, 17\frac{1}{2})$ de $\boldsymbol{\theta}$. Cela peut se calculer sous *R* :

```
> qbeta(0.025, 8 + 1 / 2, 17 + 1 / 2)
[1] 0.1644348
0.1721441
> qbeta(0.975, 8 + 1 / 2, 17 + 1 / 2)
[1] 0.5146123
```

On trouve donc que (la réalisation de) l'intervalle de confiance bayésien au risque 5 % vaut $[0,164, 0,515]$ (concernant les arrondis, confer remarque (SQ) ci-dessus). Autrement dit, au vu de l'épreuve passée par le tireur, il y a 95 % de chances que la fiabilité $\boldsymbol{\theta}$ de celui-ci soit comprise entre 0,172 et 0,515. Ce qui peut paraître étonnamment large : alors que notre tireur a subi un test sur pas moins 25 items et qu'il a touché un chouïa moins d'un plateau sur trois, l'intervalle probable pour la valeur de son niveau réel inclut des niveaux supérieurs à 50 % ! On a là un premier

[‡]. Dans ce cours, j'écrirai systématiquement l'intervalle sous forme fermée. Les pointilleux remarqueront toutefois que dans certains cas, on pourrait en fait ouvrir une ou deux des bornes de l'intervalle : la définition la plus générale consiste en effet à prendre pour I le plus petit intervalle tel que $\mathbb{P}_{\text{post}}(\forall z \in I \boldsymbol{\varphi}(\boldsymbol{\theta}) < z)$, $\mathbb{P}_{\text{post}}(\forall z \in I \boldsymbol{\varphi}(\boldsymbol{\theta}) > z) \leq \alpha/2$. Cela dit, ouvrir les bornes de l'intervalle n'a aucun intérêt du point de vue pratique, puisque cela ne fait rendre l'intervalle plus étroit que d'un *seul* point, autrement dit rien du tout par rapport à la taille de l'intervalle, et même à la précision des bornes...

exemple où l'analyse statistique contredit notre intuition, qui aurait eu tendance à nous faire sur-interpréter le résultat du test de tir passé par le chasseur.

Dans le contexte du club de chasse, on peut aussi imaginer que ce qui intéresse le club, c'est d'avoir une valeur telle qu'on est quasi-sûr que le niveau du tireur est supérieur à cette valeur (pour avoir une garantie sur son niveau), ou au contraire telle qu'on est quasi-sûr que le niveau du tireur est inférieur à cette valeur (dans le contexte où ce qu'on voudrait serait avant tout ne pas risquer de se priver d'une bonne recrue). On pourrait alors chercher les minorant et majorant bayésiens au risque 5 %, qu'on calcule similairement par :

```
> qbeta(0.05, 8 + 1 / 2, 17 + 1 / 2)
[1] 0.186193
> qbeta(0.95, 8 + 1 / 2, 17 + 1 / 2)
[1] 0.483074
```

Cela donne, après arrondi (cf. remarque (SQ)), un minorant de 18,6 %, resp. un majorant de 48,4 %. (Ces valeurs étant aussi les bornes de l'intervalle de confiance bayésien au risque 10 %).

Si on avait plutôt pris la priore uniforme $\text{Unif}^{\text{me}}(0, 1)$ pour θ , on aurait trouvé un intervalle de confiance bayésien au risque 5 % valant [17,2 %, 51,8 %].

Remarque (SS). On voit ici que les résultats obtenus sont très similaires quelle que soit la priore utilisée. C'est là un phénomène général, que nous constaterons également lorsque nous parlerons d'estimation : dès lors qu'on d'un (assez) grand nombre d'observations comme ici, l'importance du choix de la priore devient minime, toutes les priores raisonnables donnant alors des résultats comparables. (Par contre, le choix de la priore s'avère souvent crucial lorsque les données sont moins nombreuses !). \clubsuit

Pour compléter l'exemple ci-dessus, voyons comment évolue l'intervalle de confiance bayésien au risque 10 % (et la médiane de la loi à postériori, qui représente en quelque sorte le "centre" des intervalles de confiance, et dont nous verrons plus loin qu'elle correspond aussi à un *estimateur* pertinent de θ) à mesure que le nombre de plateaux lancés au cours du test passé par notre chasseur progresse, en supposant que le test va même jusqu'à 125 plateaux (autrement dit, ici la valeur de n utilisée pour l'analyse évolue avec le nombre total de plateaux lancés). Ici le niveau réel de notre tireur a été pris à $\theta_{\text{v}} = 40\%$ pour la simulation, et ses résultats successifs ont été les suivants, où θ indique un échec et 1 un succès^[§] :

```

  0 1 0 0 1 1 0 0 0 1 0 0 0 0 1 1 0 1 1 0 0 0 0 0 0 1 0 1 0 1 0 0 0 0
  0 0 ↓
... 1 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 1 1 1 1 0 0 0 0 1 0 1 0 1 0 1 0 0 0
  1 1 ↓
... 1 1 0 0 0 0 1 1 1 1 0 0 1 1 0 0 1 1 0 0 0 1 1 0 0 0 0 1 1 1 0 1 1 0
  0 0 ↓
... 0 0 0 0 0 0 1 1 1 1 0 0 1 1 1 0 1
```

Le tracé de la façon dont le club évalue le niveau du tireur est donné par la figure 9.1. Cette figure nous montre plusieurs points qui satisfont notre bon sens : notamment, après un tir réussi, l'estimateur et les deux bornes de l'intervalle de confiance bayésien pour θ se déplacent vers le haut, alors qu'après un tir raté, c'est l'inverse. On voit aussi que, à mesure qu'on collecte des données sur la performance du chasseur, l'intervalle de confiance bayésien sur θ est de plus en plus étroit, ce qui signifie que notre incertitude sur son niveau réel est de plus en plus faible. Par

[§]. Pour information, chaque ligne comporte 36 essais.

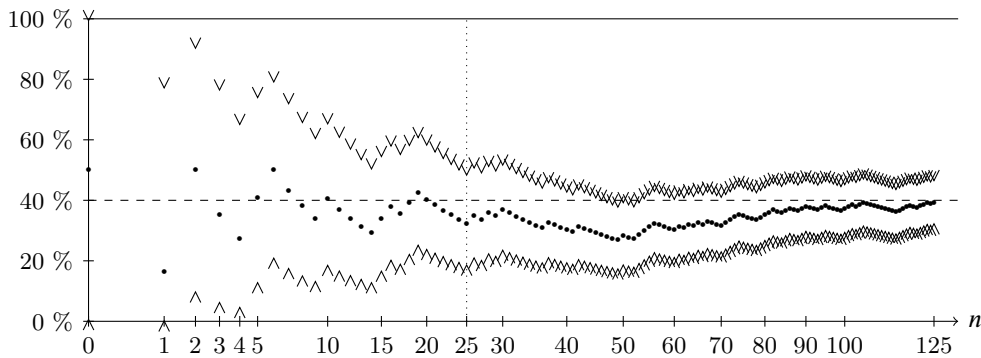


FIGURE 9.1 – Évolution de la médiane de la loi à postérieure de θ (marquée par le symbole \bullet) et de l'intervalle de confiance bayésien (entre les symboles \wedge et ∇) obtenus par la méthode bayésienne à mesure du nombre de tirs tentés par le chasseur. L'axe des abscisses a été gradué irrégulièrement pour faciliter la visualisation. La ligne tiretée horizontale indique le niveau réel de réussite du chasseur utilisé pour la simulation, qui était de 40 % ; la ligne pointillée verticale correspond à $n = 25$, pour lequel, conformément aux valeurs utilisées précédemment, le chasseur a réussi 8 tirs.

ailleurs, sauf sur une courte période, la véritable valeur θ_{v} est bien contenue dans l'intervalle de confiance. On pourrait démontrer que si on continuait à tracer ce graphique à l'infini en faisant tirer le chasseur encore et encore, les deux bornes de l'intervalle de confiance (et à fortiori la médiane à postérieure) convergeraient vers θ_{v} .

Notez qu'on voit clairement que l'intervalle de confiance bayésien ne nous donne qu'une indication solide sur la valeur de θ (fiable à 90 %, très précisément \smile), mais pas de certitude absolue : par exemple, entre $n = 46$ et $n = 53$, il se trouve que, par chance, le tireur a largement sous-performé jusqu'à présent, de sorte que le club tend alors à nettement sous-estimer son niveau réel et que θ_{v} n'est *pas* dans l'intervalle de confiance. Au demeurant, si on devait raisonner en termes de *certitudes*, il est clair que ce qui est *rigoureusement* impossible au vu d'un certain ensemble d'informations ne peut pas redevenir possible avec des informations supplémentaires : ainsi, si on devait raisonner en termes de certitudes, il est clair que l'intervalle où l'on *sait* que θ se trouve [¶] ne pourrait faire que devenir de plus en plus petit [||] à mesure qu'on récolte des données. Or ce n'est pas le cas ici : par exemple, après 4 tirs, l'intervalle de confiance [4,6 %, 65,1 %] ne contient pas la valeur 70 %, alors qu'après le succès au 5^e tir, cette valeur revient à nouveau dans l'intervalle ! (qui devient alors [12,7 %, 74,0 %]). \clubsuit

9.5 Décisions et estimateurs

Introduction aux notions d'estimation et de prédiction

La théorie de la décision optimale peut s'appliquer au cas où on considère une certaine quantité d'intérêt *explicative* $\varphi(\theta) =: \varphi$ et que notre "décision" consiste à

[¶]. Dans le cas du modèle probabiliste considéré, en fait, dans la mesure où on a exclu 0 et 1 se l'espace du paramètre caché, *toute* valeur de θ est *théoriquement* compatible avec l'observation !...

[||]. « De plus en plus petit » au sens ensembliste du terme, s'entend, pas seulement en termes de largeur !

essayer deviner quelle est la valeur de $\boldsymbol{\varphi}$ au moyen d'une statistique $\hat{\boldsymbol{\varphi}}(X) =: \hat{\boldsymbol{\varphi}}$. Dans ce cas, on parle plus spécifiquement d'*estimation* :

Définition (ST) (Estimateur). Un *estimateur* d'une quantité d'intérêt explicative $\boldsymbol{\varphi}(\boldsymbol{\theta}) =: \boldsymbol{\varphi}$ est une statistique $\hat{\boldsymbol{\varphi}}(X) =: \hat{\boldsymbol{\varphi}}$ qui "essaye de ressembler" à $\boldsymbol{\varphi}$. La réalisation $\hat{\boldsymbol{\varphi}}(\mathbf{x}_{\mathcal{J}})$ d'un estimateur est qualifiée d'*estimation* de la quantité d'intérêt $\boldsymbol{\varphi}(\boldsymbol{\theta}_{\mathcal{J}})$.

La convention est d'utiliser le même symbole pour désigner l'estimateur que celui pour désigner la quantité qu'il estime, mais agrémenté d'un chapeau.

L'*estimation statistique* est la branche de la statistique qui s'intéresse à la façon d'obtenir de (bons) estimateurs des quantités d'intérêt explicatives. ♡

Remarque (SU). Dans la définition ci-dessus, la locution « essaye de ressembler » a été mise entre guillemets, car ce n'est pas un concept mathématique : ce qui fait qu'on qualifie $\hat{\boldsymbol{\varphi}}(X)$ d'estimateur de $\boldsymbol{\varphi}(\boldsymbol{\theta})$, c'est qu'on a décidé d'introduire cette quantité avec l'*intention* d'approcher $\boldsymbol{\varphi}(\boldsymbol{\theta})$ par un statistique ! Mais, formellement, n'importe quelle statistique peut être considérée comme un estimateur de n'importe quelle quantité d'intérêt (sous réserve, tout de même, qu'elles soient à valeurs dans le même espace). En fait, c'est au niveau de l'*usage* qu'on va chercher à ne s'intéresser, autant que possible, qu'à des estimateurs ayant certaines "bonnes" propriétés. ♣

Le pendant de la notion d'estimation dans le cadre prédictif est qualifié de *prédiction* :

Définition (SV) (Prédicteur). Un *prédicteur* d'une quantité d'intérêt explicative $g(Y) =: G$ est une statistique $\hat{g}(X) =: \hat{G}$ qui "essaye de ressembler" à G . La réalisation $\hat{g}(\mathbf{x}_{\mathcal{J}})$ d'un prédicteur est qualifiée de *prédiction* de la quantité d'intérêt $g(\mathbf{y}_{\mathcal{J}})$. La *prédiction statistique* est la branche de la statistique qui s'intéresse à la façon d'obtenir de (bons) prédicteurs des quantités d'intérêt prédictives. ♡

Remarque (SW). Dans le cadre bayésien, estimation et prédiction se traiteront exactement de la même façon, *mutatis mutandis*. Si nous avons introduit des définitions clairement séparées pour ces deux concepts, c'est parce que nous verrons que dans le cadre fréquentiste, par contre, il y a des techniques bien spécifiques à l'estimation qui n'ont pas leur pendant dans le cadre prédictif. ♣

Remarque (SX). Dans la mesure où un estimateur (resp. un prédicteur) essaye de ressembler à la quantité d'intérêt qu'il estime (prédit), il va de soit qu'il faudra que les espaces dans lesquels vivent les v.a. $\boldsymbol{\varphi}$ et $\hat{\boldsymbol{\varphi}}$ (resp. G et \hat{G}) soient les mêmes.

Néanmoins, il pourra arriver qu'on se permette, dans certains cas, d'"étendre" l'espace \mathcal{S} dans lequel vit $\boldsymbol{\varphi}$ (ou G) : par exemple, dans certains cas, il peut être intéressant d'estimer un nombre entier par un nombre réel. Cela se fait d'ailleurs plus ou moins dans la vie de tous les jours : si je dis « vu l'épidémie de gastro-entérite, il devrait y avoir environ cinq ou six élèves absents cette semaine », on peut considérer que « cinq ou six » est une locution signifiant en fait « 5,5 »^[**], et donc que j'énonce une prédiction du nombre d'élèves absents (qui sera, évidemment, un entier !) comme valant 5,5, afin de rester assez bien compatible avec tous les scénarios. (C'est la même idée que pour la remarque digressive ?? un peu plus haut). Mais dans tous les cas on peut considérer que quantité d'intérêt et prédiction vivent dans le même espace : pour cela, il me suffit de décréter que la quantité d'intérêt

[**]. En tout cas, si je dis qu'un sac pèse « cinq ou six kilos », je veux bien dire qu'il pèse environ 5,5 kg, et pas qu'il pèse *soit* 5 kg, *soit* 6 kg ! (Même si on peut considérer que, dans un tel cadre, cette locution indique *aussi* l'ordre de grandeur de notre incertitude sur l'estimation du poids).

« nombre d'élèves absents » est à valeurs dans \mathbb{R}_+ : ce qui est techniquement tout à fait correct, même si, *en pratique*, seules des valeurs dans \mathbb{N} pourront effectivement être prises! ♣

Dans le contexte de l'estimation ou de la prédiction, on peut, exactement comme dans le cadre général de la théorie de la décision, introduire une fonction de perte. Cette fonction de perte $\ell(g, \hat{g})$ sera néanmoins un peu particulière, puisque ses deux arguments g et \hat{g} vivront dans le même espace \mathcal{G} . En outre, notre but étant ici que notre estimation (ou prédiction) approche au mieux la quantité d'intérêt, on veut que, pour tout $g \in \mathcal{G}$, l'application $\hat{g} \mapsto \ell(g, \hat{g})$ prenne son maximum pour $\hat{g} = g$. On choisit alors de normaliser "verticalement" la fonction de perte de sorte que $\ell(g, g) = 0$ pour tout $g \in \mathcal{G}$ ^[††]. De ce fait, les pertes seront positives : c'est pourquoi on tend à préférer ce formalisme sur celui de l'utilité, où on se retrouverait à manipuler des quantités *négatives* ∴

Convention (SY). Dans le cadre de l'estimation et de la prédiction, les fonctions de perte $(g, \hat{g}) \mapsto \ell(g, \hat{g})$ ne prennent que des valeurs positives ou nulles, et vérifient $\ell(g, g) = 0$ pour tout g ^[††]. ♡

Remarque (SZ). La convention de normalisation à zéro peut aboutir à des fonctions de pertes assez difficiles à appréhender dans le cadre industriel, que je voudrais essayer d'expliquer ici. Mettons que vous soyez un éditeur et que vous veniez d'accepter le roman d'une nouvelle autrice : vous aimeriez savoir combien de lecteurs seront intéressés par acheter ce roman pour ajuster au mieux votre tirage : l'idée étant que, si vous saviez qu'il y a exactement φ lecteurs intéressés, vous imprimeriez précisément φ exemplaires afin de rentrer au mieux dans vos frais. Nous prenons ici le modèle économique simplifié suivant (où nous assimilons gain financier et utilité) : l'impression de chaque exemplaire nous coûte a , et la vente de chaque exemplaire nous rapporte $b > a$; en outre, tous les lecteurs intéressés achèteront un exemplaire dans la limite des stocks disponibles. Dès lors, si le nombre de lecteurs réels est de φ et que nous l'avons estimé à $\hat{\varphi}$ (et donc que nous avons imprimé $\hat{\varphi}$ exemplaires), nous aurons engagé des frais $a\hat{\varphi}$ pour l'impression et récupérerons $b(\varphi \wedge \hat{\varphi})$ à la vente, ce qui suggère de prendre comme perte $\ell_{\text{brut}}(\varphi, \hat{\varphi}) := a\hat{\varphi} - b(\varphi \wedge \hat{\varphi})$. Néanmoins, on voit que cette perte n'est pas normalisée : lorsque $\hat{\varphi} = \varphi$, $\ell_{\text{brut}}(\varphi, \varphi)$ ne vaut pas du tout zéro... L'idée derrière la normalisation est de prendre

$$\ell_{\text{norm}}(\varphi, \hat{\varphi}) := \ell_{\text{brut}}(\varphi, \hat{\varphi}) - \ell_{\text{brut}}(\varphi, \varphi)^{[*]} : \quad (\text{TA})$$

ainsi, la fonction de perte décrit un *manque à gagner* lorsqu'on croit que le nombre de lecteurs intéressés est $\hat{\varphi}$, par rapport à si on avait su que ce nombre était φ : mais ce n'est pas une perte *au sens propre* : il se peut très bien que $\ell(\varphi, \hat{\varphi})$ soit positif mais qu'en réalité l'entreprise fasse un bénéfice en ayant estimé φ par $\hat{\varphi}$!

Au passage, noter que, ici, notre fonction de perte normalisée a une forme plutôt simple, mais non symétrique : on a $\ell(\varphi, \hat{\varphi}) = (b - a)(\varphi - \hat{\varphi})$ pour $\hat{\varphi} \leq \varphi$ (on aurait pu imprimer $(\varphi - \hat{\varphi})$ livres supplémentaires et faire un bénéfice de $b - a$ sur chacun d'eux) ; et de $\ell(\varphi, \hat{\varphi}) = a(\hat{\varphi} - \varphi)$ pour $\hat{\varphi} \geq \varphi$ (on a gâché la dépense provenant de l'impression de $(\hat{\varphi} - \varphi)$ livres qui ne seront pas achetés). ♣

Remarque (TB). Attention, une fonction de perte n'est pas nécessairement symétrique quand on échange les variables : par exemple, si vous cherchez à estimer

[††]. Cela ne perd pas en généralité, car cette translation de la perte ne changera rien à la comparaison des risques intégrés des différentes décisions, et donc au choix de la décision : confer remarque (HU').

[††]. La plupart des fonctions de pertes seront en outre telles que $\ell(g, \hat{g}) > 0$ dès que $\hat{g} \neq g$, mais il peut néanmoins y avoir des exceptions.

[*]. Attention! La normalisation doit se faire en retranchant $\ell(\varphi, \varphi)$, surtout pas $\ell(\hat{\varphi}, \hat{\varphi})$: en effet, c'est quand on translate la perte par une fonction de φ qu'on laisse le problème physiquement invariant ; par contre, traduire par une fonction de $\hat{\varphi}$ changerait radicalement la modélisation...!

la résistance d'un acier pour construire un pont, estimer la limite d'élasticité de votre acier à 300 MPa alors que la réalité est à 500 MPa sera certes problématique (car vous utiliserez alors plus d'acier que nécessaire, entraînant des surcouts), mais moins qu'estimer que la limite d'élasticité est de 500 MPa alors que la réalité est de 300 MPa, où votre pont risque de se détériorer rapidement après sa construction, voire de s'effondrer...! ☹

Remarque (TC). Dans de nombreux ouvrages de statistique, écrits par des mathématiciens, on ne s'intéresse pas vraiment à déterminer la fonction de perte, puisque le contexte industriel n'est pas présent à l'esprit^[†]. De ce fait, on utilise souvent des fonctions de perte génériques. En particulier, une fonction de perte très courante consiste à prendre $\ell(\varphi_0, \varphi_1) = (\varphi_1 - \varphi_0)^2$, qui a l'avantage d'être strictement convexe et très simple analytiquement : c'est ce qu'on appelle la *fonction de perte quadratique*. Cependant cette fonction de perte n'est pas une panacée dans les situations pratiques : et bien trop souvent, au motif que c'est elle qui est mentionnée dans les ouvrages de référence, on l'utilise dans des contextes où elle est tout-à-fait inadaptée...! ☹

Lorsque la quantité d'intérêt est à valeurs réelles, il y a trois fonctions de perte particulièrement classiques qu'il est utile de connaître :

Définition (TD).

- La fonction de perte $(g, \hat{g}) \mapsto \mathbf{1}_{\hat{g} \neq g}$ est appelée « fonction de perte de Hamming ». Le risque associé (autrement dit, l'espérance de la fonction de perte) est appelé *taux d'erreur*, abrégé en anglais « ER » (*error rate*). Cette fonction de perte peut en fait être définie sur n'importe quel espace, et en général on l'utilise surtout sur des espaces discrets, même s'ils ne se plongent pas naturellement dans \mathbb{R} .
- La fonction de perte $(g, \hat{g}) \mapsto |\hat{g} - g|$ est appelée *erreur absolue*, ou *erreur L¹*. Le risque associé est appelé *erreur absolue moyenne*, abrégé en anglais « MAE » (*mean absolute error*).
- La fonction de perte $(g, \hat{g}) \mapsto (\hat{g} - g)^2$ est appelée *erreur quadratique*. Le risque associé est appelé *erreur quadratique moyenne*, abrégé en anglais « MSE » (*mean squared error*). Pour des raisons d'homogénéité physique, on préfère généralement afficher la *racine carrée* de l'erreur quadratique moyenne, abrégée « RMSD » ou « RMSE » en anglais (*root-mean-square deviation*^[‡], resp. *root mean square error*).

♡

Pour les trois fonctions de perte ci-dessus, la solution du problème d'optimisation particularisé auquel nous conduit le théorème ?? a une forme particulièrement élégante :

Théorème (TE). *Soit G une quantité d'intérêt à valeurs réelles que nous cherchons à estimer ou prédire dans le cadre bayésien.*

- (i) *Si nous utilisons la fonction de perte de Hamming, la réalisation de l'estimateur (ou du prédicteur) optimal sera la valeur de g pour laquelle $\mathbb{P}_{\text{post}}(G = g)$ sera maximal (en supposant ici que G est à valeurs discrètes) : autrement*

[†]. Non tant parce que les mathématiciens se fichent des applications, mais tout simplement parce qu'ils développent des méthodes *générales* : or, il peut y avoir une foultitude de contextes industriels correspondants, avec des fonctions de pertes très différentes !

[‡]. En fait, cette idée de prendre la racine de la moyenne d'un carré est apparentée à la notion d'*écart-type* ; or « écart-type » se dit « standard deviation » en anglais, d'où le fait qu'on utilise régulièrement le sigle « RMSD ».

dit, c'est le mode de la loi à postérieure de G [§] [¶]. (Cela reste en l'occurrence valable pour G à valeurs dans un espace quelconque).

(ii) Si nous utilisons l'erreur absolue comme fonction de perte, la réalisation de l'estimateur/prédicteur optimal sera la médiane de la loi à postérieure de G .

(iii) Si nous utilisons l'erreur quadratique comme fonction de perte, la réalisation de l'estimateur/prédicteur optimal sera l'espérance de la loi à postérieure de G .

◇

Remarque (TF). Si G est à valeurs dans \mathbb{R}^d , l'espérance à postérieure de G peut également être définie : dans ce cas, elle correspond à optimiser l'espérance de la fonction de perte $\ell(\varphi, \hat{\varphi}) := Q(\hat{\varphi} - \varphi)$, pour $Q(\bullet)$ une forme quadratique définie positive [||] quelconque sur \mathbb{R}^d . (Quelle que soit la forme quadratique définie positive choisie, la solution au problème d'optimisation sera l'espérance à postérieure). ♣

! *Remarque (TG).* Ce qui est remarquable, c'est qu'avec cette approche fondée sur la notion de décision optimale, on finit par retomber sur une idée toute simple : pour estimer ou prédire notre quantité d'intérêt G , on va regarder sa distribution à postérieure et essayer de "résumer" celle-ci par une valeur "typique". Les trois indicateurs les plus classiques pour résumer une distribution de probabilité à valeurs réelles sont l'espérance, la médiane et le mode : et ce que montre le théorème ci-dessus, c'est que chacun de ces choix peut en fait être considéré comme optimal par rapport à une fonction de perte appropriée! ☺ ♣

Remarque (TH). À noter que dans le cadre où G est une variable à densité par rapport à la mesure de Lebesgue sur \mathbb{R}^d , on peut également définir son « mode à postérieure » comme le point où la densité à postérieure par rapport à la mesure de Lebesgue atteint son maximum. Il convient néanmoins de prendre garde à ce que, malgré la similarité des noms, le « mode » d'une distribution de probabilité à densité n'a pas les mêmes propriétés d'invariance que le « mode » d'une distribution discrète. ♣

Remarque (TI). Dans le cas d'une distribution à densité, le mode à postérieure peut aussi être vu comme résultant de l'optimisation de l'espérance d'une certaine perte, mais la formalisation est un peu plus subtile. Soit c une fonction continue croissante de \mathbb{R}_+ dans $[0, 1]$ vérifiant $c(0) = 0$ et $c(d) \equiv 1$ pour d suffisamment grand (mais pour le reste, la forme précise de $c(\bullet)$ importe peu) ; et pour $\varepsilon > 0$, considérons la fonction de perte $\ell^{(\varepsilon)}(\varphi, \hat{\varphi}) := c(|\hat{\varphi} - \varphi| / \varepsilon)$, où $|\bullet|$ se réfère à une norme quelconque sur \mathbb{R}^d . Alors l'estimateur ou prédicteur qui optimise le risque associé à la perte $\ell^{(\varepsilon)}$ va tendre, lorsque $\varepsilon \rightarrow 0$, vers le mode de la loi à postérieure de G (qu'on suppose ici exister et être unique). ♣

Exemple (TJ). Voyons ce que donnent les trois estimateurs classiques de θ dans pour le modèle du chasseur, avec les valeurs par défaut (en particulier $x_{\mathcal{V}} = 8$). Dans ce cas, comme nous l'avons vu, la loi à postérieure de θ est la loi Bêta($8\frac{1}{2}, 17\frac{1}{2}$). Les mathématiciens ont établi une formule exacte pour l'espérance, resp. le mode, des lois bêta : pour $\alpha, \beta > 0$, on a

$$\mathbb{E}(\text{Bêta}(\alpha, \beta)) = \frac{\alpha}{\alpha + \beta}; \quad (\text{TK})$$

[§]. On peut dire aussi « le mode à postérieure de G » ; la même remarque vaut aussi pour la médiane, l'espérance, &c.

[¶]. Dans de nombreuses références, le mode à postérieure de G est également appelé (par abus de langage) *maximum à postérieure* de G , ce qu'on abrège communément en « MAP ».

[||]. Une *forme quadratique définie positive* sur \mathbb{R}^d est une application de la forme $x \mapsto \vec{x}^T \mathbf{M} \vec{x}$, où \vec{x} est le vecteur des coordonnées de x dans la base canonique de \mathbb{R}^d est \mathbf{M} une matrice symétrique définie positive. (Pour information, il existe aussi une définition directe ne passant pas par les matrices).

et pour $\alpha, \beta \geq 1$:

$$\arg \max_{x \in]0, 1[} \frac{\mathbb{P}(\text{Bêta}(\alpha, \beta) \in dx)}{\text{vol}_1(dx)} = \frac{\alpha - 1}{\alpha + \beta - 2}. \quad (\text{TL})$$

On en déduit, en l'occurrence, que l'espérance à postériori de θ vaut 0,327, et que son mode à postériori vaut 0,312. Quant à la médiane à postériori, il n'y a pas de formule simple, mais on peut la calculer par exemple par la commande suivante sous *R* :

```
> qbeta(1 / 2, 8 + 1 / 2, 17 + 1 / 2)
[1] 0.3224241
```

Par ailleurs, si, à la place de la priore arcsinus, on avait utilisé la priore uniforme sur $]0, 1[$, on aurait obtenu comme postérieure la loi Bêta(9, 18) ; et les trois estimateurs ci-dessus auraient été remplacés respectivement par 0,333, 0,320 et 0,329.

On voit sur cet exemple qu'il n'y a pas forcément *un* estimateur parfait pour la quantité d'intérêt : selon les hypothèses de modélisation (bayésienne) et le critère de qualité visé (confer fonction de perte), on trouve des estimateurs différents ! Cependant il est aussi rassurant de constater que tous ces estimateurs donnent des résultats très proches, et cohérents avec l'idée intuitive que notre tireur devrait être en mesure de toucher à peu près 8 plateaux sur 25 (puisque c'est la performance qu'il a eue lors du test) : cela était attendu, puisqu'on est en train d'exploiter les mêmes données, et que celles-ci sont suffisamment nombreuses pour qu'on commence à se faire une idée honnête de la vraie valeur de θ ... Confer aussi remarque (SS) *supra*. ♣

Remarque (TM). Il peut aussi être naturel de s'intéresser à des fonctions de perte qui croissent très vite lorsque la distance entre g et \hat{g} augmente, ce qui revient en substance à vouloir contrôler le pire cas possible concernant cette distance : par exemple, on peut instancier cette idée en considérant la fonction de perte $\ell(g, \hat{g}) := |\hat{g} - g|^p$ pour $p \rightarrow \infty$. Dans ce cas encore, (l'asymptotique de) l'estimateur optimal est un indicateur de tendance centrale assez classique : il s'agit du milieu de la plage de valeurs que la loi à postériori de G est susceptible de prendre (plus précisément, c'est $(\text{quantile}(\text{Loi}_{\text{post}}(G); 0+) + \text{quantile}(\text{Loi}_{\text{post}}(G); 1-))/2$), ce qu'on appelle parfois la « mi-gamme » ^[**] de cette loi. Néanmoins, cet indicateur ne présente guère d'intérêt dans les situations qu'on rencontre en pratique, car la valeur centrale de $\text{Loi}(G \mid X = x)$ est très souvent indéfinie et/ou indépendante de x , de sorte qu'elle ne présente guère d'intérêt pour l'analyse statistique ! ♣

9.6 Conclusion sur l'analyse bayésienne

Le contexte bayésien présente l'avantage considérable qu'on a une réponse précise aux questions qu'on se pose sur $\varphi(\theta)$ (ou $g(Y)$), qu'on peut calculer d'une façon presque "mécanique". Les concepts de probabilité à postériori, de test, d'intervalle de croyance et d'estimateur (dont nous retrouverons des avatars dans la partie III) viennent automatiquement accompagnés d'une façon de les calculer qui est optimale dans un certain sens. En outre, on comprend très bien la mécanique du raisonnement ; bref, c'est parfait...

Le souci cependant, c'est que dans de très nombreuses situations, il n'y a pas de manière vraiment *objective*, voire pas même *neutre*, de déterminer la distribution de probabilité à priori pour θ . Prenons l'exemple de notre chasseur : nous avons

[**]. Également appelée : milieu de gamme, valeur centrale, mid-extrême, *midrange* (anglicisme).

dit que nous considérons que son taux de succès suivait à priori une loi Bêta($\frac{1}{2}, \frac{1}{2}$), mais... il n'y a pas de raison que ce soit vraiment le cas! En effet, comme nous le développerons dans la § 15, sauf cas bien particulier, il n'y a pas de façon objective de choisir la priore, pour deux raisons essentielles :

- D'une part, la priore reflète l'état des informations, formelles et informelles, dont on dispose sur θ . Deux individus disposant d'informations différentes, ou même tout simplement d'intuitions différentes sur l'état du monde, risquent donc de se retrouver en désaccord sur le choix de la priore appropriée pour un problème donné, sans qu'il y ait de façon incontestable de des départager...
- D'autre part, même lorsqu'on est d'accord sur les informations dont on dispose, il n'y a pas véritablement de méthode parfaite pour transformer celles-ci en priore : une forme de "tour de main" s'avère alors indispensable, sur lequel deux statisticiens pourront être en désaccord, et ce désaccord peut dans certains cas (en particulier lorsque la dimension de Θ est grande, voir ci-après) conduire à des conclusions fortement différentes...!

Il est certes vrai que, dans un certain nombre de cas, le choix de la priore n'est pas *si* important que cela, car les données dont on dispose sont suffisamment nombreuses pour "écraser" l'impact du choix de la priore : c'est notamment ce que nous avons observé dans les exemples (SR) et (TJ). Néanmoins, d'une part, cela requiert que les données soient suffisamment nombreuses (or dans de nombreuses situations pratiques, on est intrinsèquement limité quant aux données qu'on peut collecter); et d'autre part, ce phénomène ne se produit que très laborieusement (au sens où il faut vraiment *énormément* de données) lorsque la dimension de l'espace du paramètre caché Θ devient grande^[††] ^[‡‡]! Et il y a même des cas où Θ est de dimension infinie, par exemple si le paramètre caché est une fonction continue de \mathbb{R} dans \mathbb{R} ... La problématique du choix de la priore est donc véritablement une épine dans le pied à subir dès lors qu'on fait de la statistique bayésienne, dont il est impossible de se débarrasser complètement...

... En revanche, il y a une chose qui est parfaitement connue *sans* avoir à formuler aucun jugement subjectif concernant la priore : la probabilité que telle ou telle observation se produise selon que le paramètre caché vaille tant ou tant! (P. ex., dans le cas du chasseur, la probabilité, pour un niveau de fiabilité donné, de toucher tant ou tant de plateaux lors du test). C'est pourquoi, dans un grand nombre de situations, on va préférer recourir à des techniques *fréquentistes*, dans lesquelles on se passera complètement de l'information sur la loi à priori!

[††]. La raison en est un peu compliquée à expliquer; mais intuitivement, l'idée est que chaque dimension supplémentaire de l'espace du paramètre caché Θ peut représenter un paramètre sur lequel il y a une ambiguïté quant au choix de priore à prendre, et que quand on cumule les effets d'un grand nombre de paramètres, le grand nombre de choix différents qu'ont pu faire deux statisticiens finit par avoir un effet très important, qui ne pourrait être pallié que par un nombre gigantesque de données...

[‡‡]. Un cas extrêmement important où la dimension de Θ peut devenir *très* grande est celui des *réseaux de neurones* utilisés en intelligence artificielle. L'apprentissage automatique à l'aide d'un tel réseau de neurones consiste en substance à considérer que les données qu'on observe (par exemple, dans le cas d'un modèle de langage, il s'agit des exemples de textes récoltés pour entraîner le modèle) provient d'un certain réglage du réseau de neurones (ou du moins, peut être très bien imité par un tel réglage), dont le paramètre θ_{\vee} , décrivant *l'ensemble* des poids à utiliser pour connecter les différents neurones, est inconnu. Dans cette optique, l'apprentissage automatique peut être vu comme une *inférence statistique* sur θ , où la dimension de l'espace du paramètre caché, pour les modèles les plus avancés comme *GPT*, peut atteindre plusieurs *milliards*...!

Dans la partie suivante de ce polycopié, nous allons présenter les méthodes d'analyses qu'on peut utiliser dans un tel cadre fréquentiste.

Troisième partie

Méthodes fréquentistes

Chapitre 10

Notion de vraisemblance

Lors de la partie de ce cours consacrée à l'approche bayésienne de l'inférence statistique, nous avons mis en évidence un concept essentiel intervenant au cœur du théorème de Bayes : à savoir, la *fonction de vraisemblance* associée l'observation qui a été réalisée (définition (PE)). Cependant, nonobstant son rôle crucial en statistique bayésienne, la vraisemblance est un concept fondamentalement *fréquentiste*, puisque pouvant être définie sans faire référence à quelque choix de priore que ce soit... À vrai dire, la vraisemblance est même un des rares outils qui restent encore à notre disposition dans le paradigme fréquentiste : et pour cette raison, elle joue en analyse fréquentiste un rôle encore plus important que dans le paradigme bayésien ! C'est pourquoi, plutôt que de reléguer le concept de vraisemblance à une simple section du chapitre sur le théorème de Bayes, il m'a semblé approprié de lui consacrer un chapitre à part entière, et de placer ce chapitre au sein de la partie fréquentiste du cours. C'est ce chapitre que vous vous apprêtez à lire : parlons donc, en toute généralité, de la notion de vraisemblance ! ☺

10.1 Définition générale et type d'objet

Définition et notation

Pour commencer, rappelons la définition générale de la vraisemblance, telle qu'elle nous est apparue dans le chapitre 8. Il est à noter que, pour les applications qu'un ingénieur pourra être amené à rencontrer, il n'est pas nécessaire de connaître cette définition générale : je le mentionne donc uniquement par souci de cohérence interne ! ☺

Définition (TN) (Vraisemblance). Considérons un modèle de statistique avec les notations génériques. On appelle « fonction de vraisemblance du paramètre caché au vu de l'observation effectivement réalisée », et on note $\theta \mapsto \mathcal{L}(\theta = \theta \mid X = x_{\mathcal{J}})^{[*]}$ (ou plus simplement $\mathcal{L}(\theta = \theta)$, voire $\mathcal{L}(\theta)$, lorsqu'il n'y a pas d'ambiguïté), toute application de Θ dans \mathbb{R}_+ qui est proportionnelle à l'application

$$\theta \mapsto \mathbb{P}(X \in dx_{\mathcal{J}} \mid \theta = \theta), \quad (\text{T0})$$

où $dx_{\mathcal{J}}$ désigne un voisinage infinitésimal arbitraire de $x_{\mathcal{J}}$, à condition que cette application prenne des valeurs qui ne soient ni infinitésimalement petites, ni infini-

[*]. Le symbole ' \mathcal{L} ', universellement employé, provient de l'anglais *likelihood*.

tésimalement grandes. (Autrement dit, il doit s'agir d'une application qu'on peut évaluer numériquement).

Il s'avère que le choix du voisinage infinitésimal dx_{\checkmark} retenu est sans incidence sur la définition : en effet, si dx_{\checkmark}' est un autre voisinage infinitésimal de x_{\checkmark} , la fonction $\theta \mapsto \mathbb{P}_{\theta}(X \in dx_{\checkmark}')$ sera proportionnelle à la fonction $\theta \mapsto \mathbb{P}_{\theta}(X \in dx_{\checkmark})$, de sorte que l'ensemble des fonctions qui sont proportionnelles à l'une ou à l'autre sera le même! \heartsuit

Remarque (TP). La barre verticale dans la notation de la vraisemblance se lit « sachant que », comme pour les probabilités conditionnelles : ce qu'il y a derrière la barre indique un *contexte* (on écrit notre fonction de vraisemblance *au vu de l'observation* x_{\checkmark}), tandis que ce qu'il y a devant la barre indique *l'objet* sur lequel porte la vraisemblance (la fonction de vraisemblance ayant pour domaine Θ , $\mathcal{L}(\theta | x_{\checkmark})$ est la vraisemblance *de la valeur* θ dans le contexte x_{\checkmark} , pas « la vraisemblance du couple (θ, x_{\checkmark}) » : confer aussi remarque (TS) *infra*).

Néanmoins, malgré la similitude de notation et de sens, la barre verticale intervenant dans l'écriture de la vraisemblance ne constitue pas pour autant un « conditionnement » au sens probabiliste du terme : c'est juste une autre forme de contexte! \heartsuit

Remarque (TQ). Dans ce cours, tous les modèles que nous manipulerons admettront toujours des fonctions de vraisemblance, pour toutes les valeurs possibles de l'observation : des tels modèles sont dits « dominés ». Cependant, il existe aussi certains modèles « pathologiques » pour lesquels la vraisemblance ne peut pas être définie... Heureusement, l'existence de tels modèles « non dominés » n'est pas très gênante en pratique, dans la mesure où ils peuvent toujours être approchés par des modèles dominés : on pourra donc faire si la vraisemblance était toujours définie! \heartsuit

Comme nous l'explique la proposition ci-dessous, toutes les fonctions de vraisemblance se déduisent les unes des autres par proportionnalité ; c'est pourquoi, en pratique, on parle de « la » fonction de vraisemblance : même si, en réalité, il ne s'agit pas d'une fonction à proprement parler mais plutôt d'une *classe d'équivalence* de fonctions pour la relation de proportionnalité, classe d'équivalence que, par abus de langage, on pourra désigner par n'importe laquelle de ses représentants ^[†].

! **Proposition (TR).** *La fonction de vraisemblance est définie de façon unique à constante multiplicative près : si $\theta \mapsto L(\theta)$ est une fonction de vraisemblance de notre modèle (associée à l'observation x_{\checkmark}), alors une fonction $\theta \mapsto \tilde{L}(\theta)$ sera une fonction de vraisemblance si et seulement si il existe une constante $\alpha \in \mathbb{R}_+^*$ (ni infinitésimalement grande, ni infinitésimalement petite) telle que $\tilde{L}(\theta) = \alpha L(\theta)$ pour tout θ .*

Par conséquent, nous parlerons par abus de langage de « la » vraisemblance ; mais il faut bien comprendre que, quand on écrit « $\mathcal{L}(\theta = \theta | X = x_{\checkmark}) = L(\theta)$ », cela signifie en fait « l'application $\theta \mapsto L(\theta)$ est une fonction de vraisemblance associée à l'observation x_{\checkmark} ; et plus généralement, l'ensemble des fonctions de vraisemblance associées à cette observation est donné par les applications de la forme $\theta \mapsto \alpha L(\theta)$ pour $\alpha \in \mathbb{R}_+^$ ».* \heartsuit

[†]. En fait, vous avez déjà rencontré ce principe dans le passé, lorsqu'on vous a parlé des congruences modulo n en arithmétique. En effet, quand on écrit que, « modulo 7, on a $5 + 4 = 2$ », du fait qu'on se place modulo 7, les notations '5', '4' et '2' dans l'égalité se réfèrent en fait aux *classes de congruences* (dans $\mathbb{Z}/7\mathbb{Z}$) des (véritables) entiers resp. 5, 4 et 2 [véritables entiers pour lesquels on a évidemment $5 + 4 \neq 2\dots$], c.à.d. aux ensembles $\{5 + 7k | k \in \mathbb{Z}\}$, etc. Et la notion d'addition *entre classes de congruences*, elle, vérifie bien que la somme de la classe de 5 et de la classe de 4 est égale (*rigoureusement* égale, en tant que classe de congruence) à la classe de 2!

Type d'objet

La définition (TN) et la proposition (TR) ci-dessus nous ont dit que « la vraisemblance associée à l'observation x » est un objet de type « fonction de Θ dans \mathbb{R}_+ , vue à constante multiplicative près ». Ce type d'objet mathématique étant quelque peu inhabituel, nous allons dans cette sous-section souligner les conséquences que la nature de la vraisemblance entraîne sur la façon de la manipuler.

Un premier point important à conserver à l'esprit est la valeur numérique $\mathcal{L}(\theta | x_{\checkmark})$ prise par la vraisemblance en une valeur possible du paramètre caché, dès lors qu'elle est non nulle, n'a aucun sens dans l'absolu. Ainsi une vraisemblance n'est pas petite ou grande *en soi*, mais seulement *en comparaison* avec la vraisemblance associée à une autre valeur du paramètre caché :

Remarque (TS).

- (i) La fonction de vraisemblance n'étant définie qu'à proportionnalité près, la seule chose objective (i.e., indépendante d'un choix de normalisation) qu'on puisse dire sur une valeur spécifique $\mathcal{L}(\theta | x_{\checkmark})$ de la vraisemblance est le fait qu'elle soit nulle ou non nulle.
- (ii) En revanche, on peut donner un sens non ambigu (j'entends par là, indépendant du choix de la constante de proportionnalité) au *rapport* entre les vraisemblances de deux valeurs possibles du paramètre caché associés à une même valeur de l'observation, càd. aux expressions de la forme $\mathcal{L}(\theta_1 | x) / \mathcal{L}(\theta_0 | x)$. ♣

Attention cependant : certes, on peut comparer $\mathcal{L}(\theta_0 | x)$ et $\mathcal{L}(\theta_1 | x)$; mais par contre, il n'y a pas de sens à comparer des vraisemblances associées à deux valeurs *distinctes* de l'observation :

Remarque (TT). Lorsqu'on choisit une normalisation de la fonction de vraisemblance (càd. qu'on choisit un représentant dans la classe des fonctions vues à proportionnalité près), il n'y a pas de moyen canonique de lier les choix de représentants correspondant à deux valeurs *distinctes* de l'observation. Par conséquent, pour $x_1 \neq x_0$, les valeurs prises par des expressions comme « $\mathcal{L}(\theta | x_1) / \mathcal{L}(\theta | x_0)$ » ou « $\mathcal{L}(\theta_1 | x_1) / \mathcal{L}(\theta_1 | x_0)$ » dépendront du choix de normalisation : ces expressions n'ont donc *aucun sens objectif* ! ♣

Bien voir aussi que la vraisemblance se rapporte toujours (quoique parfois indirectement, confer définition ?? *infra*) au *paramètre caché* : c'est un concept fondamentalement *statistique*, qui n'a pas d'analogue en théories des probabilités « pures » ! Ainsi :

Remarque (TU). La vraisemblance est un concept qui concerne uniquement *le paramètre caché*, dans un contexte relatif à l'*observation* : mais à part cela, on ne peut pas parler de « vraisemblance d'une variable aléatoire sachant une autre » en général ! Notamment, il n'y aurait *pas de sens* à écrire quelque chose comme « $\mathcal{L}(Y = y | X = x_{\checkmark})$ » ou « $\mathcal{L}(X = x | \theta = \theta_{\checkmark})$ »...

Un peu plus loin, vous verrez néanmoins une « exception » apparent à cette règle : en effet, les définitions ?? et ?? *infra* étendent la notion de vraisemblance aux *quantités d'intérêt explicatives* [expressions du type « $\mathcal{L}(\gamma(\theta = \gamma_{\bullet} | X = x_{\checkmark}))$ », resp. aux *hypothèses* [« $\mathcal{L}(\theta \in \Theta' | X = x_{\checkmark})$ »]. Cependant, l'exception ne sera qu'apparente : en effet, ces notions de vraisemblance étendue n'auront de sens que dans la mesure où quantités d'intérêt explicatives et hypothèses sont elles-mêmes

définies à partir du paramètre caché ! Il ne s'agira donc pas d'autres instances d'un inexistant concept plus général de soi-disant « vraisemblance d'une v.a. par rapport à une autre », mais simplement de “raccourcis notationnels”, qui exprimeront en fait des propriétés de la fonction de vraisemblance *du paramètre caché* lui-même ! ♣

Nuances entre vraisemblance et conditionnement

Après avoir bien précisé quel type d'objet mathématique est la fonction de vraisemblance, je voudrais aussi souligner les différences subtiles, mais néanmoins importantes, que ce concept présente avec deux autres notions proches : à savoir, la probabilité fréquentiste et la probabilité à postériori.

En effet, à la lecture de la définition générale de la vraisemblance (et ce sera encore plus marqué avec le théorème ?? *infra*), on est tenté de se dire : « ah mais, en fait, “ $\mathcal{L}(\theta = \theta \mid X = x)$ ”, ça veut juste dire la même chose que $\mathbb{P}(X = x \mid \theta = \theta)$ ”, à ceci près qu'on a inversé les expressions qui apparaissaient de part et d'autres de la barre ! ». Mais si tel était le cas, il n'y aurait pas lieu de consacrer un chapitre complet à ce concept... Soulignons donc bien les nuances qu'il y a entre vraisemblance et probabilité fréquentiste :

Remarque (TV) (Différences entre vraisemblance et probabilité fréquentiste).

- (i) Alors que la fonction de masse $\mathbb{P}(X = x \mid \theta = \theta)$ est une fonction de x (qu'on peut considérer pour différentes valeurs de θ), la fonction de vraisemblance $\mathcal{L}(\theta = \theta \mid X = x)$ est au contraire une fonction de θ (qu'on peut considérer pour différentes valeurs de x) !
- (ii) En outre, il n'y a que dans le cas où \mathcal{X} est discret que la loi fréquentiste $\text{Loi}(X \mid \theta = \theta)$ peut être assimilée à sa fonction de masse : dans le cas continu, l'expression « $\mathbb{P}(X = x \mid \theta = \theta)$ » vaut systématiquement zéro, et n'a donc aucun intérêt... En réalité, les lois fréquentistes sont fondamentalement des *mesures* sur \mathcal{X} , de sorte que les seules expressions véritablement pertinentes les concernant sont celles de la forme « $\mathbb{P}(X \in dx \mid \theta = \theta)$ ». Rien de tel pour la vraisemblance, où c'est véritablement une *égalité* qu'on considère lorsqu'on écrit le contexte « $X = x$ » à droite de la barre verticale : il n'y a pas une « mesure de contextes » sur \mathcal{X} , mais bien un contexte différent pour *chaque* valeur possible de l'observation !
- (iii) On notera par ailleurs que l'expression « $\mathbb{P}(X \in dx \mid \theta = \theta)$ » se réfère le plus souvent à une valeur infinitésimalement petite ; alors que la vraisemblance $\mathcal{L}(\theta = \theta \mid X = x)$, elle, est un “véritable” nombre réel...
- (iv) Enfin, n'oublions que, alors que l'expression $\mathbb{P}(X \in dx \mid \theta = \theta)$ se réfère à une valeur *précise*, définie de façon non ambiguë, la fonction de vraisemblance $\theta \mapsto \mathcal{L}(\theta = \theta \mid X = x)$, elle, n'est définie qu'à proportionnalité près... Expliquer le pourquoi du comment de l'importance de cette définition à proportionnalité près serait malheureusement un peu trop ambitieux dans le cadre de ce cours ; mais sachez en tout cas qu'il ne s'agit pas d'un simple *choix* définitionnel qu'on aurait fait “par commodité”, mais bien d'une véritable *contrainte*, nécessaire pour que la vraisemblance puisse se manipuler convenablement ! ♣

Un autre concept qui ressemble fortement à la vraisemblance est celui de *postérieure*. (Et d'ailleurs, en l'occurrence, même la position des expressions de part et d'autre de la barre verticale est identique...). Mais là encore, il ne faut pas confondre

les deux concepts ! Les différences sont même encore plus fondamentales qu'avec les probabilités fréquentistes, comme l'explique la remarque suivante :

Remarque (TW) (Différences entre vraisemblance et postérieure).

- (i) La postérieure est un concept *bayésien* (on a besoin de spécifier une loi à priori pour lui donner un sens), tandis que la vraisemblance est un concept *fréquentiste* (elle est définie indépendamment du choix d'une priore) !
- (ii) Postérieure comme vraisemblance portent toutes les deux sur l'espace du paramètre caché Θ : néanmoins, la postérieure est une *mesure de probabilité*, tandis que la vraisemblance est une *fonction*. En particulier, lorsque Θ est continu (avec une priore diffuse), $\mathbb{P}_{\text{post}}(\Theta = \theta)$ vaudra systématiquement zéro, et même l'expression « $\mathbb{P}_{\text{post}}(\Theta \in d\theta)$ » se référera à une quantité infinitésimalement petite ; tandis que $\mathcal{L}(\Theta = \theta)$ sera (en général) non nulle...
- (iii) Nous verrons ultérieurement comment on peut étendre la vraisemblance pour donner un sens à des expressions de la forme « $\mathcal{L}(\Theta \in \Theta')$ » pour $\Theta' \subseteq \Theta$. Mais attention : du fait que la vraisemblance est une *fonction* et non pas une *mesure* la façon de passer de la vraisemblance "ponctuelle" à la vraisemblance d'une hypothèse n'aura rien à voir avec ce qui se passe pour la postérieure : dans le cas de la postérieure, la probabilité à postériori $\mathbb{P}_{\text{post}}(\Theta \in \Theta')$ s'obtient par *somme* (intégrale), puisqu'elle est égale à $\int_{\theta \in \Theta'} \mathbb{P}_{\text{post}}(\Theta \in d\theta)$; tandis que pour la vraisemblance, la vraisemblance $\mathcal{L}(\Theta \in \Theta')$ d'une hypothèse sera définie comme le *supremum* $\sup_{\theta \in \Theta'} \mathcal{L}(\Theta = \theta)$ des sous-hypothèses ponctuelles qui la composent !
- (iv) En outre, comme nous l'avons déjà signalé en remarque ??, la notion de vraisemblance est fondamentalement définie relativement au paramètre caché θ *lui-même* (par exemple, il n'y a pas de sens à parler de vraisemblance de l'observation future). La postérieure, en revanche, n'est qu'une instance de la notion plus générale de « la loi à postériori », dans le cas particulier où on applique cette notion au paramètre caché : mais on peut tout aussi bien considérer la loi à postériori de l'observation future, ou d'une quantité d'intérêt prédictive !
- (v) Pour finir, même remarque qu'en (TV)-(iv) : contrairement à l'expression « $\mathbb{P}_{\text{post}}(\Theta \in d\theta)$ », qui se réfère à une valeur bien précise, la fonction de vraisemblance $\theta \mapsto \mathcal{L}(\Theta = \theta \mid X = x)$, elle, n'est, rappelons-le, définie qu'à proportionnalité près... \clubsuit

Dans l'item (i) de la remarque ci-dessus, nous avons souligné à nouveau un point déjà mentionné dans l'introduction de ce chapitre : à savoir, que la vraisemblance est un concept fondamentalement fréquentiste. Ce point est tellement important que je voudrais le mettre en avant encore une fois, sous la forme d'une remarque spécifiquement dédiée :

Remarque (TX). Il résulte directement de la définition (TN) que la vraisemblance d'un modèle statistique dépend uniquement de la donnée de la famille de lois $\theta \mapsto \text{Loi}(X \mid \Theta = \theta)$, mais *pas* de la priore $\text{Loi}(\Theta)$: par conséquent, *la vraisemblance est un contexte fréquentiste*, qui conerve tout son sens même en-dehors du paradigme bayésien ! \clubsuit

!!

10.2 Comment calculer la vraisemblance

Comme dit au début de la § 10.1, vous n'êtes pas censés retenir la définition générale (TN) : en pratique, seule la façon de calculer la vraisemblance dans les cas les plus classiques mérite d'être retenue ! Il y a trois résultats à connaître : celui qui donne la vraisemblance dans le cas discret ; celui qui donne la vraisemblance dans le cadre à densité ; et celui qui permet de combiner des vraisemblances partielles dans le cas indépendant.

! **Théorème (TY)** (Vraisemblance dans le cas discret). *On considère les notations génériques. Supposons que l'espace de l'observation \mathcal{X} soit discret. Alors une fonction de vraisemblance est donnée par*

$$\mathcal{L}(\theta = \theta \mid X = x_{\mathcal{J}}) = \mathbb{P}(X = x_{\mathcal{J}} \mid \theta = \theta). \quad (\text{TZ})$$

◇

Démonstration. En effet, dans le cas discret, pour $dx_{\mathcal{J}}$ un voisinage infinitésimal de $x_{\mathcal{J}}$, $\mathbb{P}_{\theta}(X \in dx_{\mathcal{J}})$ est précisément égale à $\mathbb{P}_{\theta}(X = dx_{\mathcal{J}})$. ◇

! **Théorème (UA)** (Vraisemblance dans le cas à densité). *On considère les notations génériques. Supposons que l'espace de l'observation \mathcal{X} soit de la forme \mathbb{R}^n ; et supposons en outre que, pour tout $\theta \in \Theta$, la loi $\text{Loi}_{\theta}(X \mid \theta = \theta)$ possède une densité par rapport à la mesure de Lebesgue vol_n , densité notée $f_{\theta}(\bullet) : \mathcal{X} \rightarrow \mathbb{R}_+$.*

Alors une fonction de vraisemblance du modèle est donnée par

$$\mathcal{L}(\theta = \theta \mid X = x_{\mathcal{J}}) = f_{\theta}(x_{\mathcal{J}}). \quad (\text{UB})$$

◇

Démonstration. La propriété de densité énonce que, pour toute zone infinitésimale dx de \mathcal{X} (au niveau du point x), on a

$$\mathbb{P}(X \in dx \mid \theta = \theta) = f_{\theta}(x) \text{vol}_n(dx). \quad (\text{UC})$$

Par conséquent, si nous appliquons la définition générale (TN) de la vraisemblance à un voisinage infinitésimal $dx_{\mathcal{J}}$ de $x_{\mathcal{J}}$, on trouve que notre fonction de vraisemblance doit être proportionnelle à $\theta \mapsto f_{\theta}(x_{\mathcal{J}}) \text{vol}_n(dx_{\mathcal{J}})$: la fonction $\theta \mapsto f_{\theta}(x_{\mathcal{J}})$ convient donc, puisque le facteur « $\text{vol}_n(dx_{\mathcal{J}})$ » ne dépend pas de θ . [Noter, au passage, que la fonction $\theta \mapsto f_{\theta}(x_{\mathcal{J}}) \text{vol}_n(dx_{\mathcal{J}})$ n'aurait par contre pas été une vraisemblance acceptable, dans la mesure où les valeurs qu'elle prend sont infinitésimalement petites : ce qui est interdit par la définition !]. ◇

! **Théorème (UD)** (Vraisemblance pour un modèle à sous-observations indépendantes). *Soit un modèle statistique (fréquentiste) dont l'espace du paramètre caché est Θ et dont l'espace de l'observation est de la forme $\mathcal{X}_0 \times \mathcal{X}_1 \cdots \times \mathcal{X}_{n-1} =: \mathcal{X}$, l'observation s'écrivant alors $(X_0, X_1, \dots, X_{n-1}) =: \vec{X}$: autrement dit, les X_i sont les « sous-observations » dont l'observation est constituée. Supposons que, dans notre modèle, sous les différents contextes fréquentistes $\mathbb{P}_{\theta}(\bullet)$, les sous-observations soient toujours indépendantes :*

$$\forall \theta \in \Theta \quad \text{Loi}_{\theta}(\vec{X}) = \text{Loi}_{\theta}(X_0) \otimes \text{Loi}_{\theta}(X_1) \otimes \cdots \otimes \text{Loi}_{\theta}(X_{n-1}). \quad (\text{UE})$$

Supposons enfin que pour tout i , dans le modèle $\theta \mapsto \text{Loi}_{\theta}(X_i)$ (d'espace du paramètre caché Θ et d'espace de l'observation \mathcal{X}_i) consistant à “ne regarder que la sous-observation X_i ”, on dispose d'une fonction de vraisemblance associée à $x_{i\mathcal{J}}$, que nous noterons $\mathcal{L}(\theta \mid X_i = x_{i\mathcal{J}})$.

Alors, dans notre modèle de départ (celui où on considère l'ensemble de toutes les sous-observations), une fonction de vraisemblance associée à l'observation $(x_{0\mathcal{V}}, x_{1\mathcal{V}}, \dots, x_{(n-1)\mathcal{V}}) =: \vec{x}$ s'obtient en faisant le produit des vraisemblances "partielles" associées à chaque sous-observation :

$$\mathcal{L}(\theta \mid \vec{X} = \vec{x}_{\mathcal{V}}) = \mathcal{L}(\theta \mid X_0 = x_{0\mathcal{V}}) \times \mathcal{L}(\theta \mid X_1 = x_{1\mathcal{V}}) \times \dots \times \mathcal{L}(\theta \mid X_{n-1} = x_{(n-1)\mathcal{V}}). \quad (\text{UF})$$

◇

Démonstration. Fixons-nous des voisinages infinitésimaux $dx_{0\mathcal{V}}, \dots, dx_{(n-1)\mathcal{V}}$ des différents $x_{i\mathcal{V}}$. Le fait que les $\mathcal{L}(\bullet \mid X_i = dx_{i\mathcal{V}})$ soient des vraisemblances signifie alors qu'on a, pour tout $\theta \in \theta$,

$$\mathbb{P}(X_i \in dx_i \mid \theta = \theta) = \alpha_i(dx_{i\mathcal{V}}) \times \mathcal{L}(\bullet \mid X_i = dx_{i\mathcal{V}}) \quad (\text{UG})$$

pour un certain $\alpha_i(dx_{i\mathcal{V}}) \in \mathbb{R}_+^*$ (possiblement infinitésimalement petit). Mais alors on a

$$\begin{aligned} \mathbb{P}((X_0, \dots, X_{n-1}) \in dx_{0\mathcal{V}} \times \dots \times dx_{(n-1)\mathcal{V}} \mid \theta = \theta) & \stackrel{\text{d\u00e9f}}{=} \\ & \mathbb{P}_\theta(X_0 \in dx_{0\mathcal{V}} \text{ et } \dots \text{ et } X_{n-1} \in dx_{(n-1)\mathcal{V}}) \\ & \stackrel{\text{ind\u00e9p}}{=} \mathbb{P}_\theta(X_0 \in dx_{0\mathcal{V}}) \times \dots \times \mathbb{P}_\theta(X_{n-1} \in dx_{(n-1)\mathcal{V}}) \\ & = \alpha(dx_{0\mathcal{V}}) \mathcal{L}(\theta \mid X_0 = x_{0\mathcal{V}}) \times \dots \times \alpha(dx_{(n-1)\mathcal{V}}) \mathcal{L}(\theta \mid X_0 = x_{(n-1)\mathcal{V}}) \\ & = \prod_{i=0}^{n-1} \alpha(dx_{i\mathcal{V}}) \times \mathcal{L}(\theta \mid X_0 = x_{0\mathcal{V}}) \times \dots \times \mathcal{L}(\theta \mid X_0 = x_{(n-1)\mathcal{V}}). \end{aligned}$$

Cela montre bien que le produit $\mathcal{L}(\theta \mid X_0 = x_{0\mathcal{V}}) \times \dots \times \mathcal{L}(\theta \mid X_0 = x_{(n-1)\mathcal{V}})$ est une fonction de vraisemblance pour le modèle de départ (au vu de l'observation $\vec{x}_{\mathcal{V}}$) : en effet, $dx_{0\mathcal{V}} \times \dots \times dx_{(n-1)\mathcal{V}}$ est bien un voisinage infinitésimal de $\vec{x}_{\mathcal{V}}$; le facteur $\prod_{i=0}^{n-1}$ est bien indépendant de θ ; et les valeurs prises par le produit $\mathcal{L}(\theta \mid X_0 = x_{0\mathcal{V}}) \times \dots \times \mathcal{L}(\theta \mid X_0 = x_{(n-1)\mathcal{V}}$ ne sont effectivement ni infinitésimalement petites, ni infinitésimalement grandes (puisque aucune des $\mathcal{L}(\theta \mid X_i = x_{i\mathcal{V}})$ ne l'est). ◇

Remarque (UH). En pratique, il est très courant de se retrouver dans cette situation où l'observation se décompose en plusieurs sous-observations qui sont indépendantes sous les différents contextes fréquentistes : cela recouvre notamment le cas important des modèles d'échantillonnage, conf. définition (NE). ♣

Remarque (UI). Dans le cas où les « sous-modèles » $\theta \mapsto \text{Loi}_\theta(X_i)$ sont tous discrets, resp. tous à densité, on peut écrire explicitement la fonction de masse (resp. de densité) de la loi de l'observation globale comme le produit des fonctions de masse (resp. de densité) des lois de chaque sous-observation ; et alors on peut calculer la vraisemblance en utilisant uniquement le théorème (TY) (resp. le théorème (UA)). Néanmoins, le théorème (UD) s'avère incontournable si *certaines* des sous-modèles sont discrets tandis que *certaines* autres sont à densité : car dans ce cas, le modèle global n'est ni discret, ni à densité ! ♣

10.3 Extensions du concept de vraisemblance

Vraisemblance d'une hypothèse

L'importance du concept de vraisemblance dans le cadre fréquentiste justifie d'en introduire diverses variantes et extensions.

Bien que, comme nous l'avons souligné en remarque (TU), la définition de la vraisemblance porte fondamentalement sur le paramètre caché *lui-même* (et non pas sur la quantité d'intérêt — y compris en inférence explicative) ; il peut être

utile^[‡] de définir un concept de « vraisemblance d'une hypothèse ». Il s'avère qu'en l'occurrence, la façon appropriée de définir la vraisemblance d'une hypothèse est de prendre le *supremum* de la fonction de vraisemblance sur l'ensemble des valeurs individuelles du paramètre caché composant cette hypothèse :

! **Définition (UJ).** On prend les notations génériques. Pour $\Theta' \subseteq \Theta$, on définit la vraisemblance de l'hypothèse $\{\theta \in \Theta'\}$ ^[§] comme

$$\mathcal{L}(\theta \in \Theta' \mid X = x_{\mathcal{J}}) := \sup_{\theta \in \Theta'} \mathcal{L}(\theta = \theta \mid X = x_{\mathcal{J}}). \quad (\text{UK})$$

♥

Remarque (UL). Notez qu'il s'agit d'une *définition*, pas d'un théorème : jusqu'ici, la notion de « vraisemblance d'une hypothèse » n'avait en effet pas encore de sens... ! Expliquer pourquoi c'est bien un supremum qu'il convient de considérer en l'occurrence, plutôt qu'autre chose, dépasserait malheureusement de cadre pédagogique de ce cours : mais en tout cas, sachez qu'il y a une bonne raison à ce que la vraisemblance d'une hypothèse soit définie ainsi ! ☹

Remarque (UM). De même que pour la vraisemblance du paramètre caché, la vraisemblance d'une hypothèse n'est définie qu'à constante multiplicative près : seul le *rappor* de vraisemblance entre deux hypothèses, pour une *même* valeur de l'observation, est défini de façon univoque ! ♣

On pourra remarquer que, via la notion de vraisemblance pour une hypothèse, on obtient même une notion de fonction de vraisemblance pour les quantités d'intérêt explicatives :

Définition (UN). Pour $\gamma =: \gamma(\theta)$ une quantité d'intérêt explicative à valeurs dans \mathcal{G} , on pourra appeler *fonction de vraisemblance de γ* (pour l'observation $x_{\mathcal{J}}$) l'application :

$$\begin{aligned} \mathcal{G} &\rightarrow \mathbb{R}_+ \\ \gamma_* &\mapsto \mathcal{L}(\gamma = \gamma_* \mid X = x_{\mathcal{J}}) : \end{aligned} \quad (\text{UO})$$

dans cette formule, " $\mathcal{L}(\gamma = \gamma_*)$ " se réfère à la vraisemblance de l'hypothèse $\{\gamma = \gamma_*\}$, elle-même définie via la définition (UJ) : l'évènement $\{\gamma = \gamma_*\}$ étant bien une hypothèse, puisque γ désigne en fait $\gamma(\theta)$, de sorte que $\{\gamma = \gamma_*\}$ se produit si et seulement si θ appartient au sous-ensemble $\{\theta \in \Theta \mid \gamma(\theta) = \gamma_*\}$.

De même que pour la fonction de vraisemblance du paramètre caché lui-même, la fonction de vraisemblance d'une quantité d'intérêt explicative n'est définie qu'à proportionnalité près, et ne permet pas de comparaison entre deux cas relatifs à des valeurs différentes de l'observation passée. ♥

Attention néanmoins à un piège : bien que nous ayons étendu le vocabulaire de la « vraisemblance » pour l'appliquer à des hypothèses ou des quantités d'intérêt explicatives, le théorème de Bayes, lui, ne se généralise pas dans ce cadre :

Remarque (UP). Attention ! La formule de Bayes ne marche pas avec des hypothèses : dans le cadre bayésien, pour $\Theta', \Theta'' \subseteq \Theta$, on aura en général

$$\frac{\mathbb{P}_{\text{post}}(\theta \in \Theta'')}{\mathbb{P}_{\text{post}}(\theta \in \Theta')} \neq \frac{\mathcal{L}(\theta \in \Theta'') \mathbb{P}_{\text{pr}}(\theta \in \Theta')}{\mathcal{L}(\theta \in \Theta') \mathbb{P}_{\text{pr}}(\theta \in \Theta')} ! \quad (\text{UQ})$$

[‡]. La notion de vraisemblance d'une hypothèse sera en particulier utile pour construire des *statistiques de test* : confer chapitre 12 à venir.

[§]. Rappelons ici que n'importe quelle hypothèse peut de mettre sous la forme $\{\theta \in \Theta\}$: confer définition (RN).

Et cela ne fonctionne pas non plus avec les quantités d'intérêt explicatives : en général,

$$\frac{\mathbb{P}_{\text{post}}(\boldsymbol{\gamma} \in d\boldsymbol{\gamma}_1)}{\mathbb{P}_{\text{post}}(\boldsymbol{\gamma} \in d\boldsymbol{\gamma}_0)} \neq \frac{\mathcal{L}(\boldsymbol{\gamma} = \boldsymbol{\gamma}_1) \mathbb{P}_{\text{pr}}(\boldsymbol{\gamma} \in d\boldsymbol{\gamma}_1)}{\mathcal{L}(\boldsymbol{\gamma} \in d\boldsymbol{\gamma}_0) \mathbb{P}_{\text{pr}}(\boldsymbol{\gamma} \in d\boldsymbol{\gamma}_0)}. \quad (\text{UR})$$

♣

Log-vraisemblance

Dans un certain nombre de cas (voir remarque (UW) ci-dessous), plutôt que de manipuler la vraisemblance elle-même, il est préférable de manipuler son *logarithme*. On parle alors de *log-vraisemblance* :

Définition (US) (Log-vraisemblance). Prenons les notations génériques. On appelle *fonction de log-vraisemblance* (associée à l'observation effective $x_{\mathcal{J}}$) la fonction :

$$\begin{aligned} \Theta &\rightarrow \mathbb{R} \sqcup \{-\infty\} \\ \theta &\mapsto \log \mathcal{L}(\theta = \theta \mid X = x_{\mathcal{J}}). \end{aligned} \quad (\text{UT})$$

On parlera plus précisément de ln-vraisemblance, de lg-vraisemblance, etc. pour signifier que le logarithme est considéré en base e , en base 10, etc.

La définition s'étend aussi, *mutatis mutandis*, en une notion de log-vraisemblance pour une hypothèse ou une quantité d'intérêt explicative. ♡

Remarque (UU). Les vraisemblances étant définies à proportionnalité près, les log-vraisemblances sont définies à *translation verticale près* : seule la *différence* entre deux log-vraisemblances (associées à la *même* valeur du paramètre caché) possède une valeur numérique non ambiguë ! ♣

Remarque (UV). On notera que, lors du passage au logarithme, le théorème (UD) (sur la multiplicativité de la vraisemblance dans le cas de sous-observations indépendantes) se traduit par une propriété d'*additivité* concernant les log-vraisemblances dans ce cas. ♣

Remarque (UW). Parler en termes de log-vraisemblance, plutôt que de vraisemblance elle-même, s'avère particulièrement utile dans les cas suivants :

1. Lorsque les jeux de données considérés commencent à être grand, les formules numériques pour la vraisemblance vont facilement fournir des nombres absolument minuscules (comme 10^{-1000} ou gigantesques (comme 10^{1000})... Or, du point de vue de l'implémentation numérique, on est alors confronté à un grave souci, car il y a *dépassement de capacité* de la façon standard dont les ordinateurs représentent les nombres réels : lorsque l'ordinateur trouve un nombre inférieur à 10^{-324} en valeur absolue, il le représente en interne comme valant *exactement* zéro ; tandis que les valeurs supérieures à 10^{309} sont considérées comme infinies [¶]... Travailler en logarithmes évite ces soucis.
2. Plus généralement, les variations de $\mathcal{L}(\theta)$ lorsque θ se déplace dans Θ ont tendance à couvrir facilement plusieurs ordres de grandeur (voire des centaines d'ordre de grandeur) : la log-vraisemblance est alors largement préférable pour comparer efficacement les valeurs — en particulier dans une perspective de visualisation, d'une part, pour pouvoir distinguer les vraisemblances valant 10^{-4} de celles valant 10^{-8} !
3. Puisque, dans les modèles à sous-observations indépendantes (qui sont fréquents en pratique), la log-vraisemblance se comporte de façon additive, cela rend alors les calculs bien plus faciles à exprimer en termes de log-vraisemblances : du point de vue numérique, d'une part, comme évoqué dans

[¶]. Ce qui peut même faire planter la machine !

les items précédents, mais aussi pour les calculs théoriques, par exemple lorsqu'on souhaitera dériver la fonction de (log-)vraisemblance pour trouver en quel point elle est maximale! (confer § ??).

4. Enfin, à un niveau encore plus fondamental, de nombreuses notions de théorie de l'information font apparaître naturellement la *log*-vraisemblance plutôt que la vraisemblance elle-même : c'est en particulier le cas de l'*information de Fisher*, très utile en statistique asymptotique, qui se définit comme la dérivée seconde de la ln-vraisemblance. (Et on pourrait citer aussi les critères de sélection de modèles, comme AIC et BIC — que vous verrez l'année prochaine en analyse de données —, qui se formulent naturellement en termes de log-vraisemblances également).

♣

10.4 Deux propriétés remarquables de la vraisemblance

Pour terminer ce chapitre, je voudrais présenter deux propriétés que je trouve particulièrement intéressantes concernant la vraisemblance : s'il ne sera pas nécessaire de la connaître ou de les utiliser dans le cadre de ce cours, je trouve néanmoins qu'elles montrent particulièrement bien toute la puissance du concept de vraisemblance! ☺

Invariance par reparamétrisation

Proposition (UX). Soit $(\Theta, \mathcal{X}, \theta \mapsto \text{Loi}(X \mid \theta = \theta))$ un modèle statistique, et soit $f: \Xi \rightarrow \Theta$ une application. On peut alors déduire de notre modèle un autre modèle statistique dont l'espace du paramètre caché serait Ξ (le paramètre caché lui-même étant alors noté ξ), ayant également \mathcal{X} pour espace de l'observation, par

$$\text{Loi}(X \mid \xi = \xi) := \text{Loi}(X \mid \theta = f(\xi)) : \quad (\text{UY})$$

cela revient en substance à re-paramétriser le modèle par ξ plutôt que par θ , les paramètres ξ et θ étant liés par la relation $\theta = f(\xi)$. Dans ce cas, la fonction de vraisemblance pour le modèle re-paramétré se déduit immédiatement de la fonction de vraisemblance pour le modèle de base : pour tout $x \in \mathcal{X}$,

$$\mathcal{L}(\xi = \xi \mid X = x) = \mathcal{L}(\theta = f(\xi) \mid X = x). \quad (\text{UZ})$$

◇

Cette proposition, en dépit de son allure technique, découle immédiatement de la définition de la vraisemblance. La stabilité par re-paramétrisation de l'observation, en revanche, est plus inattendue ; c'est la proposition suivante :

Proposition (VA). Soit $(\Theta, \mathcal{X}, \theta \mapsto \text{Loi}(X \mid \theta = \theta))$ un modèle statistique, et soit $f: \mathcal{X} \rightarrow \tilde{\mathcal{X}}$ une application. On peut alors déduire de notre modèle un autre modèle statistique dont l'espace de l'observation serait $\tilde{\mathcal{X}}$ (l'observation elle-même étant alors notée \tilde{X}), ayant également Θ pour espace du paramètre caché, par

$$\text{Loi}(\tilde{X} \mid \theta = \theta) := f_* \text{Loi}(X \mid \theta = \theta) \stackrel{\text{déf}}{=} \text{Loi}(f(X) \mid \theta = \theta) : \quad (\text{VB})$$

cela revient en substance à re-paramétriser l'observation x en \tilde{x} , x et \tilde{x} étant liées par la relation $\tilde{x} = f(x)$. Supposons que f soit injective ; elle admet alors sur $f(\mathcal{X})$ une unique rétraction f^{-1} . Dans ce cas, la fonction de vraisemblance pour le modèle re-paramétré se déduit de la fonction de vraisemblance du modèle de base par la relation :

$$\mathcal{L}(\theta = \theta \mid \tilde{X} = \tilde{x}) = \mathcal{L}(\theta = \theta \mid X = f^{-1}(\tilde{x})), \quad (\text{VC})$$

valable pour tout $\tilde{x} \in f(\mathcal{X})$ [1].

◇

[1]. La condition « $\tilde{x} \in f(\mathcal{X})$ » n'est en réalité pas restrictive, car, quelle que soit la valeur de θ_ν , on aura presque-surement $\tilde{X} \in f(\mathcal{X})$ sous la loi $\mathbb{P}(\bullet \mid \theta = \theta_\nu)$.

Remarque (VD). Attention, l'injectivité de f est essentielle pour que ce théorème soit valide! \clubsuit

Démonstration. Moralement, l'idée de la preuve est de dire que, notant $d\tilde{x}$ un voisinage infinitésimal de \tilde{x} , l'évènement $\{\tilde{X} \in d\tilde{x}\}$ est égal à l'évènement $\{f^{-1}(\tilde{X}) \in f^{-1}(d\tilde{x})\}$, et que l'injectivité de f assure, d'une part que $f^{-1}(\tilde{X})$ suit la loi de X , et d'autre part que $f^{-1}(d\tilde{x})$ est un voisinage infinitésimal de $f^{-1}(\tilde{x})$... \diamond

Non-nullité de la vraisemblance

En toute généralité, une vraisemblance peut être positive *ou nulle*. Des situations de vraisemblance nulle peuvent se rencontrer, par exemple, lorsqu'on a un modèle dans lequel l'erreur est bornée : ainsi, dans le modèle dont le paramètre caché est $(\mu, \epsilon) \in \mathbb{R} \times \mathbb{R}_+^*$ et dont l'observation, à valeur dans \mathbb{R} , suit (sous \mathbb{P}_ν) la loi $\text{Unif}^{\text{me}}(\mu_\nu - \epsilon_\nu, \mu_\nu + \epsilon_\nu)$, alors une fonction devraisemblance est donnée par

$$\mathcal{L}(\mu, \epsilon \mid x_\nu) = \frac{\mathbf{1}_{\mu - \epsilon < x} \mathbf{1}_{\mu + \epsilon > x}}{\epsilon} : \tag{VE}$$

il y a donc tout une partie de l'espace du paramètre caché qui aura une vraisemblance nulle!

On peut se demander si « vraisemblance nulle » implique « valeur *rigoureusement* impossible pour le paramètre caché » (et pas juste « très improbable »). Bien que cela semble assez intuitif, cela n'est pas si évident, après tout. En effet, pensez au paradoxe suivant, qu'on rencontre en théorie des probabilités : chaque fois qu'on tire une loi diffuse (une loi $\text{Unif}^{\text{me}}(0, 1)$, par exemple), la valeur qu'on obtiendra sera un certain nombre réel; et ce nombre réel, en tant que valeur bien spécifique, avait *forcément* une probabilité rigoureusement nulle de sortir lorsqu'on a procédé à notre tirage...

Les vraisemblances, néanmoins, ne présentent pas ce paradoxe : lorsqu'une vraisemblance est nulle, cela veut dire qu'on peut *complètement* exclure la valeur correspondante du paramètre caché. C'est ce qu'énonce le théorème suivant :

Théorème (VF). *Considérons un modèle statistique avec les notations génériques, pour lequel, pour tout $x \in \mathcal{X}$, nous disposons d'une fonction de vraisemblance $\theta \mapsto \mathcal{L}(\theta = \theta \mid X = x)$. (Cette dernière expression sera notée plus simplement « $\mathcal{L}(\theta \mid x)$ » dans la suite de l'énoncé).*

Alors, pour tout $\theta \in \Theta$, \mathbb{P}_θ -presque-surement, on aura une observation pour laquelle θ sera de vraisemblance non nulle :

$$\forall \theta \in \Theta \quad \mathbb{P}^{x \sim \text{Loi}_\theta(X)}(\mathcal{L}(\theta \mid x) > 0) = 1^{[**]}. \tag{VG}$$

\diamond

Démonstration. Si x et θ sont des valeurs pour lesquelles $\mathcal{L}(\theta \mid x) = 0$, alors le théorème de Bayes nous assure que, dans une version bayésienne de notre modèle, si l'observation effective vaut x , la probabilité à postériori que θ vaille θ sera nulle (et ce, peu importe le choix de la priore), ce qu'on peut encore écrire en disant que, dans ce cas, il sera presque-certain à postériori que $\{\theta \neq \theta\}$:

$$\mathcal{L}(\theta \mid x) = 0 \implies \mathbb{P}(\theta \neq \theta \mid X = x) = 1. \tag{VH}$$

Maintenant, pour $\theta \in \Theta$, considérons une version bayésienne de notre modèle statistique où la priore est prise égale à δ_θ : autrement dit, dès la phase à priori, nous avons la *certitude* que $\{\theta = \theta\}$. Il est intuitivement clair que cette certitude “devrait” alors rester variable à postériori. Pour le voir formellement, on va décomposer l'évènement $\{\theta \neq \theta\}$ selon les valeurs possibles de x : on obtient que

$$\mathbb{P}_{\text{pr}}(\theta \neq \theta) = \int_{x \in \mathcal{X}} \mathbb{P}_{\text{pr}}(X \in dx) \mathbb{P}(\theta \neq \theta \mid X = x). \tag{VI}$$

[**]. Le formalisme de cette formule est un peu compliqué à comprendre... En fait, ce que (VG) veut dire, c'est la chose suivante : un véritable contexte probabiliste \mathbb{P}_ν ayant été fixé (correspondant à une vraie valeur θ_ν du paramètre caché), si dans un premier temps on tire une observation x sous la loi $\text{Loi}_{\theta_\nu}(X)$, puis que dans un second temps on regarde ce que vaut, *au point* θ_ν , la fonction de vraisemblance pour le x que nous venons de tirer, alors on obtiendra (\mathbb{P}_ν -presque)-toujours une valeur non nulle.

Appelant $N_\theta := \{x \in \mathcal{X} \mid \mathcal{L}(\theta \mid x) = 0\}$ l'ensemble des observations qui donnent une vraisemblance nulle à θ , on en déduit en particulier que

$$\begin{aligned} \mathbb{P}_{\text{pr}}(\theta \neq \theta) &\geq \int_{x \in N_\theta} \mathbb{P}_{\text{pr}}(X \in dx) \underbrace{\mathbb{P}(\theta \neq \theta \mid X = x)}_{=1} \\ &= \int_{x \in N_\theta} \mathbb{P}_{\text{pr}}(X \in dx) = \mathbb{P}_{\text{pr}}(X \in N_\theta). \quad (\text{VJ}) \end{aligned}$$

Mais si nous prenons $\text{Loi}_{\text{pr}}(\theta) \leftarrow \delta_\theta$, le contexte probabiliste à priori correspond par construction au contexte probabiliste \mathbb{P}_θ ! Le membre de gauche de notre inégalité vaut alors 0, tandis que le membre de droite vaut $\mathbb{P}_\theta(X \in N_\theta)$. On a donc $\mathbb{P}_\theta(X \in N_\theta) \leq 0$, d'où en fait $\mathbb{P}_\theta(X \in N_\theta) = 0$: ce qui est bien le résultat annoncé ! \spadesuit

Chapitre 11

Estimation et prédiction en statistique fréquentiste

11.1 Prolégomènes

Notations pour la statistique fréquentiste

Le présent chapitre s’inscrivant de plain-pied dans le paradigme fréquentiste de la statistique inférentielle, il est utile de rappeler au préalable quelques usages notationnels relatifs à ce paradigme, notamment concernant la notion de « véritable contexte probabiliste » :

Remarque (VK). Dans le cadre fréquentiste, il n’y a plus de notions de contextes probabilistes a priori ou a posteriori : seuls les contextes fréquentistes (i.e. les contextes sachant une valeur donnée du paramètre caché) restent bien définis. Pour cette raison, il n’y aura plus vraiment d’intérêt à utiliser les notations conditionnelles « $\mathbb{P}(\bullet \mid \theta = \theta)$ », « $\mathbb{E}(Z \mid \theta = \theta)$ », etc. : à la place, nous privilégierons les notations à indices « $\mathbb{P}_\theta(\bullet)$ », « $\mathbb{E}_\theta(Z)$ », etc. \clubsuit

Remarque (VL). Rappelons que, parmi les contextes probabilistes fréquentistes, le contexte $\mathbb{P}_{\theta_\mathcal{J}}(\bullet)$ associé à la *vraie* valeur du paramètre caché sera appelé « véritable contexte probabiliste » : on pourra également l’abréger par la notation ‘ $\mathbb{P}_{\mathcal{J}}(\bullet)$ ’^[*]. \clubsuit

La notation ‘ $\mathbb{P}_{\mathcal{J}}$ ’ sera en particulier utilisée pour un “truc” rédactionnel que j’utiliserai très régulièrement dans mes corrigés d’exercices, et qu’il est donc important de connaître :

Point (VM). Souvent nous serons amené à établir une certaine propriété relative au véritable contexte probabiliste $\mathbb{P}_{\mathcal{J}}(\bullet)$, puis à faire remarquer que, dans la mesure où cette propriété a été établie sans utiliser aucune connaissance spécifique sur la valeur $\theta_\mathcal{J}$, cela signifie qu’elle est en fait valable en remplaçant $\theta_\mathcal{J}$ par n’importe quelle valeur $\theta \in \Theta$. Par exemple, pour $\gamma =: \gamma(\theta)$ une certaine quantité d’intérêt, si on parvient à établir, sans rien savoir spécifiquement sur la valeur $\theta_\mathcal{J}$, qu’on a $\mathbb{E}_{\mathcal{J}}(\hat{\gamma}(X)) = \gamma_\mathcal{J}$, cela signifiera en fait qu’on a montré l’égalité $\mathbb{E}_\theta(\hat{\gamma}(X)) = \gamma(\theta)$ ^[†] !

[*]. Attention, le véritable contexte probabiliste $\mathbb{P}_{\mathcal{J}}(\bullet)$ signifie donc « $\mathbb{P}(\bullet \mid \theta = \theta_\mathcal{J})$ » : ce n’est pas la même chose que « $\mathbb{P}(\bullet \mid X = x_\mathcal{J})$ », qui correspondrait pour sa part au contexte à posteriori !

[†]. En utilisant que, puisque la v.a. γ désigne en fait $\gamma(\theta)$, sa réalisation $\gamma_\mathcal{J}$ correspond en l’occurrence à $\gamma(\theta_\mathcal{J})$.

pour tout θ . (Pour un cas concret de mise en application de cette astuce rédactionnelle, voir p.ex. la remarque (VW) *infra*). \clubsuit

Enjeux de l'estimation en contexte fréquentiste

Les notions d'estimation et de prédiction ont déjà été rencontrées dans la § 9.5 du polycopié, où le but était de les investiguer dans le cadre bayésien. Rappelons donc les définitions de ces concepts (définitions (ST) et (SV)) :

! Définition (VN).

- (i) Pour $\theta \mapsto \text{Loi}_\theta(X)$ un modèle statistique explicatif avec les notations génériques, pour $\gamma(\theta) =: \boldsymbol{\gamma}$ une quantité d'intérêt explicative à valeurs dans un espace \mathcal{E} , un *estimateur* de $\gamma(\theta)$ est une statistique $\hat{\gamma}(X) =: \hat{\boldsymbol{\gamma}}$ à valeurs dans \mathcal{E} , visant à "reconstruire" $\gamma(\theta_{\checkmark})$ à partir de l'observation x_{\checkmark} . On parle d'*estimation* (de $\gamma(\theta)$) pour désigner la réalisation d'un estimateur (autrement dit, la valeur $\hat{\gamma}(x_{\checkmark})$).
- (ii) Pour $\theta \mapsto \text{Loi}_\theta(X, Y)$ un modèle statistique prédictif avec les notations génériques, pour $g(Y) =: G$ une quantité d'intérêt prédictive à valeurs dans un espace \mathcal{E} , un *prédicteur* de $g(Y)$ est une statistique $\hat{g}(X) =: \hat{G}$ à valeurs dans \mathcal{E} , visant à "deviner" $g(y_{\checkmark})$ à partir de l'observation passée x_{\checkmark} . On parle de *prédiction* (de $g(Y)$) pour désigner la réalisation d'un estimateur (autrement dit, la valeur $\hat{g}(x_{\checkmark})$). \heartsuit

Dans le paradigme bayésien, les développements sur les estimateurs et prédicteurs se sont avérés particulièrement simples : en effet, on disposait d'une façon "balisée" de construire des estimateurs ou prédicteurs, qui permettait de trouver les statistiques optimales par rapport à des fonctions de perte données. Dans le paradigme fréquentiste, en revanche, la situation va s'avérer sensiblement plus complexe : car les critères de qualité des estimateurs et prédicteurs seront nettement plus flous ; et cela nous amènera en outre à envisager des techniques bien plus diverses pour construire des statistiques pertinentes...

C'est donc à cette approche fréquentiste de l'estimation et de la prédiction que nous allons consacrer la suite de ce chapitre. Nous présenterons tout d'abord des critères pour jauger de la "qualité" d'un estimateur ou d'un prédicteur (§§ ?? et ??) ; puis nous aborderons les grands principes qui peuvent être suivis pour trouver des estimateurs ou des prédicteurs satisfaisant de telles « bonnes » propriétés (§§ ?? à ??).

On notera au passage que, alors qu'il n'y avait aucune différence essentielle quant à la façon de définir et d'utiliser les estimateurs et les prédicteurs dans le paradigme bayésien, il n'en ira pas de même dans le cadre fréquentiste : de nombreuses méthodes de construction, en particulier, concerneront spécifiquement les estimateurs, tandis que les prédicteurs ne pourront le plus souvent être construits qu'indirectement, au travers d'estimateurs ! (confer § ??).

11.2 Fonctions de risque et de biais

Fonction de risque

Dans la § 9.5, nous avons introduit la notion de *risque* pour jauger la qualité d'un estimateur ou d'un prédicteur : le risque étant défini comme l'espérance (sous

tel ou tel contexte probabiliste, en fonction du type de risque considéré) d'une certaine *perte* quantifiant l'inadéquation entre la statistique destinée à approcher la quantité d'intérêt considérée et la quantité d'intérêt elle-même.

En l'occurrence, puisque le présent chapitre s'inscrit dans le paradigme fréquentiste, « risque » appropriée en l'occurrence sera celle de *risque fréquentiste*, qui consiste en une certaine *fonction* définie sur l'espace du paramètre caché. Rappelons la définition de cette fameuse « fonction de risque fréquentiste » :

Définition (V0) (Fonction de risque fréquentiste). Considérons un modèle avec les notations génériques. Soit $g(Y) =: G$ une quantité d'intérêt prédictive à valeurs dans un certain espace \mathcal{G} ; et soit $\ell : \mathcal{G} \times \mathcal{G} \rightarrow \mathbb{R}_+$ une fonction de perte, $\ell(g, \hat{g})$ quantifiant à quel point il est « grave » de prédire que la quantité d'intérêt va valoir \hat{g} lorsqu'elle s'avèrera en réalité valoir g .

Dans ce contexte, pour $\hat{g}(X) =: \hat{G}$ un prédicteur visant à approcher la quantité d'intérêt G , la *fonction de risque fréquentiste* du prédicteur \hat{G} est la fonction $R_{\hat{G}}^{\text{fréq}} : \Theta \rightarrow \mathbb{R}_+$ [on pourra omettre l'exposant « fréq » en l'absence d'ambiguïté] qui donne l'espérance de la perte sous les différents contextes fréquentistes :

$$R_{\hat{G}}^{\text{fréq}}(\theta) := \mathbb{E}_{\theta}(\ell(G, \hat{G})) \stackrel{\text{déf}}{=} \mathbb{E}_{\theta}(\ell(g(Y), \hat{g}(X))). \quad (\text{VP})$$

♡

J'ai donné la définition ci-dessus dans le cadre d'un prédicteur; mais on a la même définition, *mutatis mutandis*, pour un estimateur :

Définition (VQ). Toujours avec les notations génériques, si $\hat{\gamma}(X) =: \hat{\gamma}$ est un estimateur de la quantité d'intérêt explicative $\gamma(\theta) =: \gamma$ (l'espace dans lequel vivent γ et $\hat{\gamma}$ étant noté \mathcal{G}), alors la fonction de risque (fréquentiste) de cet estimateur (par rapport à la perte $\ell(\bullet, \bullet)$ sur $\mathcal{G} \times \mathcal{G}$) est

$$R_{\hat{\gamma}}^{\text{fréq}} : \Theta \rightarrow \mathbb{R}_+ \\ \theta \mapsto \mathbb{E}_{\theta}(\ell(\gamma, \hat{\gamma})) \stackrel{\text{déf}}{=} \mathbb{E}_{\theta}(\ell(\gamma(\theta), \hat{\gamma}(X))). \quad (\text{VR})$$

♡

Exemple (VS). Considérons le modèle du chasseur avec $n = 1$ (!), dans lequel nous souhaitons estimer la fiabilité θ . Dans ce modèle très simple, l'observation X ne peut prendre que deux valeurs, à savoir 0 ou 1) : par conséquent, pour décrire un estimateur, il suffit de préciser les valeurs qu'il prend resp. lorsqu'on observe $\{X = 0\}$ et lorsqu'on observe $\{X = 1\}$. Ci-après, pour $\hat{\theta}_0, \hat{\theta}_1 \in [0, 1]$, j'utiliserai la notation $\hat{\theta}^{\hat{\theta}_0, \hat{\theta}_1}$ pour désigner la statistique qui vaut resp. $\hat{\theta}_0$ lorsque $\{X = 0\}$ et $\hat{\theta}_1$ lorsque $\{X = 1\}$. Je considérerai plus spécifiquement trois estimateurs particuliers : $\hat{\theta}^{0,1}$, qui n'est autre que X elle-même ^[†]; $\hat{\theta}^{1/4, 3/4}$, qui peut aussi s'écrire $(X + 1/2)/2$; et $\hat{\theta}^{1/8, 5/8}$.

Pour jauger de la qualité de tels estimateurs, il me faut choisir une fonction de perte $\ell(\bullet, \bullet)$ sur $]0, 1[\times [0, 1]$, qui dira à quel point je considère comme « grave » d'estimer une véritable fiabilité θ par une estimation $\hat{\theta}$ qui s'en éloigne plus ou moins : en l'occurrence, je choisirai pour perte l'« erreur absolue » $\ell(\theta, \hat{\theta}) =: |\hat{\theta} - \theta|$.

[†]. Dans ce cas, l'estimateur vit dans l'intervalle *fermé* $[0, 1]$, qui est légèrement plus gros que l'espace $]0, 1[$ dans lequel vit la quantité d'intérêt θ . Cela arrive parfois : mais cela ne change fondamentalement rien à la théorie, à part qu'à certains endroits on doit remplacer l'espace \mathcal{G} de la quantité d'intérêt par le sur-espace $\tilde{\mathcal{G}}$ dans lequel vit l'estimateur ! ♡

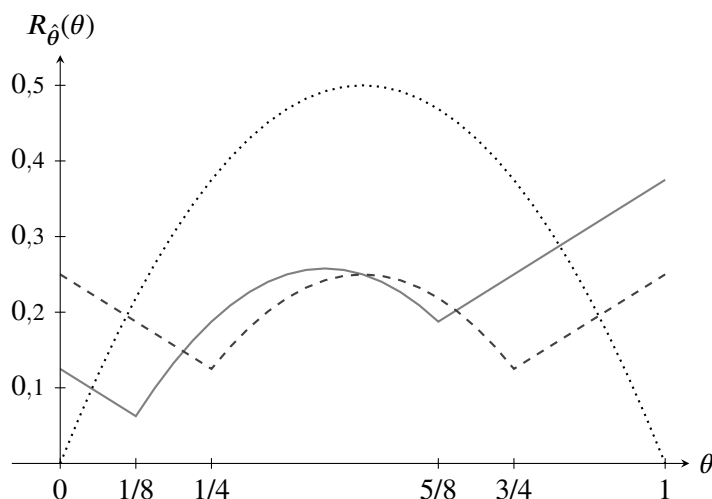


FIGURE 11.1 – Fonctions de risques de trois estimateurs de θ pour le modèle du chasseur avec un seul tir, pour la fonction de perte « erreur absolue ». La courbe en ligne continue correspond à l'estimateur $\hat{\theta}^{1/8, 5/8}$ (qui, par définition, vaut $1/8$ quand le chasseur a manqué son plateau, resp. $5/8$ quand il l'a touché); la courbe tiretée correspond à l'estimateur $\hat{\theta}^{1/4, 3/4}$; et la courbe en pointillés correspond à l'estimateur $\hat{\theta}^{0,1}$.

Cette fonction de perte étant choisie, calculons la fonction de risque de l'estimateur $\hat{\theta}^{\hat{\theta}_0, \hat{\theta}_1}$. Nous allons utiliser l'astuce de rédaction consistant à faire les calculs sous la vraie probabilité $\mathbb{P}_{\mathcal{V}}$, mais sans préciser qui est $\theta_{\mathcal{V}}$, pour sous-entendre implicitement toutes les valeurs possibles de θ (cf. point (VM) dans la § 11.1). Sous la loi $\mathbb{P}_{\mathcal{V}}$, il y a une probabilité $\theta_{\mathcal{V}}$ qu'on observe $X = 1$ et qu'on ait donc $\hat{\theta} = \hat{\theta}_1$, générant dans ce cas une perte de $|\hat{\theta}_1 - \theta_{\mathcal{V}}|$, et une probabilité qu'on observe $X = 0$, générant alors de même une perte de $|\hat{\theta}_0 - \theta_{\mathcal{V}}|$. En sommant ces deux cas, on en conclut que la perte moyenne de notre estimateur sous la vraie loi vaut $\theta_{\mathcal{V}}|\hat{\theta}_1 - \theta_{\mathcal{V}}| + (1 - \theta_{\mathcal{V}})|\hat{\theta}_0 - \theta_{\mathcal{V}}|$; et comme nous avons établi cela sans connaissance spécifique sur la valeur $\theta_{\mathcal{V}}$, cela signifie en fait que la fonction de risque de notre estimateur $\hat{\theta}^{\hat{\theta}_0, \hat{\theta}_1}$ [abrégé en ' $\hat{\theta}$ ' dans la formule ci-dessous] vaut

$$R_{\hat{\theta}}(\theta) = \theta|\hat{\theta}_1 - \theta| + (1 - \theta)|\hat{\theta}_0 - \theta|. \quad (\text{VT})$$

La figure 11.1 montre ce que vaut cette fonction de risque pour les trois estimateurs particuliers évoqués un peu plus haut. ♣

Uniforme méliorité

Maintenant que nous disposons de la notion de fonction de risque pour jauger la qualité d'un estimateur, nous aimerions nous en servir pour *comparer* deux estimateurs (ou prédicteurs). Mais, autant on sait dire, entre deux *nombres* distincts, lequel est le plus petit, autant il n'y a pas d'ordre total sur les *fonctions*...

Une idée qui semblerait naturelle ici serait de *moyenner* les fonctions de risque pour les transformer en un nombre unique, et comparer ensuite les estimateurs via les moyennes de leurs fonctions de risque. Cependant, pour procéder à une telle moyennisation, il faudrait disposer d'une mesure d'intégration sur Θ ... Ce qui est précisément ce que le fréquentisme cherche à éviter, puisqu'il part du postulat qu'on

n'a *aucune* idée de la façon dont il convient de pondérer les différentes zones de Θ entre elles! ^[§]

Nous sommes donc plus ou moins contraints de nous contenter de l'ordre *partiel* usuel entre fonctions pour comparer les fonctions de risque : si un estimateur (ou un prédicteur) a une fonction de risque en tout point inférieure à la fonction de risque d'un autre estimateur (resp. prédicteur), *alors* on pourra dire que le premier estimateur est meilleur que le second... et sinon, eh bien, nos deux estimateurs ne seront pas comparables de façon claire au sens fréquentiste du terme!

Dans le jargon, on dit qu'un estimateur (ou prédicteur) est *uniformément moins risqué* (ou « uniformément meilleur ») qu'un autre lorsqu'on a une telle comparaison entre fonctions de risque :

Définition (VU) (Estimateur ou prédicteur uniformément moins risqué). Soit un modèle statistique explicatif (resp. prédictif) avec les notations génériques et $\gamma(\theta) =: \boldsymbol{\gamma}$ (resp. $g(Y) =: \mathbf{G}$) une quantité d'intérêt à valeurs dans \mathcal{G} ; soient $\hat{\gamma}_0(X)$ et $\hat{\gamma}_1(X)$ deux estimateurs concurrents de $\boldsymbol{\gamma}$ (resp. $\hat{g}_0(X)$ et $\hat{g}_1(X)$ deux prédicteurs...); et soit $\ell(\bullet, \bullet)$ une fonction de perte sur \mathcal{G} . On dit que l'estimateur $\hat{\gamma}_0$ est *uniformément moins risqué* que l'estimateur $\hat{\gamma}_1$ (pour la fonction de perte ℓ) lorsque la fonction de risque du premier est majorée (sur tout Θ) par la fonction de risque du second :

$$\forall \theta \in \Theta \quad R_{\hat{\gamma}_0}^{\text{fréq}}(\theta) \leq R_{\hat{\gamma}_1}^{\text{fréq}}(\theta) \quad (\text{VV})$$

(et on a la même définition, *mutatis mutandis*, pour les prédicteurs). ♡

Remarque (VW). En pratique, quand nous voudrions montrer qu'un estimateur (resp. prédicteur...) est uniformément moins risqué qu'un autre, nous utiliserons l'astuce de rédaction évoquée en remarque (VM) : à savoir, nous étudierons ce qui se passe sous le véritable contexte probabiliste, i.e. au niveau de la vraie valeur du paramètre caché, mais sans rien supposer sur ce que vaut cette valeur au sein de Θ ; et notre but sera alors d'établir, sans avoir à faire aucune hypothèse sur $\theta_{\mathcal{J}}$ (j'insiste!), que

$$\mathbb{E}_{\mathcal{J}}(\ell(\gamma(\theta_{\mathcal{J}}), \hat{\gamma}_0(X))) \leq \mathbb{E}_{\mathcal{J}}(\ell(\gamma(\theta_{\mathcal{J}}), \hat{\gamma}_1(X))). \quad (\text{VX})$$

♣

Exemple (VY). Si nous reprenons l'exemple (VS), parmi les trois estimateurs dont les fonctions de risque sont tracées sur la figure 11.1, on voit aucun n'est uniformément meilleur qu'aucun autre : par exemple, entre les estimateurs $\hat{\theta}^{1/8, 5/8}$ (ligne continue) et $\hat{\theta}^{1/4, 3/4}$ (ligne tiretée), on voit que le premier a un meilleur risque lorsque $\theta < 20\%$ et que $\theta \in]50\%, 65\%[$, mais que c'est le second qui a le risque le plus faible lorsque $\theta \in]21\%, 50\%[$ ou que $\theta > 66\%$... (Et de même, il n'y a pas de comparaison uniforme entre $\hat{\theta}^{1/8, 5/8}$ et $\hat{\theta}^{0, 1}$, ni entre $\hat{\theta}^{1/4, 3/4}$ et $\hat{\theta}^{0, 1}$). ♣

Exemple (VZ). Considérons cette fois-ci le modèle du pédagogue, en version prédictive, où notre objectif est de prédire la variance empirique ^[¶] $\text{var}_{\text{emp}}(Y_0, \dots, Y_{m-1})$

[§]. Et de fait, on pourrait montrer que comparer deux estimateurs (ou prédicteurs) via l'intégrale de leurs fonctions de risques respectives contre une mesure $\mu \in \mathcal{M}(\Theta)$ serait rigoureusement équivalent à comparer les risques intégrés de ces estimateurs dans l'approche bayésienne où on prendrait pour priore sur Θ la loi de probabilité proportionnelle à μ ...

[¶]. Rappelons que, par définition de la notion de « empirique », la variance empirique d'un jeu de données (y_0, \dots, y_{m-1}) est égale à $m^{-1} \sum_{j=0}^{m-1} y_j^2 - (m^{-1} \sum_{j=0}^{m-1} y_j)^2$.

des notes de la seconde promotion. Un prédicteur classique pour cela est le prédicteur de Bessel, qui vaut

$$\hat{G}^B := \frac{(m-1)n}{m(n-1)} \text{var}_{\text{emp}}((X_i)_{i \in \llbracket 0, n \rrbracket}). \quad (\text{WA})$$

Un prédicteur concurrent est le prédicteur dit « rétréci », défini cette fois-ci par

$$\hat{G}^R := \frac{(m-1)n}{m(n+1)} \text{var}_{\text{emp}}((X_i)_{i \in \llbracket 0, n \rrbracket}) \stackrel{\text{déf}}{=} \frac{n-1}{n+1} \hat{G}^B. \quad (\text{WB})$$

Les fonctions de risques de ces deux prédicteurs sont respectivement

$$R_{\hat{G}^B}^{\text{fréq}}(\mu, \sigma) = 2(1 - 1/m)^2 (1/(n-1) + 1/(m-1)) \sigma^4; \quad (\text{WC})$$

$$R_{\hat{G}^R}^{\text{fréq}}(\mu, \sigma) = 2(1 - 1/m)^2 (1/(n+1) + 1/(m-1)) \sigma^4. \quad (\text{WD})$$

On voit ainsi que, pour tout $(\mu, \sigma) \in \Theta$, le risque du prédicteur rétréci est strictement inférieur au risque correspondant du prédicteur de Bessel : le prédicteur rétréci est donc uniformément moins risqué que le prédicteur de Bessel! ♣

Remarque (WE). En pratique, il est assez rare que le critère d'uniforme méliorité permette de trancher entre deux estimateurs ou prédicteurs : notamment, on peut montrer que, sauf cas trivial, il n'existe aucun estimateur (ou prédicteur) qui soit uniformément moins risqué que tous les autres! ♣

Optimalité minimax

Comme nous l'avons vu ci-dessus, la notion d'uniforme méliorité ne permet que rarement de trancher entre deux estimateurs (ou prédicteurs) concurrents, ce qui est assez frustrant. On aimerait donc résumer les fonctions de risques en des nombres réels, qu'on pourrait comparer entre eux plus efficacement... Certes, comme nous l'avons souligné plus haut, il serait inapproprié, si on tient au paradigme fréquentiste, de le faire via l'introduction d'une mesure sur l'espace du paramètre caché. Néanmoins l'intégration n'est pas la seule fonctionnelle permettant de résumer une fonction à un nombre unique : puisqu'on peut aussi, par exemple, considérer l'*infimum* ou le *supremum* de cette fonction : ce qui, pour le coup, peut être défini sans avoir à introduire aucune mesure arbitraire sur le domaine de la fonction! ☺

Pour ce qui est de jauger la qualité d'un estimateur ou d'un prédicteur, prendre l'infimum de la fonction de risque fréquentiste n'aurait pas grand sens, dans la mesure où cela reviendrait à regarder ce que devient le risque si on a beaucoup de chance : alors que, justement, le point est que nous ne savons pas ce que vaut la quantité d'intérêt, et que nous ne voulons pas compter sur la chance pour la deviner! En revanche, le supremum de la fonction de risque présente un véritable intérêt en termes d'ingénierie, puisqu'il nous fournit une *garantie* sur la valeur que le risque véritable est assuré de ne pas dépasser. Ce supremum est appelé, tout simplement, « risque maximal » :

! **Définition (WF).** Avec les notations génériques, on appelle *risque maximal* d'un estimateur $\hat{\gamma}(X)$ (resp. d'un prédicteur $\hat{g}(X)$) le supremum de sa fonction de risque sur Θ . ♣

Exemple (WG). Si nous reprenons l'exemple (VS) (confer en particulier figure 11.1), le risque maximal pour l'estimateur $\hat{\theta}^{0,1}$ vaut $1/2$, atteint en $\theta = 1/2$; le risque maximal pour l'estimateur $\hat{\theta}^{1/8,5/8}$ vaut $3/8$, atteint asymptotiquement pour $\theta \rightarrow 1$, et le risque maximal pour l'estimateur $\hat{\theta}^{1/4,3/4}$ vaut $1/4$, atteint à la fois pour $\theta \rightarrow 0$, $\theta = 1/2$ et $\theta \rightarrow 1$. Au sens du risque maximal, le meilleur de ces trois estimateurs est donc $\hat{\theta}^{1/4,3/4}$, suivi par $\hat{\theta}^{1/8,5/8}$ puis par $\hat{\theta}^{0,1}$.

Et si nous reprenons les prédicteurs de l'exemple (VZ), cette fois-ci, nous voyons qu'aussi bien le prédicteur de Bessel que le prédicteur rétréci ont un risque maximal infini : ces deux prédicteurs sont donc à égalité de ce point de vue-là! ♣

La notion de risque maximal nous donnant un (pré)ordre sur l'ensemble des estimateurs (ou prédicteurs) d'une quantité d'intérêt donnée, on va cette fois-ci pouvoir trouver des estimateurs qui soient *optimaux* de ce point de vue. On parlera alors d'estimateur (resp. prédicteur) « minimax » :

Définition (WH). Étant donné un modèle statistique, une quantité d'intérêt et une fonction de perte, un estimateur (resp. prédicteur) de la quantité d'intérêt est dit *minimax* lorsque son risque maximal est aussi petit que possible, au sens où aucun autre estimateur (resp. prédicteur) de cette quantité d'intérêt, quel qu'il soit, n'a de risque maximal qui soit strictement plus faible. ♡

Exemple (WI). Dans le cadre de l'exemple (VS), on peut montrer que l'estimateur $\hat{\theta}^{1/4,3/4}$ est l'(unique) estimateur minimax de θ (pour la fonction de perte considérée). ♣

Remarque (WJ). Attention néanmoins, il existe d'assez nombreuses situations où aucun estimateur (ou prédicteur) n'a de risque maximal fini [de manière générale, c'est susceptible de se produire dès lors que la fonction de perte n'est pas bornée] : dans ce cas, tous les estimateurs sont minimax, ce qui rend alors la notion sans intérêt...! ∴ C'est notamment le cas pour la problématique considérée dans l'exemple (VZ) (à savoir : dans le modèle du pédagogue, prédiction de la variance empirique de la seconde promotion, avec perte quadratique). ♣

Fonctions de biais

La fonction de risque nous explique, en substance, *de combien* un estimateur ou un prédicteur s'écarte de sa valeur cible. Dans le cas où la cible en question vit dans \mathbb{R} , il y a deux directions opposées dans lesquelles on peut se tromper : “vers le haut” ou “vers le bas”. Il est alors naturel de se demander si l'erreur de l'estimateur (resp. prédicteur) est “équilibrée” entre ces deux façons de se tromper, ou s'il existe un *biais* systématique vers le haut ou vers le bas : dans ce dernier cas, en effet, cela pourrait suggérer une façon de “corriger” l'estimateur en le “tirant” vers le bas ou vers le haut selon le cas! ∩

L'objet qui décrit la déviation moyenne dans un tel cas est ce qu'on appelle la *fonction de biais* :

Définition (WK). Soit un modèle statistique avec les notations génériques, une quantité d'intérêt à valeurs réelles, et un estimateur (ou prédicteur) de cette quantité d'intérêt. La *fonction de biais* associée à cet estimateur (resp. prédicteur) est alors la fonction de Θ dans \mathbb{R} qui nous donne, en fonction de θ , l'erreur *signée* moyenne commise par l'estimateur sous le contexte fréquentiste $\mathbb{P}_\theta(\bullet)$. ♡

Autrement dit, pour un estimateur $\hat{\gamma}(X)$ de $\gamma(\theta)$, c'est

$$\theta \mapsto \mathbb{E}_\theta(\hat{\gamma}(X)) - \gamma(\theta), \quad (\text{WL})$$

et pour un prédicteur $\hat{g}(X)$ de $g(Y)$, c'est

$$\theta \mapsto \mathbb{E}_\theta(\hat{g}(X) - g(Y)) \quad (= \mathbb{E}_\theta(\hat{g}(X)) - \mathbb{E}_\theta(g(Y))). \quad (\text{WM})$$

♡

Selon le comportement de la fonction de biais, on dit que les estimateurs (resp. prédicteurs) sont *biaisés* (ou pas) de telle ou telle façon :

! **Définition (WN).** Nous utilisons les notions génériques. Un estimateur ou prédicteur d'une quantité d'intérêt réelle est dit :

- *Sans biais* lorsque sa fonction de biais est identiquement nulle.
- *Biaisé* lorsqu'il n'est pas sans biais, autrement dit s'il existe au moins une valeur de θ pour laquelle le biais est non nul. *Biaisé vers le haut* (resp. *vers le bas*), lorsqu'il est biaisé et que sa fonction de biais est partout positive ou nulle (resp. négative ou nulle).

♡

Remarque (W0). Il est tout à fait possible qu'un estimateur biaisé ne soit ni biaisé vers le bas, ni biaisé vers le haut ; en fait, c'est même une situation assez courante : pour certaines valeurs de θ , on aura alors $\mathbb{E}_\theta(\hat{\gamma}(X)) < \gamma(\theta)$, mais pour d'autres, on aura l'inégalité inverse. ♣

Exemple (WP). Pour le modèle du chasseur, nous avons considéré dans l'exemple (TJ) l'estimateur de θ défini par $\hat{\theta}^B := (X + 1/2) / (n + 1)$, que nous avons introduit comme l'espérance à posteriori de θ lorsqu'on prend la priore arcsinus. Nous allons maintenant "oublier" les considérations bayésiennes impliquées dans la dérivation de cet estimateur et analyser les propriétés de biais ce dernier dans une perspective purement fréquentiste.

Il est facile de vérifier que l'espérance (sous le véritable contexte probabiliste) de l'estimateur $\hat{\theta}^B$ vaut $(\theta_{\mathcal{J}} + 1/2) / (n + 1)$, qui n'est pas égal à $\theta_{\mathcal{J}}$ en général : ainsi cet estimateur est biaisé. On observe en outre que l'espérance de notre estimateur sera strictement supérieure à $\theta_{\mathcal{J}}$ dans le cas où $\theta_{\mathcal{J}} < 1/2$, resp. strictement inférieure à $\theta_{\mathcal{J}}$ dans le cas où $\theta_{\mathcal{J}} > 1/2$: ainsi l'estimateur n'est ni biaisé vers le bas ni biaisé vers le haut ; par contre, on pourrait dire à la rigueur qu'il est "biaisé vers 1/2".

Un autre estimateur possible de θ est l'estimateur "naïf" $\hat{\theta}^{\text{naïf}} := X / n$. Cet estimateur-là est sans biais, puisque sous la loi $\mathbb{P}_{\mathcal{J}}$, l'espérance de $\hat{\theta}^{\text{naïf}}$ vaut, par linéarité, n^{-1} fois l'espérance d'une loi Bernoulli($n, \theta_{\mathcal{J}}$), ce qui vaut bien $\theta_{\mathcal{J}}$ (et que nous avons établi ce résultat sans aucune connaissance spécifique sur la valeur $\theta_{\mathcal{J}}$). ♣

Remarque (WQ). Attention, la notion d'« être sans biais » n'est pas stable par mesure-image par une fonction continue. Par exemple, nous avons vu ci-dessus que, pour le modèle du chasseur, $\hat{\theta}^{\text{naïf}}$ était un estimateur sans biais de θ ; néanmoins, l'estimateur de $2\theta - \theta^2$ défini par $2\hat{\theta}^{\text{naïf}} - (\hat{\theta}^{\text{naïf}})^2$ est quant à lui biaisé vers le bas (son espérance sous la vraie loi vaut $2\theta_{\mathcal{J}} - \theta_{\mathcal{J}}^2 - n^{-1}\theta_{\mathcal{J}}(1 - \theta_{\mathcal{J}})$). ♣

Remarque (WR). Attention, un estimateur (ou prédicteur) peut être sans biais tout en étant de qualité excécrable... Par exemple, dans le modèle du pédagogue, l'estimateur de μ fourni par X_0 (autrement dit, on regarde juste la note de la première copie corrigée) peut facilement être montré comme étant sans biais; pour autant, on se doute bien qu'un tel estimateur ne sera pas satisfaisant du tout [1]! ♣

Remarque (WS). Dans la littérature classique sur la statistique fréquentiste, beaucoup d'auteurs "s'excitent" sur le fait d'être sans biais, comme si cela était un critère de qualité essentiel des estimateurs et prédicteurs. Mais cela est doublement exagéré :

- D'une part, parce qu'il peut y avoir des estimateurs ou prédicteurs biaisés qui sont en fait *meilleurs* que les estimateurs "corrigés" qu'on obtiendrait en les dé-biaisant. C'est par exemple le cas pour les prédicteurs de l'exemple (VZ) : en effet, l'estimateur de Bessel est sans biais; le prédicteur rétréci, qui lui est proportionnel, est biaisé vers le bas; et pourtant le second prédicteur est uniformément moins risqué (en risque quadratique) que le premier!
- D'autre part, parce qu'il est très fréquent de rencontrer des quantités d'intérêt pour lesquelles il est *impossible* de trouver un estimateur (resp. prédicteur) non biaisé (alors qu'on peut néanmoins en trouver des estimateurs d'excellente qualité!). Pour ne citer qu'un exemple, dans le modèle du chasseur, on pourrait montrer qu'il n'existe aucun estimateur sans biais de la quantité d'intérêt $2/(1+\theta)$.

Il n'en demeure pas moins que, conformément à l'intuition, un estimateur ou prédicteur qui aurait un biais "vraiment" trop fort ne sera pas pertinent : on peut d'ailleurs montrer que, dans le cas de la fonction de perte quadratique, le risque est toujours au moins égal au carré du biais! ♣

En conclusion de cette section, les fonctions de risque et de biais nous fournissent, dans le cadre fréquentiste, des critères numériques pour jauger la qualité d'un estimateur (ou d'un prédicteur). Cependant, on est loin d'avoir une "pierre de touche" aussi décisive que ne l'était la notion de risque intégré dans le cadre bayésien...

C'est pourquoi il est fréquent de recourir aussi à des critères *qualitatifs* pour jauger la qualité des estimateurs et des prédicteurs. Dans la section suivante, nous allons parler du plus important de ces critères : à savoir, la notion de *convergence* des estimateurs en statistique asymptotique.

11.3 Convergence des estimateurs

La présente section traite de la notion de *convergence* des estimateurs. Il s'agit là d'un critère qualitatif asymptotique extrêmement important pour jauger si un estimateur est raisonnablement pertinent. La définition est la suivante :

Définition (WT). Considérons un modèle statistique explicative avec les notations génériques, dans lequel il y a un paramètre du modèle λ , pour lequel on considère un certain régime asymptotique noté « $\lambda \rightarrow \infty$ ». (Lorsque nous voudrions souligner la dépendance d'un certain objet en le paramètre du modèle, nous l'indiquerons par un exposant « (λ) »). Soient, pour un tel modèle, une quantité d'intérêt $\gamma(\theta)$ et un estimateur $\hat{\gamma}^{(\lambda)}(X) =: \hat{\gamma}^{(\lambda)}$ de cette quantité d'intérêt.

Dans un tel contexte, on dit que « l' »estimateur $\hat{\gamma}$ [ou, plus exactement, la famille d'estimateurs $(\hat{\gamma}^{(\lambda)})_{\lambda}^{[**]}$] est *convergent* lorsque, pour tout $\theta \in \Theta$, la loi

[1]. En particulier, cet estimateur n'est pas « convergent » lorsque $n \rightarrow \infty$: conférer le vocabulaire introduit dans la section suivante.

Loi $_{\theta}(\hat{\gamma}^{(\lambda)}(X^{(\lambda)}))$ (autrement dit, la loi de l'estimateur $\hat{\gamma}^{(\lambda)}$ dans le contexte probabiliste fréquentiste $\mathbb{P}_{\theta}(\bullet)$) converge, lorsque $\lambda \rightarrow \infty$, vers $\delta_{\gamma(\theta)}$ (autrement dit, vers la masse de Dirac en la valeur constante prise par γ dans le contexte probabiliste $\mathbb{P}_{\theta}(\bullet)$). \heartsuit

Remarque (WU). Vu que la convergence en loi vers une masse de Dirac est équivalente à la convergence en probabilité vers la constante correspondante, dans la définition précédente, si nous supposons \mathcal{E} muni d'une structure métrique, la propriété de convergence de l'estimateur $\hat{\gamma}$ revient à dire que :

$$\forall \theta \in \Theta \quad \forall \varepsilon > 0 \quad \mathbb{P}_{\theta}(\text{dist}(\gamma(\theta), \hat{\gamma}^{(\lambda)}(X^{(\lambda)})) > \varepsilon) \xrightarrow{\lambda \rightarrow \infty} 0. \quad (\text{WV})$$

\clubsuit

Notez au passage que la notion de convergence s'applique *uniquement* aux estimateurs, pas aux prédicteurs :

Remarque (WW). Il n'y a pas de notion de convergence pour les prédicteurs ; ou plus exactement, si on essayait de transposer la notion de convergence des estimateurs aux prédicteurs, on trouverait qu'aucun prédicteur n'est jamais convergent (sauf cas trivial)... En effet, autant se placer dans l'asymptotique d'une observation de plus en plus riche peut permettre d'obtenir une information arbitrairement fine sur la valeur du paramètre caché, autant l'aléa inhérent à l'observation future, lui, ne pourra pas être éliminé par l'accumulation d'observations ! \clubsuit

Voyons à présent quelques exemples d'estimateurs convergents (ou pas) :

Exemple (WX).

- (i) Dans le modèle du chasseur, considérons l'estimateur « naïf » de θ donné par $\hat{\theta}^{\text{naïf}} := X/n$. Notez déjà que, bien que je parle de « l' » estimateur, il s'agit en réalité d'une *famille* d'estimateurs, puisque l'application $x \mapsto x/n$ qui dit comment on crée l'estimateur à partir de l'observation dépend du paramètre du modèle n ! Considérons maintenant l'asymptotique $n \rightarrow \infty$. Sous le véritable contexte probabiliste $\mathbb{P}_{\nu}(\bullet)$ ^[††], $\hat{\theta}^{\text{naïf}(n)}$ apparaît comme la moyenne de n v.a.i.d. Bernoulli(θ_{ν}) : il s'ensuit, par la loi des grands nombres (au sens de votre cours du semestre 5), que la loi de cet estimateur tend vers $\delta_{\theta_{\nu}}$ lorsque $n \rightarrow \infty$. Et puisque nous avons écrit cela sans rien spécifier sur θ_{ν} , on a en fait que $\text{Loi}_{\theta}(\hat{\theta}^{\text{naïf}(n)}) \xrightarrow{n \rightarrow \infty} \delta_{\theta}$ pour tout $\theta \in]0, 1[$: ainsi l'estimateur $\hat{\theta}^{\text{naïf}}$ est convergent ! \smile
- (ii) Toujours dans le modèle du chasseur, considérons maintenant l'estimateur $\hat{\theta}^{\text{Bay}} := (X + 1/2)/(n + 1/2)$ ^[‡‡]. Des calculs élémentaires nous montrent qu'on a toujours $|\hat{\theta}^{\text{Bay}} - \hat{\theta}^{\text{naïf}}| \leq 1/n$. Or nous avons vu au point précédent que, sous le contexte probabiliste $\mathbb{P}_{\nu}(\bullet)$, l'estimateur $\hat{\theta}^{\text{naïf}}$ converge en loi vers la constante θ_{ν} . Par conséquent, l'estimateur converge aussi au sens de la convergence en probabilité : en particulier, des seuils arbitrairement petits $\varepsilon, \varepsilon' > 0$ étant supposés fixés, on aura pour n suffisamment grand que $\mathbb{P}_{\nu}(|\hat{\theta}^{\text{naïf}} - \theta_{\nu}| \geq \varepsilon/2) \leq \varepsilon'$. Dans ce régime, il s'ensuit, par l'inégalité triangulaire, qu'on aura

[**]. En pratique néanmoins, on dira toujours « l'estimateur » plutôt que « la famille d'estimateurs », la dépendance en le paramètre du modèle étant évidente.

[††]. En fait, là aussi, on devrait parler, en toute rigueur, de la *famille* de contextes probabilistes $\mathbb{P}_{\nu}^{(n)}(\bullet)$ associés à la vraie valeur θ_{ν} du paramètre caché.

[‡‡]. Cet estimateur correspond à l'espérance a posteriori de θ dans le paradigme bayésien ; mais ici nous allons le regarder en termes purement fréquentistes.

$\mathbb{P}_{\mathcal{J}}(|\hat{\theta}^{\text{Bay}} - \theta_{\mathcal{J}}| \geq \varepsilon/2 + 1/n) \leq \varepsilon'$. Or, si n est suffisamment grand, la borne $\varepsilon/2 + 1/n$ devient $\leq \varepsilon$: par conséquent, on a alors que $\mathbb{P}_{\mathcal{J}}(|\hat{\theta}^{\text{Bay}} - \theta_{\mathcal{J}}| \geq \varepsilon) \leq \varepsilon'$, ce qui montre que $\hat{\theta}^{\text{Bay}}$ converge vers $\theta_{\mathcal{J}}$ au sens de la convergence en probabilité. Dès lors, $\text{Loi}_{\mathcal{J}}(\hat{\theta}^{\text{Bay}})$ converge vers $\delta_{\theta_{\mathcal{J}}}$; et puisque nous avons établi tout cela sans rien dire sur $\theta_{\mathcal{J}}$, on en déduit que l'estimateur $\hat{\theta}^{\text{Bay}}$ est convergent lui aussi !

- (iii) Encore et toujours dans le modèle du chasseur, considérons à présent la quantité d'intérêt explicative $\tau := 2\theta - \theta^2$; et proposons-nous d'estimer celle-ci par $\hat{\tau}^{\text{naïf}} := 2\hat{\theta}^{\text{naïf}} - (\hat{\theta}^{\text{naïf}})^2$. Dans la suite, pour simplifier les notations, nous poserons $2\theta - \theta^2 =: \tau(\theta)$. Pour établir que l'estimateur $\hat{\tau}^{\text{naïf}}$ est convergent, on voudrait montrer que, sous $\mathbb{P}_{\mathcal{J}}(\bullet)$, on ait convergence de $\hat{\tau}^{\text{naïf}}$ vers $\tau_{\mathcal{J}}$ au sens de la convergence en probabilité, c.à.d. établir que, pour tous $\varepsilon, \varepsilon' > 0$, on ait $\mathbb{P}_{\mathcal{J}}(|\hat{\tau}^{\text{naïf}} - \tau_{\mathcal{J}}| \geq \varepsilon) \leq \varepsilon'$ pour n suffisamment grand. Or l'évènement $\{|\hat{\tau}^{\text{naïf}} - \tau_{\mathcal{J}}| \geq \varepsilon\}$ peut se ré-écrire comme $\{|\tau(\hat{\theta}^{\text{naïf}}) - \tau(\theta_{\mathcal{J}})| \geq \varepsilon\}$; et puisque la fonction $\tau(\bullet)$ est continue en $\theta_{\mathcal{J}}$, cet évènement implique que $\{|\hat{\theta}^{\text{naïf}} - \theta_{\mathcal{J}}| \geq \delta\}$ pour peu que nous choissions $\delta > 0$ suffisamment petit... Et comme la convergence de l'estimateur $\hat{\theta}^{\text{naïf}}$ vers θ nous assure que, quel que soit $\delta > 0$, on a $\mathbb{P}_{\mathcal{J}}(|\hat{\theta}^{\text{naïf}} - \theta_{\mathcal{J}}| \geq \delta) \leq \varepsilon'$ pour n suffisamment grand, on obtient bien le contrôle dont nous avons besoin pour établir la convergence de l'estimateur $\hat{\tau}^{\text{naïf}}$! ☘

Exemple (WY).

- (i) Dans le modèle du pédagogue, considérons l'estimateur de μ consistant à prendre, tout simplement, la moyenne de la promotion, autrement dit, $\hat{\mu}^{\text{naïf}} := (X_0 + \dots + X_{n-1})/n$. (Notez bien qu'il s'agit en fait, là encore, d'une *famille* n'estimateurs : non seulement parce que n intervient dans la formule définissant $\hat{\mu}^{\text{naïf}}$, mais aussi parce que le *domaine* de l'application $(x_0, \dots, x_{n-1}) \mapsto (x_0 + \dots + x_{n-1})/n$ qui permet de déduire $\hat{\mu}^{\text{naïf}}$ de l'observation dépend lui-même de n !). Sous le véritable contexte probabiliste $\mathbb{P}_{\mathcal{J}}(\bullet)$, les X_i sont i.i.d. Normale($\mu_{\mathcal{J}}, \sigma_{\mathcal{J}}^2$), de sorte que $\hat{\mu}^{\text{naïf}}$ suit la loi Normale($\mu_{\mathcal{J}}, \sigma_{\mathcal{J}}^2/n$) (en vertu des propriétés des lois normales) : laquelle loi converge bien vers $\delta_{\mu_{\mathcal{J}}}$ lorsque $n \rightarrow \infty$. Cela ayant été établi sans rien spécifier sur $\theta_{\mathcal{J}}$, on a bien montré la convergence de l'estimateur $\hat{\mu}^{\text{naïf}}$.
- (ii) Bien entendu, dans ce modèle, $\hat{\mu}^{\text{naïf}}$ n'est pas le seul estimateur convergent (dans l'asymptotique $n \rightarrow \infty$) de μ ... Par exemple, l'enseignant pourrait considérer un estimateur $\hat{\mu}^{\text{pond}}$ dans laquelle les notes de ses ouailles seraient *pondérées* de sorte à donner plus d'importance aux étudiants de milieu de classement qu'à ceux des extrêmes, de la façon suivante : à partir des notes $x_{0\mathcal{J}}, \dots, x_{(n-1)\mathcal{J}}$, il commencerait par les trier de la plus petite à la plus grande, obtenant ainsi des notes ordonnées $r_{0\mathcal{J}}, \dots, r_{(n-1)\mathcal{J}}$ (avec $\{r_{0\mathcal{J}}, \dots, r_{(n-1)\mathcal{J}}\} = \{x_{0\mathcal{J}}, \dots, x_{(n-1)\mathcal{J}}\}$ et $r_{0\mathcal{J}} \leq r_{1\mathcal{J}} \leq \dots \leq r_{(n-1)\mathcal{J}}$) ; puis il définirait

$$\hat{\mu}^{\text{pond}} := \frac{x_{0\mathcal{J}} + 2x_{1\mathcal{J}} + 3x_{2\mathcal{J}} + \dots + 3x_{(n-3)\mathcal{J}} + 2x_{(n-2)\mathcal{J}} + x_{(n-1)\mathcal{J}}}{1 + 2 + 3 + \dots + 3 + 2 + 1}. \quad (\text{WZ})$$

La v.a. $\hat{\mu}^{\text{pond}}$ ainsi obtenue est alors bien une statistique [ou plus exactement une *famille* de statistiques indexée par n], puisque l'opération pour l'obtenir à partir de l'observation, bien que compliquée, est entièrement déterministe ; et on pourrait démontrer^[*] que sa loi sous $\mathbb{P}_{\mathcal{J}}(\bullet)$ converge, lorsque $n \rightarrow \infty$, vers

[*]. Ce serait cependant sensiblement plus difficile que dans le cas de $\hat{\mu}^{\text{naïf}}$... La démonstration

la masse de Dirac $\delta_{\mu_{\mathcal{V}}}$: il s'agit donc bien, là aussi, d'un estimateur convergent de μ .

- (iii) Voyons enfin, dans le cas du pédagogue, un exemple d'estimateur de μ qui ne soit *pas* convergents. Si l'enseignant est paresseux, il peut considérer l'estimateur $\hat{\mu}^{\text{par}}$ de μ obtenu en faisant la moyenne de ses copies *en s'arrêtant à 25 élèves s'ils sont trop nombreux* : $\hat{\mu}^{\text{par}}$ coïncide donc avec $\hat{\mu}^{\text{naïf}}$ pour $n \leq 25$, et vaut $(X_0 + X_1 + \dots + X_{23})/25$ pour $n > 25$ [†]. Dans ce cas, dans l'asymptotique $n \rightarrow \infty$, on a à partir d'un certain rang que $\text{Loi}_{\mathcal{V}}(\hat{\mu}^{\text{par}}) = \text{Normale}(\mu_{\mathcal{V}}, \sigma_{\mathcal{V}}^2/25)$: cette fois-ci on n'a pas convergence vers $\delta_{\mu_{\mathcal{V}}}$, donc l'estimateur $\hat{\mu}^{\text{par}}$ n'est pas convergent. (Même si, en pratique, cela demeure un assez bon estimateur malgré tout ☺).

Remarque (XA). La comparaison entre les items numéros (i) et (iii) de l'exemple (WY) ci-dessus met en valeur le fait que c'est fondamentalement à une *famille* d'estimateurs que la notion de convergence s'applique : en effet, si dans le modèle du pédagogue on ne regarde que le cas $n = 22$ (par exemple), on voit que les deux estimateurs $\hat{\mu}^{\text{naïf}}$ et $\hat{\mu}^{\text{par}}$ sont alors rigoureusement identiques : il n'y aurait donc aucun sens, dans une telle situation, à dire que l'un est convergent, mais pas l'autre...!

Remarque (XB). En pratique, le critère de convergence des estimateurs est plutôt "facile" à satisfaire. Plus précisément : soit on sera dans une asymptotique qui ne fournit pas suffisamment d'information sur la quantité d'intérêt à la limite, et alors il n'existera aucun estimateur convergent de toutes façons ; soit l'asymptotique nous permet de construire des estimateurs convergents, et alors essentiellement tous les estimateurs "intelligents" auxquels on pourra penser seront effectivement convergents ! ☺

À l'issue de ces trois premières sections, nous comprenons à présent bien ce que signifie « être un bon estimateur (ou prédicteur) ». Néanmoins, cela ne dit pas pour autant *comment* trouver de tels bons estimateurs ou prédicteurs... C'est la question à laquelle nous allons maintenant nous attaquer ; dans la suite de ce chapitre, nous allons présenter différentes méthodes permettant de trouver des estimateurs ou prédicteurs intelligents — lesquels s'avèreront effectivement, en général, être « bons » au sens des critères présentés ci-devant ! ☺

11.4 Estimateur du maximum de vraisemblance

La méthode d'estimation par maximum de vraisemblance, que cette section va présenter, est la "star" des méthodes d'estimation : elle a en effet le triple mérite d'être définie sans ambiguïté, de fonctionner pour essentiellement tous les types de modèles, et d'avoir de très bonnes propriétés d'optimalité ! Elle présente néanmoins un inconvénient, c'est que les calculs qu'elle requiert sont compliqués, et pas toujours menables à bien : d'où l'intérêt des autres méthodes d'estimations présentées par ailleurs dans ce chapitre.

Définissons donc la méthode :

!

 serait à votre portée sous forme d'exercice avec questions intermédiaires ; mais sans aide, elle reste trop ardue pour des élèves-ingénieurs, je pense ! ☺

[†]. Notons que dans ce cas, bien que la formule définissant l'estimateur à partir des X_i soit toujours la même, il s'agit *quand même, techniquement*, d'estimateurs *différents* pour des valeurs de n différentes, dans la mesure où il s'appliquent à des *modèles* différents !

Définition (XC) (Estimateur du maximum de vraisemblance). Soit $\Theta \ni \theta \mapsto \text{Loi}_\theta(X) \in \mathcal{M}_1(\mathcal{X})$ un modèle statistique explicatif pour lequel on peut définir une fonction de vraisemblance $\mathcal{L}(\theta \mid x)$. Alors :

- (i) $\hat{\theta}(X)$ est qualifié d'*estimateur du maximum de vraisemblance* pour le paramètre caché θ lorsque, pour tout $x \in \mathcal{X}$, $\hat{\theta}(x)$ est le (ou un) point où $\theta \mapsto \mathcal{L}(\theta \mid x)$ atteint son maximum ;
- (ii) Pour $\gamma : \Theta \rightarrow \mathcal{E}$ une fonction, on dit que $\gamma(\hat{\theta})$ est un estimateur du maximum de vraisemblance pour la quantité d'intérêt $\gamma(\theta)$. ♥

Autrement dit, l'estimation par maximum de vraisemblance consiste à estimer θ par la valeur qui nous donnait la plus grande chance de tomber sur les données que nous avons effectivement observées ; et on en déduit, par composition, l'estimation par maximum de vraisemblance de n'importe quelle quantité d'intérêt explicative.

Remarque (XD). Dans le cas de l'estimation par maximum de vraisemblance, on remarquera que la *réalisation* de l'estimateur du maximum de vraisemblance (ce que j'appelle l'« estimation » du maximum de vraisemblance) se détermine plus directement que l'estimateur lui-même : en effet, il suffit alors de considérer la fonction de vraisemblance *pour l'observation qu'on a effectivement obtenue*, et à regarder où elle atteint son maximum. Et dans nombre de situations pratiques, en fait, seule la *réalisation* de l'estimateur nous intéressera : dans de tel cas, on se s'embêtera évidemment pas à écrire l'estimateur du maximum de vraisemblance en version « variable aléatoire » ! ♣

Remarque (XE). Attention, le vocabulaire est un peu glissant : nonobstant son nom, l'estimation du maximum de vraisemblance n'est pas le maximum de la vraisemblance, mais le lieu où ce maximum est atteint : autrement dit, c'est en fait *l'argument du maximum* de la fonction de vraisemblance... ♣

Remarque (XF). On voit que la méthode d'estimation par maximum de vraisemblance ne laisse pas de place à l'arbitraire dans la façon de le construire (en effet, le statisticien n'a rien à choisir), ce qui est évidemment un point très intéressant ! ♣

Remarque (XG). L'estimateur du maximum de vraisemblance n'existe pas forcément (le supremum de la fonction de vraisemblance pouvant n'être atteint qu'asymptotiquement) ; et quand il existe, il n'est pas forcément unique (car le maximum de la fonction de vraisemblance peut, à l'inverse, être atteint en plusieurs valeurs). Heureusement, on aura quand même unique existence la plupart du temps ♣

On peut montrer que, dans les modèles avec paramètre du modèle disposant d'un régime asymptotique, les estimateurs par maximum de vraisemblance sont en général convergents. On a notamment des résultats précis dans ce sens pour les modèles d'échantillonnage^[‡], comme le théorème ci-dessous, que je mentionne pour la culture :

Théorème (XH). *Dans un modèle d'échantillonnage, sous certaines hypothèses de régularité toujours vérifiées en pratique, l'estimateur du maximum de vraisemblance est convergent.* ♦

Remarque (XI). On notera que la méthode du maximum de vraisemblance relève fondamentalement de l'inférence *explicative* : il n'existe pas de notion de soi-disant « prédicteur par maximum de vraisemblance » ! De manière générale, en statistique

[‡]. Ceux d'entre vous qui suivront le cours électif « Problèmes inverses » au semestre 8 (réf. 8KU-CEN07) auront l'occasion d'en apprendre davantage à ce sujet lors de la partie statistique de ce cours-là (dispensée par M^{me} FRISTCH).

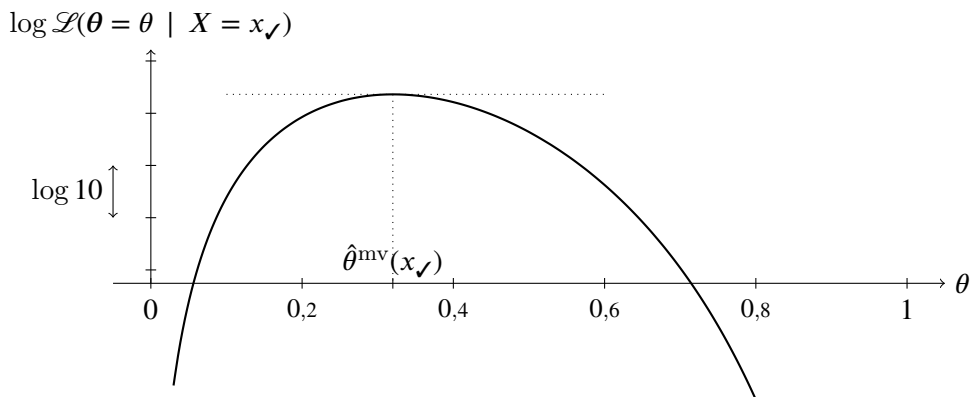


FIGURE 11.2 – Détermination du maximum de vraisemblance pour le modèle du chasseur.

fréquentiste, la plupart des méthodes de ce chapitre ne concerneront que des quantités d'intérêt explicatives : pour ce qui est de la prédiction, on procèdera généralement en deux temps, confer § ?? [§]. ♣

Maintenant que nous avons expliqué *ce qu'est* l'estimation par maximum de vraisemblance, regardons quelques exemples ! Ci-dessous nous allons calculer l'estimateur par maximum de vraisemblance pour nos deux exemples favoris ☺

Exemple (XJ) (Maximum de vraisemblance pour le chasseur). Nous avons vu que la log-vraisemblance du modèle du chasseur est

$$\ln \mathcal{L}(\theta) = x_{\mathcal{J}} \ln(\theta) + (n - x_{\mathcal{J}}) \ln(1 - \theta) : \quad (\text{XK})$$

on voit sur la figure 11.2 que, lorsque θ varie (on a pris ici $n = 25$ et $x_{\mathcal{J}} = 8$), cette fonction de vraisemblance semble bien prendre des valeurs très faibles pour des valeurs “invraisemblables” de θ (car il apparaît invraisemblable qu'un tireur ayant touché 8 cibles et en ayant manqué 17 ait un taux de succès véritable inférieur à 1 % ou supérieur à 90 %, par exemple), et passer par un maximum autour de 0,3. Calculons plus précisément où ce maximum est atteint, en annulant le dérivée de la fonction de log-vraisemblance : on calcule que

$$\frac{\partial}{\partial \theta} \ln \mathcal{L}(\theta) = \frac{x_{\mathcal{J}}}{\theta} - \frac{n - x_{\mathcal{J}}}{1 - \theta}, \quad (\text{XL})$$

qui s'annule lorsque $\theta = n^{-1}x_{\mathcal{J}}$: ainsi, l'estimation du maximum de vraisemblance pour θ est $n^{-1}x_{\mathcal{J}}$.

Et si on cherche l'estimateur du maximum de vraisemblance en tant que *variable aléatoire*, il faut alors simplement donner la v.a. dont $n^{-1}x_{\mathcal{J}}$ est la réalisation : en l'occurrence, c'est simplement $n^{-1}X$. Nous pourrions noter cette statistique ' $\hat{\theta}^{\text{mv}}$ ' pour indiquer qu'il s'agit de l'estimateur du maximum de vraisemblance. ♣

Exemple (XM) (Maximum de vraisemblance pour le pédagogue). Nous avons calculé que la log-vraisemblance, dans l'exemple du pédagogue, était

$$\begin{aligned} \ln \mathcal{L}(\theta) = & \ln \mathcal{L}(\theta = (\mu, \sigma) \mid X = (x_{1\mathcal{J}}, \dots, x_{n\mathcal{J}})) = \\ & -n \ln \sigma - \frac{1}{2} n \sigma^{-2} \left((\mu - \text{moy}(x_{1\mathcal{J}}, \dots, x_{n\mathcal{J}}))^2 + \text{var}_{\text{emp}}(x_{1\mathcal{J}}, \dots, x_{n\mathcal{J}}) \right) : \end{aligned}$$

[§]. Section non encore rédigée à l'heure actuelle

il s'agit d'ajuster μ et σ pour maximiser cette quantité. En fait, on s'aperçoit qu'à σ fixé, la vraisemblance est maximisée pour $\mu = \text{moy}(x_{1\checkmark}, \dots, x_{n\checkmark})$; de sorte que, si on admet qu'il existe un maximum de vraisemblance, il est nécessairement atteint en un point $(\hat{\mu}_{\checkmark}, \hat{\sigma}_{\checkmark})$ tel que $\hat{\mu}_{\checkmark} = \text{moy}(x_{1\checkmark}, \dots, x_{n\checkmark})$: ainsi, sans même connaître le maximum de vraisemblance $\hat{\theta}(x_{\checkmark})$ pour θ , nous pouvons déjà conclure que (la réalisation de) le maximum de vraisemblance pour μ est égal à $\text{moy}(x_{1\checkmark}, \dots, x_{n\checkmark})$.

Avec un peu de travail supplémentaire, on montrerait qu'il y a effectivement un (unique) maximum de vraisemblance pour θ , et que (la réalisation de) celui-ci est égal à $(\text{moy}(x_{1\checkmark}, \dots, x_{n\checkmark}), \text{var}_{\text{emp}}^{1/2}(x_{1\checkmark}, \dots, x_{n\checkmark}))$: en l'occurrence donc, l'estimateur du maximum de vraisemblance pour $(\mu_{\checkmark}, \sigma_{\checkmark})$ coïncide avec les moyenne et écart-type empiriques! (confer § 11.5 *infra*). \clubsuit

11.5 Estimateurs empiriques

La notion d'*estimation empirique* nous allons présenter dans cette section est une des notions d'estimation plus simples à saisir et à appliquer; néanmoins son usage sera assez restreint, car on ne peut l'appliquer qu'à certains modèles particuliers... Cela englobera néanmoins une des classes de modèles les plus importantes, à savoir, celle des modèles d'échantillonnage.

Les modèles où nous nous intéressons dans le cadre de cette section sont ceux relevant du point ci-dessous :

Point (XN). Dans cette section, nous nous intéresserons à des modèles dont l'observation se décompose en $n > 1$ sous-observations (n étant un paramètre du modèle) qui sont toutes de même loi lorsqu'on se place sous un contexte probabiliste fréquentiste: autrement dit, avec les notations génériques, on a n sous-observations X_0, \dots, X_{n-1} qui vérifient, pour tout $\theta \in \Theta$, $\text{Loi}_{\theta}(X_j) = \text{Loi}_{\theta}(X_i) \forall i, j \in \llbracket 0, n \rrbracket$. On pourra alors aussi noter l'observation globale \vec{X} pour rappeler qu'il s'agit en fait d'un "vecteur" de sous-observations. L'espace dans lequel vivent chacune de ces sous-observations sera alors noté $\mathcal{X}^{(1)}$. (L'espace de l'observation globale \mathcal{X} étant alors $(\mathcal{X}^{(1)})^n$, qu'on pourra aussi noter ' $\mathcal{X}^{(n)}$ ').

La situation décrite par ce point englobe notamment tous les modèles d'échantillonnage, qui correspondent au cas où, en plus de ce qui précède, les X_i sont indépendants sous les contextes $\mathbb{P}_{\theta}(\bullet)$. Le modèle du pédagogue, par exemple, rentre dans un tel cadre. \clubsuit

Dans un tel modèle, sous le véritable contexte probabiliste, les sous-observations $x_{0\checkmark}, \dots, x_{(n-1)\checkmark}$ constituent un *échantillon* de tirages de la loi $\text{Loi}_{\checkmark}(X_0)$. Si, en outre, on est dans un modèle d'échantillonnage, ces tirages sont indépendants: par conséquent, si n est suffisamment grand, on s'attend d'après la loi des grands nombres à ce que la distribution de cet échantillon ressemble à celle de la loi dont ils sont tirés! On peut même espérer que cela reste vrai pour des modèles généraux où il n'y aurait pas indépendance *stricto sensu*, dès lors que les réalisations X_i de la loi $\text{Loi}_{\checkmark}(X_0)$ soient suffisamment bien « mélangées ».

Par exemple, dans le modèle du pédagogue, on s'attend à ce que pour n grand, la loi (sous le véritable contexte probabiliste) de la variable aléatoire $n^{-1} \text{card}\{i \mid X_i \geq 86\}$ soit très concentrée autour de la valeur constante $\mathbb{P}(\text{Normale}(\mu_{\checkmark}, \sigma_{\checkmark}^2) \geq$

86)^[¶]; de même, $n^{-1} \sum_i X_i$ devrait être proche de la constante $E(\text{Normale}(\mu_{\mathcal{V}}, \sigma_{\mathcal{V}}^2))$ (qui, en l'occurrence, n'est autre que $\mu_{\mathcal{V}}$); et le quantile de niveau 80 % de l'(multi)ensemble des X_i (i.e. la valeur telle que 80 % des notes des élèves sont tombées en-dessous et 20 % au-dessus : par exemple, si $n = 22$, c'est la note de l'élève classé 5^e) devrait être proche du quantile de niveau 80 % de la loi Normale($\mu_{\mathcal{V}}, \sigma_{\mathcal{V}}^2$)^[||].

En fait, ces trois cas peuvent tous être reformulés dans un même cadre, en disant qu'une certaine fonction de la *loi empirique* des (sous-)observations converge vers la valeur correspondante de la véritable distribution $\text{Loi}_{\mathcal{V}}(X_0)$ de ces observations. Rappelons ci-dessous le concept de « loi empirique » (déjà introduit en définition (DR)) :

! **Définition (X0)** (Loi empirique). Soient x_0, \dots, x_{n-1} un jeu de n valeurs appartenant à un même espace $\mathcal{X}^{(1)}$; alors la *loi empirique* des x_i est la loi de probabilité qui donne le poids $1/n$ à chaque x_i (en comptant plusieurs fois les éventuelles valeurs identiques), soit

$$\frac{1}{n} \sum_{i=0}^{n-1} \delta_{x_i} : \quad (\text{XP})$$

en d'autres termes, c'est la distribution qui donne à $A \subseteq \mathcal{X}^{(1)}$ une probabilité égale à la *proportion* des $x_{i\mathcal{V}}$ qui sont tombés dans A .

Dans le cas d'un modèle relevant du point (XN), on appelle *loi empirique des (sous-)observations* la variable aléatoire à valeurs dans $\mathcal{M}_1(\mathcal{X}^{(1)})$ dont la réalisation est la loi empirique des $x_{i\mathcal{V}}$. ♡

Remarque (XQ). Bien comprendre que la loi empirique des observations est une loi de probabilité, mais que cette loi dépend directement de l'observation, et constitue donc (si on se place avant de faire l'expérience) un objet *aléatoire*! Autrement dit, la loi empirique des observations est une *variable aléatoire à valeurs dans les distributions de probabilité* (outch!). ♣

Avec ce vocabulaire, on voit bien que, dans les trois exemples cités avant la définition (X0), on était à chaque fois en train de comparer une certaine propriété de la loi empirique des observations avec la propriété correspondante de la loi véritable. Par exemple, $n^{-1} \text{card}\{i \mid x_{i\mathcal{V}} \geq 86\}$ n'était autre que la masse que la distribution empirique attribuait aux valeurs ≥ 86 , et nous disions que cela devait tendre, pour n grand, vers la masse que la loi des observations Normale($\mu_{\mathcal{V}}, \sigma_{\mathcal{V}}^2$) attribue aux valeurs ≥ 86 : rien de plus logique, puisque la distribution empirique des observations est censée ressembler à leur véritable loi! Et les deux autres exemples relevaient du même schéma, en remplaçant « masse attribuée aux valeurs ≥ 86 » par resp. « espérance » et « quantile de niveau 80 % ».

On peut généraliser cette idée *ad libitum* : c'est en cela que consiste la notion d'*estimation empirique* :

! **Définition (XR)** (Estimateur empirique). Plaçons-nous dans le cadre d'un modèle relevant du point (XN), et supposons qu'il existe une relation fonctionnelle

[¶]. Qui, en l'occurrence, peut être exprimée comme valant $1/2 + \frac{1}{2} \text{erf}((\mu_{\mathcal{V}} - 86) / \sqrt{2}\sigma_{\mathcal{V}})$, où 'erf' désigne la fonction spéciale appelée « fonction d'erreur ».

[||]. Qui, en l'occurrence, vaut $\mu_{\mathcal{V}} + \text{erf}^{-1}(0,6) \times \sqrt{2}\sigma_{\mathcal{V}}$.

entre $\text{Loi}_\theta(X_1)$ et la quantité d'intérêt $\gamma(\theta) \in \mathcal{G}$, exprimée par une certaine fonctionnelle^[**] $F(\bullet)$

$$\forall \theta \in \Theta \quad F(\text{Loi}_\theta(X_1)) = \gamma(\theta), \quad (\text{XS})$$

où l'expression « $F(P)$ » a du sens pour n'importe quelle distribution de probabilité P sur $\mathcal{X}^{(1)}$ (y compris celles qui ne sont pas de la forme $\text{Loi}_\theta(X_0)$). Alors une estimation de $\gamma(\theta)$ est donnée par

$$F\left(\frac{1}{n} \sum_{i=0}^{n-1} \delta_{x_{i\checkmark}}\right) : \quad (\text{XT})$$

dans ce cas, on parle d'*estimation empirique* (et on appelle « estimateur empirique » l'estimateur dont notre estimation est la réalisation).

En d'autres termes, pour « bidule » un concept pouvant s'appliquer à des distributions de probabilité, une procédure d'estimation empirique consiste à estimer le *bidule* de la distribution des observations par le *bidule empirique* (cf. définition (DW)) des observations. ♡

Remarque (XU). Il peut tout à fait exister plusieurs fonctionnelles différentes F_0, F_1, \dots vérifiant chacune la relation $\gamma(\theta) = F_i(\text{Loi}_\theta(X_0))$, menant à autant d'estimateurs empiriques différents. Par exemple, dans le cas du pédagogue, μ peut être certes être vue comme la moyenne de $\text{Normale}(\mu, \sigma^2)$, mais c'est aussi la médiane de cette loi : si on préfère utiliser cette fonctionnelle-là, on obtient donc un autre estimateur empirique de μ ! Par exemple, dans le cas où $n = 5$ et $(x_{0\checkmark}, \dots, x_{4\checkmark}) = (52, 79, 88, 91, 70)$, l'estimation empirique de μ par moyenne vaut 76, tandis que l'estimation empirique de μ par médiane vaut 79... ♣

On l'aura compris, les estimateurs empiriques sont très simples à fabriquer ! En outre, pour peu que le modèle présente de suffisamment bonnes propriétés de « mélange », ils auront en général le bon goût d'être convergents lorsque le nombre d'observations tend vers l'infini. Citons ainsi, pour la culture, le résultat suivant :

Théorème (XV). *Lorsqu'on est dans le cas d'un modèle d'échantillonnage, sous des hypothèses assez souples de régularité de la fonctionnelle $F(\bullet)$ ^[††], l'estimateur empirique associé à la fonctionnelle F est convergent lorsque la taille de l'échantillon tend vers l'infini.* ◇

Remarque (XW). De même que pour le maximum de vraisemblance, la méthode empirique est une technique spécifique aux questions d'*estimation* : il n'existe pas de notion de soi-disant « prédictor empirique » ! ♣

[**]. Note de vocabulaire : Lorsqu'on considère une fonction $F(\bullet)$ dont l'espace de définition est lui-même un espace de fonctions (ou quelque chose d'apparenté, comme un espace de lois de probabilités en l'occurrence), l'usage des mathématiciens est de parler de « fonctionnelle », plutôt que de « fonction », pour désigner $F(\bullet)$, afin de limiter les risques de confusion. (Par exemple, une fonction $\varphi : \mathcal{X} \rightarrow \mathbb{R}$ étant donnée, l'application « espérance contre φ », définie de $\mathcal{M}_1(\mathcal{X})$ dans \mathbb{R} par $P \mapsto \mathbb{E}^{X \sim P}(\varphi(X))$, sera volontiers qualifiée de « fonctionnelle », afin d'éviter les confusions avec la « fonction » $\varphi(\bullet)$). En général, quand on parle de « fonctionnelle », on sous-entend en outre que l'espace d'arrivée (qui, dans le cadre de la présente définition, est l'espace \mathcal{G} de la quantité d'intérêt) est quelque chose de « simple », comme \mathbb{R} .

[††]. Mais qu'il serait fort compliqué d'explicitier ici... ☹

11.6 Méthode des moindres carrés

Nous allons maintenant présenter la méthode dite « des moindres carrés ». Bien que cette méthode ne soit pas aussi universelle que l'estimation par maximum de vraisemblance, elle peut tout de même être appliquée à une très grande variété de modèles. La méthode des moindres carrés sera tout particulièrement adaptée à tous les modèles où l'observation est vue comme la somme d'une « tendance » déterministe et d'un « bruit » centré aléatoire, ce qui la rend très appréciée dans le domaine de l'*analyse de données*, qui vous sera enseigné au semestre 7.

La catégorie de modèles à laquelle s'applique la méthode des moindres carrés est décrite dans le point ci-dessous :

Point (XX). Dans cette section, nous considérerons un modèle statistique dont nous noterons l'espace du paramètre caché Θ , et dont l'observation X se décompose en n sous-observations réelles X_0, \dots, X_{n-1} , qui sont supposées d'intégrabilité L^1 sous les différents contextes probabilistes (autrement dit, on peut donner du sens aux $\mathbb{E}_\theta(X_i)$).

On notera que, par contre, on ne requiert absolument pas que les X_i aient la même loi sous $\mathbb{P}_\theta(\bullet)$, ni qu'ils soient indépendants ! \clubsuit

Au fondement de la méthode des moindres carrés est l'observation suivante : parmi les différentes valeurs possibles θ du paramètre caché, c'est lorsqu'on aura $\theta = \theta_{\mathcal{J}}$ que les $\mathbb{E}_\theta(X_i)$ seront, *en moyenne* (sous le véritable contexte probabiliste), les plus proches des X_i au sens de l'erreur quadratique ! Formellement, il s'agit de la proposition suivante :

Proposition (XY). *Considérons un modèle statistique relevant du point (XX) ; supposons en outre que pour tout $\theta \in \Theta$, les X_i soient toutes d'intégrabilité L^2 sous $\mathbb{P}_\theta(\bullet)$. Introduisons alors la fonction suivante, dite fonction de contraste quadratique :*

$$c(\bullet, \bullet) : \mathbb{R}^n \rightarrow \mathbb{R}_+ \quad (XZ)$$

$$\left((x_0, \dots, x_{n-1}), \theta \right) \mapsto \sum_{i=0}^{n-1} (x_i - \mathbb{E}_\theta(X_i))^2.$$

alors, pour une valeur $\theta_{\mathcal{J}}$ définissant un véritable contexte probabiliste $\mathbb{P}_{\mathcal{J}}(\bullet)$, la fonction

$$\Theta \ni \theta \mapsto \mathbb{E}_{\mathcal{J}}(c(\vec{X}, \theta)) \quad (YA)$$

atteint son minimum en $\theta_{\mathcal{J}}$. \diamond

Démonstration. En effet, pour $\theta \in \Theta$,

$$\begin{aligned} \mathbb{E}_{\mathcal{J}}((X_i - \mathbb{E}_\theta(X_i))^2) &= \mathbb{E}_{\mathcal{J}}(X_i - \mathbb{E}_\theta(X_i))^2 + \text{Var}_{\mathcal{J}}(X_i - \mathbb{E}_\theta(X_i))^2 \\ &= (\mathbb{E}_{\mathcal{J}}(X_i) - \mathbb{E}_\theta(X_i))^2 + \text{Var}_{\mathcal{J}}(X_i), \end{aligned}$$

de sorte que

$$\mathbb{E}_{\mathcal{J}}\left(\sum_{i=1}^n (X_i - \mathbb{E}_\theta(X_i))^2\right) = \sum_{i=1}^n \text{Var}_{\mathcal{J}}(X_i) + \sum_{i=1}^n (\mathbb{E}_{\mathcal{J}}(X_i) - \mathbb{E}_\theta(X_i))^2, \quad (YB)$$

où le premier terme est indépendant de θ tandis que le second, qui est évidemment toujours positif, est nul pour $\theta = \theta_{\mathcal{J}}$, de sorte que globalement l'expression est bien minimisée pour $\theta = \theta_{\mathcal{J}}$. \spadesuit

La proposition (XY) suggère alors l'idée suivante : puisque θ_{\checkmark} est la valeur qui minimise la fonction de contraste quadratique *en moyenne*, alors elle ne devrait pas se situer trop loin de la valeur qui minimisera la fonction de contraste quadratique *pour l'observation effective* ! La valeur $\hat{\theta}_{\checkmark}^{\text{mc}}$ réalisant cette minimisation est appelée *estimateur des moindres carrés* de notre modèle :

Définition (YC) (Estimateur des moindres carrés). Pour un modèle statistique relevant du point (XX), l'*estimation des moindres carrés* pour θ associée à l'observation effective $(x_{0\checkmark}, \dots, x_{(n-1)\checkmark})$ est définie comme la valeur de θ minimisant la fonction de contraste des moindres carrés pour l'observation effective :

$$\hat{\theta}_{\checkmark}^{\text{mc}} := \arg \min_{\theta} \sum_{i=1}^n (x_{i\checkmark} - \mathbb{E}_{\theta}(X_i))^2. \quad (\text{YD})$$

On peut ensuite éventuellement en déduire, par composition, une estimation de $\gamma(\hat{\theta}_{\checkmark}^{\text{mc}})$ d'une quantité d'intérêt de la forme $\gamma(\theta)$, qu'on qualifiera là encore d'« estimation des moindres carrés » le cas échéant.

Remarque (YE). En fait, il ne sera pas rare que le minimum de la fonction de contraste soit atteint sur un ensemble de *plusieurs* valeurs θ , mais que toutes ces valeurs correspondent à la même valeur $\gamma(\theta)$ en ce qui concerne la quantité d'intérêt considérée : auquel cas, on considèrera qu'on a bien un estimateur des moindres carrés défini de façon non ambiguë pour cette quantité d'intérêt malgré tout ! C'est notamment ce qui se produira dans le modèle de la *régression linéaire* (confer chapitre ??), où les paramètres de régression α et β auront des estimateurs de moindres carrés bien définis, mais pas le paramètre de bruit σ . ☺

Remarque (YF). L'estimateur des moindres carrés est très apprécié en analyse des données, car il ne fait intervenir que les « valeurs de régression » $\mathbb{E}_{\theta}(X_i)$, peu importe la loi précise suivie par les X_i : or, dans la mesure où il est difficile de modéliser pertinemment les détails de la loi des observations pour des données « de la vraie vie » un peu compliquées, cela constitue un avantage considérable ! ☺

Remarque (YG). En fait, l'estimateur des moindres carrés n'est qu'un cas particulier d'une méthode beaucoup plus générale, appelée « méthode de contraste », où le principe consistera à trouver d'autres « fonctions de contraste » vérifiant la propriété de la proposition (XY), puis à optimiser ces fonctions de contraste de la même façon que nous l'avions fait pour le contraste quadratique. Cependant, ces méthodes dépassent le bagage qu'on est raisonnablement en droit d'attendre d'un ingénieur généraliste, et nous ne les développerons donc pas dans ce cours ☹

Finissons cette section avec un exemple, en calculant l'estimateur des moindres carrés pour le modèle du chasseur :

Exemple (YH) (Moindres carrés pour le chasseur). Bien que l'observation du modèle du chasseur ne se décompose pas en sous-observations, les valeurs qu'elle prend sont réelles : on peut donc essayer d'appliquer la définition (YC) avec $n \leftarrow 1$ [attention, le « n » de la définition (YC) n'est pas le même que le « n » du modèle du chasseur : en l'occurrence, c'est le *premier* de ces deux « n » que nous prenons égal à 1 ici !]. On calcule facilement qu'on a $\mathbb{E}_{\theta}(X) = n\theta$ [††] ; par conséquent, la fonction de contraste est $(x, \theta) \mapsto (x - n\theta)^2$; et l'estimation des moindres carrés pour θ sera donc la valeur de θ minimisant $(x_{\checkmark} - n\theta)^2$: on tombe ainsi sur $\hat{\theta}_{\checkmark}^{\text{mc}} = x_{\checkmark} / n$, ce qui coïncide en l'occurrence avec l'estimateur « naïf » dont nous avons déjà parlé ! ☺

[††]. Où, cette fois-ci, « n » se réfère au nombre de tirs dans le test du candidat...

11.7 Estimation par tendance

Dans cette section nous allons présenter une méthode d'estimation très générale, que j'appellerai *méthode par tendance*. De manière assez étrange, cette méthode ne semble pas avoir de nom ni même être identifiée comme telle dans les ouvrages statistiques de référence ; elle est cependant au cœur de la très classique « méthode des moments » que nous présenterons un peu plus loin.

L'idée est la suivante : on va chercher une statistique $\hat{\gamma}(X)$ qui “ressemble” à la quantité d'intérêt $\gamma(\theta_{\mathcal{J}})$, au sens où, peu importe la valeur $\theta_{\mathcal{J}}$, son espérance (sous la véritable probabilité $\mathbb{P}_{\mathcal{J}}$) vaille $\gamma(\theta_{\mathcal{J}})$ (et, de préférence, que sa loi soit concentrée autant que possible autour de $\gamma(\theta_{\mathcal{J}})$) ; et on parlera alors d'« estimateur par tendance ».

! **Définition (YI)** (Estimateur par tendance). Soit un modèle statistique (fréquentiste) avec les notations habituelles, et soit $\hat{\gamma}(X)$ une statistique à valeurs réelles telle que

$$\forall \theta \in \Theta \quad \mathbb{E}_{\theta}(\hat{\gamma}(X)) = \gamma(\theta) \quad (\text{YJ})$$

pour une certaine fonction $\gamma: \Theta \rightarrow \mathbb{R}$. Alors $\hat{\gamma}(X)$ est qualifié d'*estimateur par tendance* de la quantité d'intérêt $\gamma(\theta)$. \heartsuit

Remarque (YK). On pourrait, pour le coup, étendre la notion d'« estimateur par tendance » en une notion de « prédicteur par tendance » ; cependant, ce ne serait pas très intéressant en pratique \heartsuit

Remarque (YL). En fait, dire que $\hat{\gamma}(X)$ est un estimateur par tendance de $\gamma(\theta)$ est rigoureusement synonyme à dire qu'il s'en agit d'un estimateur *sans biais*... (confer définition (WN)). Par rapport à la section 11.2, j'ai juste changé l'ordre dans lequel on raisonne : alors qu'à l'époque, l'idée était de dire « considérons un certain estimateur obtenu par telle ou telle méthode ; calculons sa fonction de biais : oh, l'estimateur est sans biais, c'est une bonne nouvelle ! », cette fois-ci, la démarche est : « considérons une quantité d'intérêt qu'on voudrait estimer ; si on arrive à trouver une statistique qui n'a pas de biais par rapport à cette quantité d'intérêt, alors il sera légitime de considérer qu'elle en constitue un estimateur ! » \heartsuit

Remarque (YM). En pratique, toute la difficulté est que la quantité d'intérêt $\gamma(\theta)$ est imposée par le contexte de notre analyse, et que nous devons trouver $\hat{\gamma}$ pour retomber précisément sur cette quantité d'intérêt ! Comme cela est particulièrement difficile à faire, nous devons souvent partir d'une statistique “inadaptée” $\hat{\psi}$, en déduire un estimateur par tendance d'une *autre* quantité d'intérêt $\psi(\theta)$, et utiliser ensuite l'idée de *substitution* que nous présenterons plus loin. \heartsuit

Exemple (YN). Voyons un exemple très simple d'estimateur par tendance pour le modèle du chasseur : le nombre de tirs X réussis lors du test est censé avoir pour moyenne $\mathbb{E}(\text{Binom}^{\text{le}}(n, \theta_{\mathcal{J}})) = \theta_{\mathcal{J}}n$; on a donc $\theta_{\mathcal{J}} = \mathbb{E}(n^{-1}X)$, ce qui suggère l'estimateur par tendance $n^{-1}X$ pour θ . Avec les paramètres pris en exemple, on estimera donc θ par 0,32. (Incidentement, on retombe en l'occurrence sur l'estimateur du maximum de vraisemblance). \heartsuit

Exemple (YO). Dans la modèle du pédagogue, on calcule que la variance empirique $\text{var}_{\text{emp}}(X_1, \dots, X_n)$ des observations (autrement dit, la quantité $n^{-1} \sum_{i=1}^n X_i^2 - (n^{-1} \sum_{i=1}^n X_i)^2$) a pour espérance $\frac{n-1}{n} \sigma_{\mathcal{J}}^2$; par conséquent, la variance bessélisée $\text{var}_{\text{B}}(X_1, \dots, X_n) := \frac{n}{n-1} \text{var}_{\text{emp}}(X_1, \dots, X_n)$ a pour espérance $\sigma_{\mathcal{J}}^2$: comme cette variance bessélisée est une statistique, elle constitue donc un estimateur par tendance de la variance σ^2 intrinsèque à la nouvelle méthode. \heartsuit

11.8 Technique de substitution

Principe de substitution

Nous venons de présenter un certain nombre de méthodes d'estimation ; mais en pratique, dès que le modèle devient un peu subtil, il est rare qu'aucune de ces méthodes nous donne directement (et de façon simple à calculer) accès à un estimateur de la quantité d'intérêt qui nous intéresse... C'est ici qu'intervient une idée aussi simple que puissante, appelée *substitution* :

Définition (YP) (Obtention de nouveaux estimateurs par substitution). Soit un modèle statistique avec les notations habituelles, pour lequel nous nous intéressons à la quantité d'intérêt $\gamma(\theta) \in \mathcal{G}$. Soient $(\hat{\psi}_i(X))_{i \in I}$ des estimateurs (qui peuvent avoir été obtenus par n'importe quelles techniques) des quantités d'intérêt respectives $(\psi_i(\theta))_{i \in I}$ (nous appellerons \mathcal{G}_i l'espace où vit la quantité d'intérêt $\psi_i(\theta)$) ; et supposons maintenant qu'il existe une fonction (pas trop biscornue) $F: \prod_{i \in I} \mathcal{G}_i \rightarrow \mathcal{G}$ telle que, pour tout $\theta \in \Theta$, on ait

$$\gamma(\theta) = F((\psi_i(\theta))_{i \in I}). \quad (\text{YQ})$$

On qualifie alors la statistique

$$\hat{\gamma}(X) := F((\hat{\psi}_i(X))_{i \in I}) \quad (\text{YR})$$

d'*estimateur par substitution* de $\gamma(\theta)$. ♡

Remarque (YS). J'ai écrit ce qui précède dans le cadre de l'estimation, mais le principe de substitution s'étend directement au cas des prédicteurs ! ♣

Remarque (YT). Ce qui rend la technique de substitution tellement puissante est que, au lieu d'être obligés de partir de notre quantité d'intérêt $\gamma(\theta)$ et de devoir "cuisiner" pour trouver *le bon* estimateur $\hat{\gamma}(X)$ qui se trouvera par miracle estimer précisément notre quantité d'intérêt, nous pouvons procéder "dans l'autre sens" en choisissant plus ou moins arbitrairement des statistiques $\hat{\psi}_1(X), \hat{\psi}_2(X), \dots$ qui sont faciles à calculer, en regardant de quelles quantités d'intérêt $\psi_1(\theta), \psi_2(\theta), \dots$ ces statistiques se trouvent être des estimateurs respectifs (en particulier, si on veut que $\hat{\psi}_i(X)$ soit un estimateur par tendance, la fonction ψ_i sera déterminée par l'égalité $\psi_i(\theta) = \mathbb{E}_\theta(\hat{\psi}_i(X))$; et, *seulement alors*, on regardera s'il y a un moyen de combiner les $\psi_i(\theta)$ pour en déduire notre quantité d'intérêt $\gamma(\theta)$ (via une relation du type $F(\psi_1(\theta), \psi_2(\theta), \dots) = \gamma(\theta)$), ce qui nous donnera alors par substitution l'estimateur $F(\hat{\psi}_1(X), \hat{\psi}_2(X), \dots)$ pour $\gamma(\theta)$. ♣

Exemple (YU). Voici un exemple d'utilisation de la substitution pour le modèle du pédagogue. On prend ici $n = 4$. (On a alors $\mathcal{X} = \mathbb{R}^4$). Nous allons considérer les deux fonctions suivantes sur \mathcal{X} :

- $\hat{\psi}_1(x_1, x_2, x_3, x_4) := \max_{i \in \{1, \dots, 4\}} x_i$, autrement dit la meilleure note obtenue parmi les quatre élèves ;
- $\hat{\psi}_2(x_1, x_2, x_3, x_4) := \max(\{x_1, x_2, x_3, x_4\} - \{\max_{i \in \{1, \dots, 4\}} x_i\})$, autrement dit la seconde meilleure note (qui peut éventuellement être identique à la première s'il y a ex-æquo).

On peut montrer que, dans le cas de notre modèle, lorsque X suit la loi Normale($\mu_{\mathcal{J}}, \sigma_{\mathcal{J}}^2$)^{⊗4}, les espérances respectives de $\hat{\psi}_1(X)$ et $\hat{\psi}_2(X)$ sont

$$\begin{aligned} \mathbb{E}_{\mathcal{J}}(\hat{\psi}_1(X)) &= \mu_{\mathcal{J}} + C_1 \sigma_{\mathcal{J}} ; \\ \mathbb{E}_{\mathcal{J}}(\hat{\psi}_2(X)) &= \mu_{\mathcal{J}} + C_2 \sigma_{\mathcal{J}} , \end{aligned}$$

où C_1 et C_2 sont des constantes numériques qu'on peut calculer et dont les valeurs respectives sont 1,029375 et 0,297011. Par conséquent, $\hat{\psi}_1(X)$ et $\hat{\psi}_2(X)$ sont des estimateurs par tendance resp. $\psi_1 := \mu + C_1\sigma$ et $\psi_2 := \mu + C_2\sigma$. Maintenant, supposons que notre quantité d'intérêt soit μ . On observe qu'on peut récupérer $\mu_{\mathcal{J}}$ à partir de $\psi_{1\mathcal{J}}$ et $\psi_{2\mathcal{J}}$ par la formule

$$\mu_{\mathcal{J}} = \frac{C_1}{C_1 - C_2} \psi_{2\mathcal{J}} - \frac{C_2}{C_1 - C_2} \psi_{1\mathcal{J}}. \quad (\text{YV})$$

Or, puisque nous savons estimer $\psi_{1\mathcal{J}}$ et $\psi_{2\mathcal{J}}$ par resp. $\hat{\psi}_1(X)$ et $\hat{\psi}_2(X)$, nous en déduisons par substitution l'estimateur suivant pour μ :

$$\hat{\mu}(X) := \frac{C_1}{C_1 - C_2} \hat{\psi}_2(X) - \frac{C_2}{C_1 - C_2} \hat{\psi}_1(X). \quad (\text{YW})$$

Cet estimateur illustre à merveille l'idée derrière la substitution, puisque nous avons *décidé* de la construire à partir de $\hat{\psi}_1(X)$ et $\hat{\psi}_2(X)$ (qui sont ce qu'on appelle les *statistiques d'ordre*, de co-rangs respectifs 1 et 2, de notre observation) n'est à vrai dire pas très performant ; mais il a l'avantage de ne demander que la connaissance des deux meilleurs notes et pas des autres ! ♣

Remarque (YX). L'estimateur présenté dans l'exemple ci-dessus semble assez grossièrement sous-optimal, dans la mesure où il n'utilise que les deux meilleures notes de la promotion : et de fait, ce n'est pas un estimateur ayant de très bonnes performances ! ∴ Néanmoins, ce genre d'estimateurs peut s'avérer très utiles dans certains domaines (je pense par exemple aux questions de performance sportive) où il est facile de trouver les données concernant les champions, mais difficile de trouver des données précises concernant la population générale... ♣

Méthode des moments

Un cas particulièrement important de technique de substitution, qui permet d'obtenir des estimateurs faciles à calculer dans un grand nombre de cas, est la *méthode des moments*, exposée ci-dessous :

Définition (YY). Supposons que nous soyons en présence d'un modèle relevant du point (XN) (càd. où l'observation est constituée de sous-observations qui sont de même loi sous les contextes fréquentistes), et pour lequel les sous-observations X_i sont en outre à valeurs réelles. Appelons k la *dimension* de l'espace Θ du paramètre caché de ce modèle, càd. le nombre de coordonnées dont on a besoin pour se repérer dans Θ [*]. Soit $\gamma := \gamma(\theta)$ une quantité d'intérêt : la *méthode des moments* consiste alors à estimer γ en procédant par substitution à partir des k premiers *moments empiriques* $\hat{\psi}_1, \dots, \hat{\psi}_k$ des observations, lesquels sont définis par

$$\hat{\psi}_r(\vec{X}_{\llbracket 0, n \rrbracket}) := \frac{1}{n} \sum_{i=0}^{n-1} X_i^r : \quad (\text{YZ})$$

le moment empirique $\hat{\psi}_r(\vec{x}_{\mathcal{J}})$ pouvant être vu comme une estimation ^[†] de (la variable aléatoire dont la réalisation est) $E_{\mathcal{J}}(X_0^r)$, qui est ce qu'on appelle le *r-ième moment* de la loi $\text{Loi}_{\mathcal{J}}(X_0)$. ♣

[*]. Par exemple, la surface de la Terre est de dimension 2 puisqu'on peut s'y repérer par les latitude et longitude, et ce, malgré que la Terre n'est pas plane.

[†]. En l'occurrence, il s'agit à la fois d'une estimation empirique et d'une estimation par tendance.

Remarque (AA'). Dans certains cas, il peut arriver qu'à cause des symétries du problème, un des moments de $\text{Loi}_{\mathcal{J}}(X_0)$ soit en fait une fonction des moments précédents (et ce, indépendamment de la valeur $\theta_{\mathcal{J}}$, s'entend : dans ce cas, ce moment n'apporte aucune information pour la substitution, et il faut donc le "sauter" — et aller alors jusqu'au $(k + 1)$ -ième moment. ♣

Remarque (AB'). Je n'exigerai pas de vous que vous reteniez par cœur la méthode des moments : cependant, comme il s'agit d'un outil très classique en statistique, il est bon que vous l'ayez vue à l'occasion de ce cours, et que vous sachiez la reconnaître si elle intervient dans vos exercices ! ☺ ♣

Exemple (AC'). Voyons ce que donne la méthode des moments dans l'exemple du pédagogue, en prenant pour quantité d'intérêt $\mu - \sigma$. Le premier moment empirique de nos observations est

$$\hat{\psi}_1(\vec{X}) := \frac{1}{n} \sum_{i=0}^{n-1} X_i = \text{moy}(X_0, \dots, X_{n-1}) : \quad (\text{AD}')$$

il s'agit d'un estimateur empirique de la quantité d'intérêt (dont la réalisation est) $\mathbb{E}_{\mathcal{J}}(X_0) = \mathbb{E}(\text{Normale}(\mu_{\mathcal{J}}, \sigma_{\mathcal{J}}^2)) = \mu_{\mathcal{J}}$. (C'en serait aussi un estimateur par tendance, du reste).

Le second moment empirique des nos observations est

$$\hat{\psi}_2(X) := \frac{1}{n} \sum_{i=1}^n X_i^2 = \text{var}_{\text{emp}}(X_1, \dots, X_n) + \text{moy}(X_1, \dots, X_n)^2 : \quad (\text{AE}')$$

il s'agit d'un estimateur empirique (et par tendance) de $\mathbb{E}_{\mathcal{J}}(X_1^2) = \mathbb{E}(\text{Normale}(\mu_{\mathcal{J}}, \sigma_{\mathcal{J}}^2)^2) = \mu_{\mathcal{J}} + \sigma_{\mathcal{J}}^2$.

Nous avons donc trouvé deux estimateurs des quantités d'intérêt resp. μ et $\mu + \sigma^2$, à savoir, $\hat{\psi}_1$ et $\hat{\psi}_2$. Nous voulons maintenant procéder par *substitution* en trouvant une moyen de "combiner" ces quantités d'intérêt pour obtenir $\mu - \sigma$ (qui est la quantité d'intérêt qui nous intéresse vraiment). Dit en termes de réalisation, il s'agit donc de trouver un moyen de déduire $\mu_{\mathcal{J}} - \sigma_{\mathcal{J}}$ à partir de $\mu_{\mathcal{J}}$ et $\mu_{\mathcal{J}} + \sigma_{\mathcal{J}}^2$. (Pour rendre les notations plus efficaces, dans la suite, nous poserons resp. $\mu_{\mathcal{J}} =: \psi_{1\mathcal{J}}$ et $\mu_{\mathcal{J}} + \sigma_{\mathcal{J}}^2 =: \psi_{2\mathcal{J}}$).

En substance, il s'agit donc d'"inverser" la fonction $\Theta \ni (\mu, \sigma) \mapsto (\mu, \mu + \sigma^2)$: comme aussi bien l'espace de départ que l'espace d'arrivée sont de dimension 2, cela a à priori de bonnes chances d'être possible (et, de fait, ce sera bien le cas ☺). On trouve immédiatement que $\sigma_{\mathcal{J}}^2 = \psi_{2\mathcal{J}} - \psi_{1\mathcal{J}}^2$, d'où $\sigma_{\mathcal{J}} = (\psi_{2\mathcal{J}} - \psi_{1\mathcal{J}}^2)^{1/2}$ (puisque $\sigma_{\mathcal{J}}$ est positif), et finalement $\mu_{\mathcal{J}} - \sigma_{\mathcal{J}} = \psi_{1\mathcal{J}} - (\psi_{2\mathcal{J}} - \psi_{1\mathcal{J}}^2)^{1/2}$. Par substitution, on en déduit finalement notre estimateur des moments pour $\mu - \sigma$: c'est

$$\hat{\gamma}(X) := \hat{\psi}_1(X) - (\hat{\psi}_2(X) - \hat{\psi}_1(X)^2)^{1/2} = \text{moy}(X_1, \dots, X_n) - \text{var}_{\text{emp}}^{1/2}(X_1, \dots, X_n). \quad (\text{AF}')$$

(Noter qu'en l'occurrence, on voit qu'il s'agit aussi de l'estimateur empirique obtenu à partir de la fonctionnelle « moyenne moins écart-type »). ♣

Chapitre 12

Tests d'hypothèses

12.1 Motivation du concept

Exemple introductif

Dans cette sous-section, nous reprenons le contexte et les notations du modèle du pédagogue : en particulier, nous nous intéressons à la nouvelle méthode pédagogique mise au point par un enseignant, et à l'impact de celle-ci sur les résultats des élèves. À ce stade, considérons que la nouvelle méthode vient juste d'être conçue, et n'a pas encore été testée en conditions réelles. Notre enseignant expose sa méthode à ses collègues, et un désaccord entre les deux camps apparaît :

- D'un côté, notre enseignant pense que la nouvelle méthode va globalement faire progresser les élèves par rapport à l'ancienne (ou, au pire, ne rien changer) : autrement dit, il pense que $\mu_{\mathcal{V}}$, la véritable espérance des scores pour sa nouvelle méthode est supérieure ou égale à $\mu_{\text{réf}}$;
- À l'inverse, les collègues estiment que la nouvelle méthode va s'avérer strictement nuisible : autrement dit, selon eux, la valeur $\mu_{\mathcal{V}}$ est donc strictement inférieure à $\mu_{\text{réf}}$!

Par ailleurs, dans le cadre de cette section, nous allons considérer qu'aussi bien l'enseignant que ses collègues sont d'accord pour dire que la variabilité des résultats d'un élève à l'autre ne sera pas impactée par le changement pédagogique : autrement dit, tous sont d'accord pour dire que la valeur $\sigma_{\mathcal{V}}$ est égale à $\sigma_{\text{réf}}$. De la sorte, l'espace dans lequel vit le paramètre caché (μ, σ) n'est pas l'espace originel $\mathbb{R} \times \mathbb{R}_+^*$ (que nous noterons dorénavant Θ^{orig} lorsqu'il faudra éviter toute ambiguïté), mais le sous-espace $\mathbb{R} \times \{\sigma_{\text{réf}}\} =: \Theta^{\text{hsc}} [^*]$, que nous noterons simplement Θ dans la suite de cette sous-section. (Gardez néanmoins à l'esprit que ce n'est pas exactement le même ' Θ ' que dans l'exemple (MC) d'origine : ici, seul μ est réellement inconnu).

La question qui se pose est : comment trancher le débat ? Par l'expérience, évidemment ! On va expérimenter la nouvelle méthode sur une promotion, et s'appuyer sur les résultats obtenus par les élèves pour dire si on a plutôt envie de pencher dans le sens de l'enseignant, ou plutôt dans celui de ses collègues. Mais comment, concrètement, procéder : quel traitement précis allons-nous appliquer aux notes pour déterminer qui a raison ?

Dans le contexte considéré, où les deux thèses soutenues ne diffèrent que par le résultat moyen espéré pour les élèves, il est intuitivement évident que la "bonne"

[*]. 'hsc' est une abréviation pour « homoscédastique » : ce que signifie cet adjectif vous sera expliqué plus tard.

façon de procéder pour trancher entre ces deux thèses devrait être la suivante : on va regarder quelle est la moyenne $m_{\mathcal{J}}$ effectivement obtenue par les élèves [il s'agit donc de la réalisation de la statistique $M := (\sum_{i=0}^{n-1} X_i) / n$] : plus $m_{\mathcal{J}}$ sera élevée, plus on aura envie de trancher en faveur de l'enseignant ; et plus elle sera basse, plus on aura envie de trancher en faveur de ses collègues ! À ce stade, on a même envie d'ajouter (mais nous ne le ferons pas) qu'on voudrait trancher en faveur de l'enseignant si $m_{\mathcal{J}}$ s'avère plus grande que $\mu_{\text{réf}}$, et en sa défaveur dans le cas contraire...

... Cela dit, autant notre enseignant est d'accord avec l'idée d'utiliser la moyenne de la promotion pour jauger l'efficacité de sa méthode, autant il considère que fixer le seuil à $\mu_{\text{réf}}$ n'est pas la bonne chose à faire. « En effet, explique-t-il, une moyenne effective de 98 serait certes plus en phase avec la théorie des collègues qu'avec la mienne ; mais après tout, elle peut aussi tout à fait s'expliquer par le fait qu'on ait en réalité $\mu_{\mathcal{J}} = 103$, et que le hasard ait fait qu'on soit tombé sur un promotion un peu faible cette année-là : dès lors, je refuserais de remettre en cause ma méthode sur foi de ce seul résultat... ! ». Têtu, notre enseignant ? Assurément ! ☺ Néanmoins, ses collègues reconnaissent que, en l'occurrence, il est effectivement légitime de traiter les deux thèses de façon *dissymétrique*. En effet, *in fine*, le point n'est pas seulement de savoir qui a raison, mais de savoir quelle méthode pédagogique il convient d'utiliser. Or, à supposer que la nouvelle méthode pédagogique soit bel et bien meilleure, l'abandonner à tort porterait un préjudice à *toutes* les générations futures d'élèves qui auraient pu en bénéficier ! Alors que si la nouvelle méthode est moins bonne, mais qu'on reconduit l'expérimentation quelques années avant de l'abandonner définitivement, cela aura certes été nuisible aux quelques promotions "cobayes", mais c'est tout de même moins grave... ☹ Bref ; pour prendre en compte ces considérations, on souhaite avoir une procédure qui conclura en faveur des collègues *seulement* si on a une *forte conviction* qu'ils ont raison, et en faveur de l'enseignant dès lors qu'il reste *plausible*, au vu des notes des élèves, que ce soit lui qui ait raison ☺

Après discussion, l'enseignant et ses collègues choisissent de fixer leur seuil de décision à $95 =: \mu_{\text{seuil}}$: si la moyenne de la promotion s'avère supérieure ou égale cette valeur, ils laisseront l'enseignant continuer d'expérimenter sa nouvelle méthode (fût-ce "au bénéfice du doute" si la moyenne est comprise entre 95 et 100) ; mais si elle s'avère strictement inférieure, l'enseignant conviendra que c'est une preuve forte que sa méthode est en fait mauvaise, et qu'il convient donc d'arrêter là les frais... ! Pour fixer ce seuil à 95, l'enseignant et ses collègues ont raisonné ainsi. « Ce qu'on veut, c'est que si la thèse de l'enseignant est vraie, on n'ait qu'un risque très faible de se tromper en jugeant en sa défaveur. Plaçons-nous dans le cas où la thèse de l'enseignant ne serait vraie que "de justesse", avec, mettons, $\mu_{\mathcal{J}} = 100,1$. Dans ce cas, la v.a. $M := (X_1 + \dots + X_n) / n$ suivra (sous le contexte probabiliste $\mathbb{P}_{\mathcal{J}}$ correspondant au $\mu_{\mathcal{J}}$ en question) la loi Normale($\mu_{\mathcal{J}}, \sigma_{\text{réf}}^2 / n$). Or la probabilité qu'une telle v.a. soit inférieure à μ_{seuil} peut être calculée comme valant environ 5 % :

```
> n = 22
> sigma = 15
> muref = 100
> mucoche = muref + 0.1
> museuil = 95
> pnorm(museuil, mucoche, sqrt(sigma ^ 2 / n))
[1] 0.05538504
```

Ainsi, si notre enseignant a effectivement raison, il y a au pire environ 5 % de risque qu'on rejette erronément sa thèse. C'est faible, ce qui est voulu : en effet, comme nous l'avons fait remarquer ci-dessus, on veut avoir une conviction *forte* que ce sont les collègues qui ont raison avant de trancher en leur faveur ; or, on n'aurait pas pu raisonnablement qualifier notre conviction de "forte" s'il y avait eu, mettons, 15 % de risque de rejeter erronément la thèse de l'enseignant. D'un autre côté, ce seuil n'est pas non plus *trop* faible : ce qui est souhaitable aussi, car il faut tout de même donner une chance à la procédure de trancher en faveur des collègues si ce sont eux qui ont raison ! ». [Ainsi, on aurait certes pu fixer le seuil à 90, auquel cas il y aurait eu moins d'un risque sur 1 000 de rejeter à tort la thèse de l'enseignant ; mais ç'aurait été là requérir un niveau de preuve peu raisonnable : pas suffisamment raisonnable en tout cas pour justifier qu'on continue d'accorder à l'enseignant le bénéfice du doute pour une moyenne comprise entre 90 et 95, dans la mesure où une moyenne aussi faible à peine 5 % de chances d'arriver si l'enseignant a raison...].

Introduction du vocabulaire

Dans cette sous-section nous allons reprendre la situation discutée dans la sous-section précédente, mais en soulignant explicitement cette fois-ci les concepts qui y ont été utilisés sans le dire, et en introduisant le vocabulaire que les statisticiens utilisent pour désigner ces concepts :

- ☛ L'affirmation de l'enseignant et celle de ces collègues correspondent à des *hypothèses* (ce vocabulaire a déjà été introduit dans la § 9.3), autrement dit, à des quantités d'intérêt explicatives booléennes. Plus précisément, l'hypothèse défendue par l'enseignant est que $\{\mu \geq \mu_{\text{réf}}\}$, tandis que celle défendue par ses collègues est que $\{\mu < \mu_{\text{réf}}\}$. Dans la suite, nous noterons ces hypothèses resp. \mathcal{H}_0 et \mathcal{H}_1 .

- ☛ On a plus précisément considéré deux hypothèses *complémentaires* l'une de l'autre : \mathcal{H}_0 et \mathcal{H}_1 sont chacune la négation de l'autre (ce qu'on note formellement « $\mathcal{H}_{1-i} = \neg \mathcal{H}_i$ »).

- ☛ L'hypothèse de l'enseignant est celle qu'on veut ne récuser que lorsqu'on a de *fortes* preuves en sa *défaveur* : dans ce cas, on la qualifie d'*hypothèse nulle* (d'où la notation ' \mathcal{H}_0 ').^[†]

- ☛ Complétement, l'hypothèse des collègues est celle qu'on ne veut valider que lorsqu'on a de *fortes* preuves en sa *faveur* : dans ce cas, on la qualifie d'*hypothèse alternative*. [Les hypothèses alternatives sont généralement notées à l'aide de l'indice '1'].

- ☛ La moyenne M des élèves (dont m_{\checkmark}) est la réalisation est l'indicateur qu'on a choisi de regarder pour trancher entre des hypothèses, indicateur qui correspond à une *statistique* (puisque'il est obtenu directement à partir des données) à valeurs *réelles*. Un tel indicateur est qualifié de *statistique de test*.

- ☛ Ici on choisit de trancher en faveur de l'hypothèse alternative lorsqu'on se situe *en-deçà* d'un certain seuil : on dit dans ce cas qu'on procède à un « *test à gauche* ». (Si on avait au contraire utilisé un seuil pour trancher en direction de l'hypothèse alternative *au-delà* dudit seuil, on aurait parlé de « *test à droite* »). De

[†]. La raison que nous avons invoquée dans la § 12.1 pour traiter de façon dissymétrique les hypothèses nulle et alternative est l'idée qu'il serait plus grave de rejeter la thèse de l'enseignant à tort, que de l'accepter à tort. En réalité, c'est une raison possible au fait de traiter les hypothèses nulle et alternative de façon dissymétrique, mais en réalité la raison *fondamentale* pour laquelle la théorie des tests est fondamentalement dissymétrique n'est pas celle-là ! Cela sera discuté ultérieurement, conférer § ??.

manière équivalente, on pourra dire que notre « *critère de suspicion* » consiste à considérer les valeurs de la statistique de test comme d'autant plus suspectes (par rapport à l'hypothèse nulle, s'entend) qu'elles sont situées vers la gauche.

☛ La procédure à laquelle on est arrivé correspond à trancher soit dans un sens, soit dans l'autre, selon ce que vaut l'observation (via la valeur de la statistique de test qui s'en déduit). On peut donc la décrire comme une statistique à valeurs dans $\{\text{VRAI}, \text{FAUX}\}$, où FAUX signifie qu'on penche vers l'hypothèse nulle, tandis que VRAI signifie qu'on penche vers l'hypothèse alternative^[‡] : formellement, cette statistique s'écrirait donc $\{M < \mu_{\text{seuil}}\}$. Une telle statistique booléenne est qualifiée de *test (booléen)*.

☛ En fait, ce qu'on qualifie de « test (booléen) », c'est plutôt la *procédure de jugement* décrite par la variable aléatoire que vous venons d'évoquer, plutôt que la variable aléatoire elle-même : autrement dit, notre « test » consiste ici à « conclure en faveur des collègues si $M < \mu_{\text{seuil}}$, et en faveur de l'enseignant sinon ».

☛ Lorsqu'on conclut en faveur de l'hypothèse alternative, on dit qu'on *rejette l'hypothèse nulle* (dans la mesure où on a trouvé un résultat qui n'est vraiment pas compatible avec celle-ci) ; on peut aussi dire, éventuellement, qu'on *valide l'hypothèse alternative* (au sens où on a apporté une quasi-preuve de ce que c'est elle qui est correcte). Dans un tel cas, on dit que le test est *positif* (au sens où sa réalisation vaut VRAI).

☛ Lorsque, au contraire, on conclut en faveur de l'hypothèse nulle, on dit qu'on *accepte* l'hypothèse nulle. L'emploi du mot « accepter » signifie que cela n'implique pas pour autant qu'il *faill*e rejeter l'hypothèse alternative, juste que l'hypothèse nulle est acceptable. Pour bien souligner cet aspect non exclusif, on peut aussi dire qu'on « ne rejette pas » l'hypothèse nulle, ou qu'on « n'est pas en mesure de valider » l'hypothèse alternative^[§]. Dans ce cas, on dit que le test est *négatif*^[¶].

☛ La valeur d'environ 5 % que nous avons trouvée correspond à la probabilité qu'on a de conclure positivement (donc en rejetant l'hypothèse nulle) lorsque l'hypothèse nulle est vraie, et plus précisément, dans le « pire » cas où l'hypothèse nulle est vraie [ou presque : pour nos calculs, nous avons en effet pris une valeur de 100,1, alors que le pire cas aurait été de prendre *précisément* $\mu_{\checkmark} = \mu_{\text{réf}}$]. C'est (modulo le fait que nous n'avons pas *exactement* considéré le pire cas) il s'agit de ce que les statisticiens qualifient de *niveau* du test.

[‡]. Mnémotechniquement, rappelons qu'il est d'usage d'associer FAUX et VRAI aux valeurs respectives 0 et 1.

[§]. Attention : « Ne pas être en mesure de valider » \mathcal{H}_1 , cela ne signifie pas pour autant qu'on *invalide* cette hypothèse : juste qu'on n'est pas en état d'en avoir une confirmation forte !

[¶]. Rappelez-vous que l'idée est de voir si on a assez d'éléments pour rejeter l'hypothèse nulle avec conviction. Vu sous cet angle, le fait de ne pas être en mesure de rejeter est une *absence* de conclusion ferme, d'où l'adjectif « négatif » ; tandis que le fait de rejeter est une affirmation désignant une véritable conclusion, d'où l'adjectif « positif ». Par ailleurs, il est logique de parler de résultat « positif » quand on penche en faveur de \mathcal{H}_1 (puisque 1 est un nombre strictement positif !), resp. de résultat « négatif » quand on penche en faveur de \mathcal{H}_0 (qui, au sens large, est bien un nombre négatif 😊).

12.2 Formalisation du concept dans le cas unilatéral

Définition formelle

Le moment est à présent venu de définir formellement la notion de test (booléen) :

Définition (AG') (Test booléen). Avec les notations standard, soit $\Theta_1 \subseteq \Theta$ un sous-ensemble de l'espace du paramètre caché, permettant de définir une hypothèse associée $\{\theta \in \Theta_1\} =: \mathcal{H}_1$. Dans ces conditions, on appelle *test de l'hypothèse alternative* \mathcal{H}_1 une statistique à valeurs booléenne, autrement dit une v.a. *Test* de la forme $test(X)$ à valeurs dans $\{\text{VRAI}, \text{FAUX}\}$, qui “essaye” de dire si l'hypothèse \mathcal{H}_1 est correcte ou non.

De façon équivalente, on dit aussi que *Test* est un *test de l'hypothèse nulle* \mathcal{H}_0 , où \mathcal{H}_0 est l'hypothèse complémentaire de \mathcal{H}_1 , autrement dit, $\mathcal{H}_0 := \neg \mathcal{H}_1 = \{\theta \in \Theta_0\}$, où $\Theta_0 := \Theta \setminus \Theta_1$ est le sous-ensemble complémentaire de Θ_1 .

Lorsque la réalisation de $test(x_{\checkmark})$ du test vaut VRAI, on interprète ce résultat en disant que « l'hypothèse nulle est rejetée », ou que « l'hypothèse alternative est validée », ou encore que « le test est positif ». Lorsque le résultat du test vaut FAUX, on interprète ce résultat en disant que « l'hypothèse nulle est acceptée », ou que « l'hypothèse alternative demeure indécise », ou encore que « le test est négatif ».

On appelle *niveau* du test (traditionnellement noté ' α ') la probabilité supremale que le test soit positif alors que l'hypothèse nulle est vraie :

$$\alpha := \sup_{\theta \in \Theta_0} \mathbb{P}_{\theta}(Test) \quad \text{[III]}. \quad (\text{AH}')$$

On dit que le test est « de niveau α » lorsque son niveau est $\leq \alpha$ [**]. ♡

Cette définition appelle un certain nombre de commentaires :

Remarque (AI'). Les hypothèses nulle et alternative étant par définition complémentaires l'une de l'autre, en général on ne définira qu'une des deux lorsqu'on parlera d'un test, l'autre s'en déduisant automatiquement. ♣

Remarque (AJ'). Comme vous le voyez, les hypothèses nulle et alternative jouent des rôles bien distincts au niveau de l'interprétation du résultat et de la définition

[III]. Notez que, dans la mesure où v.a. *Test* est *booléenne*, on peut bien parler de la probabilité de la v.a. *elle-même*, puisque dans un tel cas l'évènement *Test* est strictement équivalent à $\{Test = \text{VRAI}\}$: en effet, l'évènement $\{Test = \text{VRAI}\}$ est l'évènement qui est vrai lorsque *Test* vaut VRAI, et faux sinon (sachant que, en tant que v.a. booléenne, lorsque *Test* ne vaut pas VRAI, elle vaut nécessairement FAUX) : par conséquent, cet évènement vaut VRAI lorsque *Test* vaut VRAI, resp. FAUX lorsque *Test* vaut FAUX : c'est donc bien *Test* elle-même ! Bien sûr, si cela vous dérange trop d'écrire une formule avec une variable aléatoire “nue”, sans signe de relation, vous pouvez écrire « $\mathbb{P}_{\theta}(Test = \text{VRAI})$ » plutôt que « $\mathbb{P}_{\theta}(Test)$ » à l'intérieur d'une probabilité ; mais notez l'aspect pléonastique de la chose ! N'oubliez pas en effet que, même si on ne vous l'a pas dit ainsi dans vos cours précédents, un évènement n'est autre qu'une variable aléatoire booléenne : on peut certes *aussi* le voir comme un sous-ensemble de l'univers Ω ; néanmoins, il est plus cohérent pour l'intuition de dire que, quand on parle de l'évènement {la somme des deux dés vaut au moins 10}, on parle bien de quelque chose qui vaut soit VRAI, soit FAUX, selon ce que vaut la somme des deux dés : c'est là une façon de voir les choses qui est beaucoup plus naturelle que de prétendre qu'on serait en fait en train de parler de l'ensemble des éventualités ω de l'univers pour lesquelles la valeur de la variable aléatoire « somme des deux dés » atteint ou dépasse 10...

[**]. Attention donc : dire « le niveau du test est α » n'est pas exactement synonyme de « le test est de niveau α ».

du niveau. Par conséquent, veillez à ne jamais parler de « test de l'hypothèse \mathcal{H} », mais à toujours préciser si \mathcal{H} joue le rôle de l'hypothèse nulle ou de l'hypothèse alternative : en disant soit « test de l'hypothèse *nulle* \mathcal{H} » (ce qui est équivalent à « test de l'hypothèse alternative $\neg\mathcal{H}$ »), soit « test de l'hypothèse *alternative* \mathcal{H} » (ce qui est équivalent à « test de l'hypothèse nulle $\neg\mathcal{H}$ »). ♣

!! **Remarque (AK')** (Dissymétrie des conclusions d'un test). Attention ! « Accepter » l'hypothèse nulle ne signifie pas qu'on a de quoi affirmer avec confiance que l'hypothèse nulle est vraie : cela signifie juste qu'on trouve qu'il est *acceptable* de croire en l'hypothèse nulle au vu des résultats obtenus — mais peut-être que l'hypothèse alternative est tout aussi acceptable ! À l'inverse, « rejeter » l'hypothèse nulle signifie qu'on trouve que cette hypothèse n'est vraiment pas crédible au vu des données. Il y a donc une forte dissymétrie : « accepter l'hypothèse nulle » n'est absolument pas synonyme de « rejeter l'hypothèse alternative » ! Pour cette raison, certains auteurs préfèrent écrire qu'on « ne rejette pas » l'hypothèse nulle (c'est pourquoi on parle de test *néгатif* dans ce cas), le terme « accepter » étant jugé trop ambigu. ♣

! **Remarque (AL')**. Notez bien que, dans la formule définissant le niveau d'un test, seule l'hypothèse *nulle* intervient ! Ce qui se passe dans le cas où θ_{\checkmark} est dans Θ_1 n'est tout simplement pas considéré ici. ♣

!! **Remarque (AM')**. En fait, la remarque précédente provient de ce que, en substance, le test répond à la question « l'observation est-elle *compatible* avec l'hypothèse nulle », *sans* se poser la question de l'hypothèse alternative ! (Ou alors, seulement de façon indirecte, via le choix d'une statistique de test pertinente : confer § 12.4 *infra*). ♣

Remarque (AN'). La notion de niveau du test est fondamentalement *fréquentiste* : et, de manière générale, toute la théorie des tests d'hypothèse nulle se place entièrement dans le paradigme fréquentiste ! On peut d'ailleurs comparer la notion de niveau d'un test avec le « niveau bayésien » dans le cadre bayésien (définition (RZ)) :

- Le niveau *bayésien* était le supremum, sur les valeurs possibles de l'*observation* correspondant à un test *positif*, de la probabilité que le test se trompe (autrement dit, qu'on soit en réalité sous l'hypothèse *nulle*), probabilité qui est considérée dans le contexte *sachant la valeur correspondante de l'observation*. (Il s'agit donc d'un contexte de nature bayésienne : le contexte probabiliste $\mathbb{P}(\bullet \mid X = x)$ étant, en quelque sorte, « la probabilité à postériori dans le cas où l'observation vaut x »).
- Le niveau *fréquentiste* est aussi un supremum de la probabilité que le test se trompe en répondant « positif » dans un cas où c'est l'hypothèse nulle qui est vraie : néanmoins, cette fois-ci on considère le supremum sur les valeurs possibles du *paramètre caché* correspondant à l'hypothèse *nulle*, et la probabilité considérée est celle que le test réponde *positif*, probabilité qui est considérée dans le contexte *sachant la valeur correspondante du paramètre caché*. (Il s'agit donc d'un contexte de nature *fréquentiste* : c'est le contexte \mathbb{P}_{θ} , qui correspond, en quelque sorte, à « la véritable probabilité dans le cas où le paramètre caché vaut θ »).

Ainsi, en passant du paradigme bayésien au paradigme fréquentiste, on a été obligés de considérer des contextes probabilistes différents (dans la mesure où on ne peut plus parler du contexte à postériori), ce qui nous a amenés à inverser les rôles joués par le paramètre caché et l'observation. ♣

Remarque (AO'). Ici j'ai simplement qualifié α de « niveau » ; mais il faudrait en toute rigueur parler de « niveau *de risque* » : et on utilise, à l'inverse, l'expression « niveau *de confiance* » pour désigner $1 - \alpha$. Néanmoins, en pratique, on choisit toujours des niveaux de risque inférieurs ou égaux à $1/2$, de sorte qu'il n'y a aucun

risque d'ambiguïté : si on parle d'un test « de niveau 92 % », cela se référera forcément à un niveau de confiance (auquel cas $\alpha = 8\%$); alors que, lorsqu'on parle d'un test « de niveau 20 % », il s'agira à l'inverse d'un niveau *de risque*.

À noter, au passage, que la notation traditionnelle ' α ' sous-entend qu'on parle d'un niveau *de risque* : les niveaux de confiance, eux, sont traditionnellement notés $1 - \alpha$. \clubsuit

Remarque (AP'). Pour déterminer le niveau d'un test, il est essentiel d'être capable de définir l'hypothèse nulle \mathcal{H}_0 sous la forme $\{\theta \in \Theta_0\}$ (ce qui est toujours possible, par définition de ce qu'est une hypothèse), puisque le sous-espace Θ_0 intervient dans le supremum définissant le niveau. Et notez bien que ce supremum portera sur le paramètre caché *lui-même*, par sur une quantité d'intérêt explicative : par exemple, si nous prenons le modèle du pédagogue avec l'espace originel $\Theta^{\text{orig}} = \mathbb{R} \times \mathbb{R}_+^*$ du paramètre caché, et que nous considérons l'hypothèse nulle $\{\sigma \geq 0,7 \sigma_{\text{réf}}\}$, alors le supremum pour la définition du niveau de sera pas un supremum sur la seule variable σ , mais bien sur le *couple* (μ, σ) constituant le paramètre caché, quand bien même μ n'intervient pas explicitement dans la définition de l'hypothèse nulle ! — Plus précisément, on aura alors que

$$\alpha \stackrel{\text{déf}}{=} \sup \{ \mathbb{P}_{\mu, \sigma}(Test) \mid \mu \in \mathbb{R}, \sigma \geq 0,7 \sigma_{\text{réf}} \}. \quad (\text{AQ}')$$

\clubsuit

Remarque (AR'). Dans la suite du cours, il sera fréquent que j'introduise seulement les notations \mathcal{H}_0 et \mathcal{H}_1 , puis que je me mette à utiliser les notations ' Θ_0 ' et ' Θ_1 ' sans les avoir définies formellement : il faudra alors comprendre que Θ_i est le sous-ensemble de Θ caractérisé par le fait que $\mathcal{H}_i = \{\theta \in \Theta_i\}$. Parfois, l'association implicite se fera en sens inverse : je n'introduirai formellement que les ensembles Θ_0 et Θ_1 , puis utiliserai les notations ' \mathcal{H}_0 ' et ' \mathcal{H}_1 ' sans prendre le temps de préciser que \mathcal{H}_i est défini par $\{\theta \in \Theta_i\}$.

En fait, j'assimilerai très même souvent l'hypothèse en tant que telle et le sous-ensemble de Θ correspondant : il pourra ainsi m'arriver d'écrire des choses comme « \mathcal{H}_0 est topologiquement fermée », auquel cas je voudrai en réalité que c'est l'ensemble Θ_0 qui est un sous-ensemble fermé de Θ . \clubsuit

Remarque (AS'). Notez que, de même que le concept d'estimateur (ou de prédicteur), le concept de « test » est très creux au niveau de sa définition « nue » : « être un test », en tant que tel, cela ne présente quasiment aucun intérêt... Ce qui est intéressant, c'est d'être un test *pertinent* ! De ce point de vue, « être de niveau α » est un tel critère de pertinence, de même que « être convergent » était un critère de pertinence pour un estimateur. Nous verrons néanmoins que le fait d'atteindre un niveau donné n'est pas en soi un critère de qualité suffisant pour un test : ce qui compte aussi, c'est qu'il ait une bonne *puissance*... Ce concept de « puissance » sera introduit un peu plus loin (définition (AU')). \clubsuit

Puissance d'un test

Nous avons expliqué que, lorsqu'on procède à des tests fréquentistes, les hypothèses nulle et alternative sont traitées de façon dissymétrique : le « niveau » du test, en particulier, se réfère uniquement au cas où on se trompe en répondant VRAI alors qu'on est en réalité dans l'hypothèse nulle. Néanmoins, si on ne considérait que ce critère, la théorie des tests deviendrait complètement triviale : il suffirait de

ne considérer que les tests répondant systématiquement FAUX, et ce serait imbattable...! ☹

Bien entendu, en pratique, on veut *non seulement* que notre test ne se trompe pas trop souvent lorsqu'on est dans l'hypothèse nulle, mais *aussi* qu'il parvienne à identifier correctement l'hypothèse alternative dans certains cas... Ce sont les notions d'*erreur de seconde espèce* et de *puissance*, que nous allons à présent présenter, qui permettent de quantifier cela.

Définition (AT') (Erreurs de première et de seconde espèce). On dit qu'un test commet une *erreur de première espèce* lorsqu'il tranche en faveur de l'hypothèse alternative alors que c'est l'hypothèse nulle qui était correcte ; et qu'il commet une *erreur de seconde espèce* lorsqu'il tranche en faveur de l'hypothèse nulle alors que c'est l'hypothèse alternative qui était correcte. ♥

! **Définition (AU')**. Lorsqu'on utilise un test booléen *Test* pour tester une hypothèse nulle $\{\theta \in \Theta_0\}$ contre l'hypothèse alternative $\{\theta \in \Theta_1\}$ (avec $\Theta_1 := \Theta \setminus \Theta_0$) :

- La *fonction de risque de première espèce* est l'application, définie sur Θ_0 , qui associe à une valeur de θ le risque de première espèce correspondant : $\Theta_0 \ni \theta \mapsto \mathbb{P}_\theta(\text{Test})$.
- En pratique, on s'intéresse uniquement au *supremum* de la fonction de risque de seconde espèce, qu'on appelle *niveau* du test, et qu'on note traditionnellement α (on retombe donc sur la définition (AG')) :

$$\alpha := \sup\{\mathbb{P}_\theta(\text{Test}) \mid \theta \in \Theta_0\}. \quad (\text{AV}')$$

- La *fonction de risque de seconde espèce* est l'application, définie sur Θ_1 , qui associe à une valeur de θ le risque de seconde espèce correspondant. On utilise généralement la lettre β pour la désigner :

$$\Theta_1 \ni \theta \mapsto \beta(\theta) := \mathbb{P}_\theta(\neg \text{Test}). \quad (\text{AW}')$$

On considère également fréquemment le complément à 1 de la fonction de risque de seconde espèce, qu'on appelle *fonction de puissance* du test :

$$\Theta_1 \ni \theta \mapsto 1 - \beta(\theta) = \mathbb{P}_\theta(\text{Test}). \quad (\text{AX}')$$

- En pratique, contrairement à ce qu'on fait pour le risque de première espèce, le supremum du risque de seconde espèce (ou, de manière équivalente, l'infimum de la fonction de puissance) n'est jamais considéré. ♥

Remarque (AY'). Vous vous demandez certainement pourquoi on considère le supremum en ce qui concerne le risque de première espèce, mais qu'on ne le fait pas en ce qui concerne le risque de seconde espèce. La raison en est que, dans les modèles qu'on est amené à considérer en pratique, il est possible de rendre un de ces deux risques uniformément petit, mais pas les deux simultanément... L'explication à ce phénomène sera abordée un peu plus loin dans ce polycopié : néanmoins, à ce stade, vous avez déjà suffisamment de nouveautés à intégrer dans ce chapitre pour qu'il soit préférable de simplement accepter l'affirmation ci-dessus telle quelle dans un premier temps! ☺

En tout cas, cela justifie ce traitement dissymétrique des hypothèses complémentaires : puisqu'il est impossible de contrôler *uniformément* les deux risques à la

fois, on décide qu'il y aura une hypothèse pour laquelle le risque sera contrôlé uniformément, et l'autre pour laquelle on se devra se contenter d'une fonction de risque selon la valeur du paramètre caché : et c'est celle qu'on choisit pour le contrôle uniforme qu'on qualifie d'« hypothèse nulle », et l'autre qu'on qualifie d'« hypothèse alternative » ! \clubsuit

12.3 Construction d'un test à partir d'une statistique

Notion de statistique de test

La définition (AG') nous définit ce qu'est, de manière générale, un test. Pour *construire* un test, la procédure de base (dont nous verrons une extension un peu plus loin [Procédure (CQ')]) est la suivante :

Définition (AZ') (Statistique de test unilatérale). Lorsqu'on cherche à tester l'hypothèse nulle \mathcal{H}_0 contre l'hypothèse alternative \mathcal{H}_1 , une *statistique de test à droite* est une statistique $t(X) =: T$ à valeurs réelles dont on « espère » qu'elle a « tendance » à prendre des valeurs plus grandes lorsque \mathcal{H}_1 est vérifiée que lorsque \mathcal{H}_0 est vérifiée. Dans ce cas, on dit aussi qu'on « soupçonne les grandes valeurs de T » (en partant toujours du principe que « être suspect », cela signifie « être de nature à nous inciter à rejeter l'hypothèse nulle »).

À l'inverse, une *statistique de test à gauche* est une statistique $t(X) =: T$ à valeurs réelles dont on « espère » qu'elle a « tendance » à prendre des valeurs plus petites (au sens de « plus proches de $-\infty$ », quand bien même la valeur absolue deviendrait plus grande) lorsque \mathcal{H}_1 est vérifiée que lorsque \mathcal{H}_0 est vérifiée. Dans ce cas, on dit aussi qu'on « soupçonne les petites valeurs de T ». (On peut aussi dire « les valeurs les plus à gauche » pour souligner que « petit » signifie ici « petit en valeur *signée* », pas en valeur absolue).

Collectivement, les statistiques de test à gauche et à droite sont qualifiées de *statistiques de test unilatérales*. Une statistique de test unilatérale étant donnée, on parlera du *critère de suspicion* associé à cette statistique pour dire si on en soupçonne les grandes valeurs ou au contraire les petites valeurs (autrement dit, pour dire s'il s'agit d'une statistique à droite ou à gauche) : un « critère de suspicion » ne peut donc, à ce stade, que prendre deux modalités : « on soupçonne les valeurs les plus à droite » ou « on soupçonne les valeurs les plus à gauche ». (Nous verrons un peu plus tard une troisième modalité, conférer définition (CO')).

Une statistique de test unilatérale étant donnée, on dira qu'un test est associé à cette statistique de test lorsque :

- Dans le cas d'une statistique de test à droite : le test est de la forme $\{T > t_{\text{seuil}}\}$ pour un certain $t_{\text{seuil}} \in [-\infty, +\infty]$;
- Dans le cas d'une statistique de test à gauche : le test est de la forme $\{T < t_{\text{seuil}}\}$ pour un certain $t_{\text{seuil}} \in [-\infty, +\infty]$.

\heartsuit

Remarque (BA'). Bien noter que, de même que la définition d'« être un estimateur » et « être un test », la définition d'« être une statistique de test » est, en tant que telle, « creuse » : quand on dit qu'on « espère » qu'une statistique de test à droite prendra de plus grandes valeurs sous \mathcal{H}_1 que sous \mathcal{H}_0 , non seulement cela ne veut, formellement parlant, rien dire ; mais en plus, quand bien même la statistique prendrait des valeurs comparables sous \mathcal{H}_1 et sous \mathcal{H}_0 , voire aurait tendance à prendre

de plus *petites* valeurs sous \mathcal{H}_1 que sous \mathcal{H}_0 , cela n'empêcherait pas, *techniquement*, de la traiter comme une statistique de test à droite : le cas échéant, cette statistique ne serait *pas intelligente* (au sens où on n'arriverait pas à construire un test *pertinent* à partir d'elle) ; mais du point de vue de la définition formelle il n'y aurait pas de problème à dire quand même que c'est une statistique de test à droite ! Prendre une statistique (réelle) donnée comme statistique de test, puis décider d'en soupçonner les grandes ou les petites valeurs, est donc un *choix* du statisticien : tous les choix sont techniquement possibles, et aucun ne sera formellement faux : en revanche, certains seront plus pertinents que d'autres... ! ♣

Remarque (BB'). Notez que j'ai décidé de toujours utiliser des inégalités *strictes* pour définir les test : lorsque la statistique de test est exactement égale à la valeur-seuil, on penche systématiquement pour l'hypothèse nulle. Techniquement parlant, ce choix n'est pas obligatoire ; néanmoins cette convention s'avèrera très pratique dans la suite (en particulier pour définir un test à partir d'une statistique, d'un critère de suspicion et d'un niveau, cf. procédure (BD') *infra*) : je vous conseille donc de toujours la respecter. (Car la violer n'apporterait quasiment rien et risquerait de rendre certaines des constructions ci-dessous fausses...). ♣

Remarque (BC'). En pratique, on construit quasiment toujours les tests en passant par une statistique de test (qui peut néanmoins être, dans certains cas, une statistique de test *bilatérale*, confer § 12.5). (Et d'ailleurs, même dans les rares cas où aucune statistique de test n'est formellement introduite, on pourrait toujours considérer que le test *Test* est en fait associé à la statistique $\mathbf{1}_{Test}$ avec critère de suspicion à droite... ! ☺). ♣

Construction d'un test de niveau imposé

Dans la définition (AZ') ci-dessus, ainsi que dans l'exemple introductif, nous avons utilisé la statistique de test en fixant *d'abord* un seuil pour la statistique, *puis* en étudiant le niveau associé à ce seuil. En pratique, lorsqu'on cherche à construire un test booléen, on procède généralement dans l'autre sens : on *fixe* d'abord le niveau qu'on désire pour notre test ; puis on *cherche* le seuil permettant d'atteindre ce niveau (mais sans pousser le seuil trop loin pour garder une chance de conclure en faveur de l'hypothèse alternative lorsqu'elle est vraie). La procédure pour trouver le niveau adéquat est toujours la même, et vous devez la connaître :

! **Procédure (BD')**. Une partition de l'espace du paramètre caché en hypothèse nulle et hypothèse alternative étant donnée, un niveau $\alpha \in]0, 1/2]$ étant fixé, et une statistique de test unilatérale (avec son critère de suspicion) étant choisie, on construit canoniquement, à partir de ces éléments, « le » test de niveau α de la façon suivante :

— Dans le cas d'une suspicion à gauche, on fixe le seuil à

$$t_{\text{seuil}} = \inf\{\text{Qtile}(\text{Loi}_\theta(T); \alpha) \mid \theta \in \Theta_0\} \quad (\text{BE}')$$

(puis on considèrera le test $\{T < t_{\text{seuil}}\}$) ;

— Dans le cas d'une suspicion à droite, on fixe le seuil à

$$t_{\text{seuil}} = \sup\{\text{Qtile}(\text{Loi}_\theta(T); 1 - \alpha) \mid \theta \in \Theta_0\} \quad (\text{BF}')$$

(puis on considèrera le test $\{T > t_{\text{seuil}}\}$). ♡

Pour comprendre d'où vient cette procédure (et pour la retenir), voyons comment on démontre que cela donne bien un test de niveau α . Nous ne traiterons ici que le cas des tests à droite, mais celui des tests à gauche serait similaire.

Démonstration de la validité de la procédure. Plaçons-nous donc dans le cas d'un test à droite. Soit t_* un seuil envisagé, et demandons-nous à quelle condition est-ce que le test $Test := \{T > t_*\}$ de l'hypothèse nulle $\{\theta \in \Theta_0\}$ sera de niveau α .

Pour que le test soit de niveau α , par définition du niveau, il faut et il suffit qu'on ait, pour tout $\theta \in \Theta_0$, que $\mathbb{P}_\theta(Test) \leq \alpha$, autrement dit, que $\mathbb{P}_\theta(T > t_*) \leq \alpha$. Cela peut encore se réécrire, par passage au complémentaire, en disant que la fonction de répartition de $Loi_\theta(T)$, évaluée en t_*+ , doit valoir au moins $1 - \alpha$:

$$\text{Répart}(Loi_\theta(T); t_*+) \geq 1 - \alpha. \quad (\text{BG}')$$

Mais dire qu'il y a au moins une fraction $(1 - \alpha)$ de la masse de $Loi_\theta(T)$ à gauche de t_* , c'est équivalent à dire que le seuil auquel on sépare une fraction $(1 - \alpha)$ à gauche d'une fraction α à droite est situé *avant* t_* , autrement dit, que le quantile de niveau $(1 - \alpha)$ de $Loi_\theta(T)$ est plus petit que t_* ! La condition que nous avons écrite sur la fonction de répartition est ainsi équivalente^[††] à :

$$\text{Qtile}(Loi_\theta(T); 1 - \alpha) \leq t_*. \quad (\text{BH}')$$

Ainsi, pour qu'on ait bien un test de niveau α , le seuil t_* doit être (et il suffit qu'il le soit) supérieur ou égal à $\text{Qtile}(Loi_\theta(T); 1 - \alpha)$; et ce, pour tous les $\theta \in \Theta_0$: ce qui revient bien à dire que t_* doit être au moins égal au supremum de ces quantiles lorsque θ parcourt Θ_0 . Parmi l'ensemble de ces choix possible pour fixer le seuil, il est alors naturel de prendre le plus petit. En effet, comme nous le verrons plus loin, la condition sur le niveau contrôle simplement le risque de répondre « positif » alors qu'on est sous l'hypothèse nulle ; mais, autant que possible, on voudrait *aussi* limiter le risque de répondre « négatif » à tort : or, plus t_* est grand, plus le test $\{T > t_*\}$ répondra facilement « négatif » ; et dès lors, pour contrôler le second type de risque, il est logique de chercher à prendre t_* aussi petit que possible parmi les valeurs compatibles avec le niveau requis ! Ce qui nous conduit à la définition du t_{seuil} de la procédure. \heartsuit

Remarque (BI'). La preuve ci-dessus est écrite en termes assez formels ; mais une fois que vous l'avez comprise, son idée doit vraiment être simple à retenir (de sorte que vous pourrez immédiatement écrire la formule pour t_{seuil}) : pour qu'on soit de niveau α , il faut que, sous les différents contextes $\mathbb{P}_\theta(\bullet)$ correspondant à l'hypothèse nulle, la probabilité de dépasser le seuil n'excède pas α ; autrement dit, que le seuil soit au-delà du quantile de niveau $(1 - \alpha)$ de $Loi_\theta(T)$: et puisqu'on veut que cela soit le cas pour tous les θ de Θ_0 , on va prendre le supremum de ces quantiles lorsque θ décrit l'hypothèse nulle ! \heartsuit

Remarque (BJ'). Lorsqu'on détermine le seuil comme un supremum ou infimum de quantiles, il est fréquent qu'on doive procéder à une approximation à un moment donné, ne serait-ce que pour arrondir la valeur qui est, en toute généralité,

[††]. En toute rigueur, pour que l'équivalence soit parfaite, il faudrait préciser qu'on prend dans (??) la version *continue à droite* de la fonction de quantile, c.à.d. que le quantile de niveau " $1 - \alpha$ " est, en toute rigueur, le quantile de niveau " $(1 - \alpha)-$ "... Néanmoins, conformément à notre habitude, nous négligerons ces subtilités en ce qui concerne les fonctions de quantile \heartsuit

un nombre réel non décimal. (Il peut aussi arriver qu'on ne sache pas déterminer *exactement* les quantiles, mais seulement en donner un encadrement). Dans ce cas, il convient de toujours faire l'approximation dans le sens qui assure que le quantile reste bien de niveau α : ainsi, dans le cas d'un test à droite, on prend une approximation du seuil dont on sache qu'elle est située un peu *au-dessus* du seuil critique (et pas en-dessous, auquel cas on augmenterait la probabilité que le test soit positif, ce qui risquerait d'enfreindre la condition de niveau α) ; tandis que dans le cas d'un test à gauche, on doit prendre une approximation située un peu *en-dessous* de la vraie valeur. Typiquement, lorsqu'il s'agit de faire des arrondis, on fera l'arrondi du seuil par excès pour un test à gauche, resp. par défaut pour un test à droite. Pour souligner le fait que les approximations opérées garantissent qu'on conserve l'exigence d'« être de niveau α », on parle d'approximations *conservatives*. Cette idée de conservativité est essentielle dans la science de la maîtrise des risques (industriels, sanitaires, financiers), où nos modèles sont entachés d'une certaine incertitude, mais où on veut néanmoins avoir une garantie solide sur le niveau de risque résiduel !

Dans un sens, c'est aussi l'idée de conservativité qui explique pourquoi, quand la statistique de test est pile au niveau du seuil, on décide que le test sera négatif : pour respecter le contrat « être de niveau α », on peut en effet se permettre de répondre un peu moins souvent « positif » que le maximum requis, mais en aucun cas de répondre « positif » *trop* souvent ! Donc, dans le cas-limite, mieux vaut répondre « négatif » ☺^[††] ♣

Exemples de constructions de tests

Nous allons conclure cette section en montrant des exemples de constructions de tests. Mais avant de nous lancer dans ces exemples, il convient de souligner un point qui intervient fréquemment dans ce genre de constructions :

! **Remarque (BK ')**. Le supremum qui intervient dans la définition du niveau d'un test peut fréquemment être restreint à un certain sous-ensemble de Θ_0 : moralement, l'idée est que, dans le cas (mettons) d'un test à droite, pour étudier le supremum des valeurs $\mathbb{P}_\theta(T > t_{\text{seuil}})$, il y a certaines valeurs de θ qu'on peut se permettre d'omettre, parce qu'il y a d'autres valeurs θ' , situées également dans Θ_0 , pour lesquelles T prend des valeurs « encore plus grandes » sous $\mathbb{P}_{\theta'}$ que sous \mathbb{P}_θ , de sorte que ces θ -là ne contribueront pas au supremum. En général, pour étudier le supremum, on peut en fait se limiter à un sous-ensemble de l'hypothèse nulle correspondant plus ou moins au « bord » de Θ_0 ♣

Remarque (BL '). En continuation de la remarque précédente, il est également fréquent que, au moins sur le sous-ensemble de Θ_0 considéré, la loi la statistique de test sous \mathbb{P}_θ soit en fait la même *pour tous les θ à considérer*^[*] : dès lors, le

[††]. Pour être plus précis, il faut distinguer deux cas : lorsque les lois $\text{Loi}_\theta(T)$ sont diffuses, en fait, la probabilité que la statistique de test vaille *exactement* t_{seuil} est nulle (sous la vraie loi), de sorte que définir le test par une inégalité large ou stricte ne change absolument rien en pratique. En revanche, si les lois $\text{Loi}_\theta(T)$ sont discrètes, il devient très important de bien définir le test par une inégalité stricte ! En effet dans ce cas, il se passe (génériquement) la chose suivante : le test défini, conformément aux explications ci-dessus, par $\{T > t_{\text{seuil}}\}$ aura un niveau strictement inférieur à α (et restera donc « de niveau α ») ; par contre, le test défini par $\{T \geq t_{\text{seuil}}\}$ (avec inégalité *large*) aura un niveau *strictement supérieur* à α , et ne constituera donc pas un test de niveau α ... !

[*]. Attention : Je ne suis pas en train de dire que la statistique T aurait la même loi sous absolument tous les $\mathbb{P}_\theta(\bullet)$: une telle statistique de test ne présenterait en effet aucun intérêt, puisqu'elle aurait le même comportement sous l'hypothèse nulle et sous l'hypothèse alternative !

calcul du supremum peut en fait se faire en considérant *une seule* loi... Ce qui est évidemment bien plus simple! (En fait, un certain nombre de test classiques sont justement construits de façon à faire en sorte que la statistique de test ait toujours la même loi pour θ_{\vee} dans Θ_0 : confer principe (CN') *infra*).

En outre, cette idée que le supremum définissant le niveau peut être restreint à un sous-ensemble de valeurs plus simples, voire être ramené à un calcul sur une seule loi, marchera exactement de la même façon lorsqu'il s'agira de déterminer le seuil définissant le test à partir d'un certain niveau. (Confer les exemples ci-dessous). ♣

Remarque (BM'). En fait, la remarque précédente avait déjà été plus ou moins utilisée dans l'exemple de la § 12.1, lorsqu'on avait convenu que, pour évaluer le risque que le test tranche en faveur des collègues alors que c'est l'enseignant qui avait raison, il fallait considérer un cas où l'enseignant avait raison "de justesse" (et, en l'occurrence, vu que l'hypothèse nulle revenait à affirmer que $\{\mu \geq \mu_{\text{réf}}\}$, nous avons choisi de nous placer dans la situation où $\mu_{\vee} = \mu_{\text{réf}} + 0,1$). En fait, ce qu'on aurait pu montrer en étant plus rigoureux, c'est que, quel que soit le seuil t_{seuil} choisi pour définir le test, le supremum des valeurs $\mathbb{P}_{\mu, \sigma_{\text{réf}}}(M < t_{\text{seuil}})$ pour $\mu \geq \mu_{\text{réf}}$ serait atteint précisément pour $\mu = \mu_{\text{réf}}$ (ce qui correspond bien au cas situé "au bord" de l'hypothèse nulle) : le véritable niveau du test se calculant donc, en l'occurrence, en évaluant $\mathbb{P}_{\mu_{\text{réf}}, \sigma_{\text{réf}}}(M < 95)$. (Et cette probabilité est très proche, en pratique, de la probabilité $\mathbb{P}_{\mu_{\text{réf}}+0,1, \sigma_{\text{réf}}}(M < 95)$ que nous avons calculée dans la § 12.1 : elles valent respectivement 5,9 % et 5,6 %). ♣

Exemple (BN'). Reprenons, pour commencer, l'exemple de notre sous-section introductive : on considère donc le modèle du pédagogue où l'espace du paramètre caché est restreint à $\Theta^{\text{hsc}} := \mathbb{R} \times \{\sigma_{\text{réf}}\}$, et on cherche à tester l'hypothèse nulle $\{\mu \geq \mu_{\text{réf}}\} =: \mathcal{H}_0$ contre l'hypothèse alternative $\{\mu < \mu_{\text{réf}}\} =: \mathcal{H}_1$. On choisit d'utiliser comme statistique de test la moyenne de la première promotion testée, autrement dit, $M := (\sum_{i=0}^{n-1} X_i) / n$; avec suspicion pour les petites valeurs. La seule différence que nous allons prendre par rapport à la § 12.1 est sur la façon de choisir le seuil du test : contrairement à ce que nous avons fait précédemment, nous allons ici imposer qu'on souhaite avoir un test de niveau $\alpha := 5 \%$, et regarder comment on en déduit le seuil approprié : cela nous permettra ainsi de mettre en pratique la procédure (BD') et la remarque (BK') ci-dessus.

En vertu de la procédure (BD'), notre test devra être de la forme $\{M < m_{\text{seuil}}\}$, avec

$$m_{\text{seuil}} = \inf_{\mu \geq \mu_{\text{réf}}} \text{Qtile}(\text{Loi}_{\mu, \sigma_{\text{réf}}}(M); \alpha)^{[\dagger]} = \inf_{\mu \geq \mu_{\text{réf}}} \text{Qtile}(\text{Normale}(\mu, \sigma_{\text{réf}}^2 / n); \alpha). \quad (\text{BO}')$$

Nous allons maintenant voir comment le calcul de cet infimum peut être grandement simplifié. D'après les propriétés des lois normales, $\text{Normale}(\mu, \sigma_{\text{réf}}^2 / n) = \mu + \text{Normale}(0, \sigma_{\text{réf}}^2 / n)$. Or, puisque l'application $x \mapsto \mu + x$ est croissante, pour

Ce que je dis, c'est seulement que notre statistique aura fréquemment la même loi sous tous les $\mathbb{P}_{\theta(\bullet)}$ pour θ dans Θ_0 — ou du moins, pour θ dans l'ensemble "critique", marquant le "bord" de Θ_0 , auquel nous avons fait allusion à l'instant.

[†]. Notez ici que, même s'il *semble* qu'on ait pris notre infimum sur μ plutôt que sur θ — alors que j'ai pourtant bien expliqué [remarque (AP') *supra*] qu'il fallait toujours prendre l'infimum sur θ *lui-même*! —, la forme spécifique de l'espace du paramètre caché *dans cet exemple* fait que μ est en bijection avec $\theta = (\mu, \sigma_{\text{réf}})$

toute loi P , pour tout $p \in]0, 1[$, on a $Q_{\text{tile}}(\mu + P; p) = \mu + Q_{\text{tile}}(P; p)$: en l'occurrence, on a donc

$$Q_{\text{tile}}(\text{Normale}(\mu, \sigma_{\text{réf}}^2 / n); \alpha) = Q_{\text{tile}}(\text{Normale}(0, \sigma_{\text{réf}}^2 / n); \alpha) + \mu, \quad (\text{BP}')$$

d'où

$$\begin{aligned} \inf_{\mu \geq \mu_{\text{réf}}} Q_{\text{tile}}(\text{Normale}(\mu, \sigma_{\text{réf}}^2 / n); \alpha) &= Q_{\text{tile}}(\text{Normale}(0, \sigma_{\text{réf}}^2 / n); \alpha) + \mu_{\text{réf}} \\ &= Q_{\text{tile}}(\text{Normale}(\mu_{\text{réf}}, \sigma_{\text{réf}}^2 / n); \alpha) : \quad (\text{BQ}') \end{aligned}$$

on voit ainsi qu'en fait, l'infimum des quantiles pour toutes les valeurs de θ dans Θ_0 correspond au quantile pour θ égal à la valeur "critique" $(\mu_{\text{réf}}, \sigma_{\text{réf}})$ située au bord de Θ_0 .

Il ne reste plus qu'à utiliser un logiciel numérique pour calculer la fonction de table que nous avons fait apparaître :

```
> qnorm(5 / 100, muref, sqrt(sigma ^ 2 / n))
[1] 94.73974
```

cela nous suggère, en fin de compte, de prendre pour test $\{M < 94,7\}$, en prenant garde d'arrondir la valeur numérique trouvée *par dessous*^[‡] pour être sûr que, même après arrondi, le test reste bien de niveau α . \clubsuit

Exemple (BR'). Pour bien illustrer la remarque (BK'), reprenons l'exemple précédent, et regardons quel est l'impact de l'arrondi final sur le véritable niveau du test (dans la suite, nous utiliserons la notation $94,7 =: \tilde{\mu}_{\text{seuil}}$, et nous noterons $\tilde{\alpha}$ le niveau du test correspondant). Par définition, on a que

$$\tilde{\alpha} \stackrel{\text{déf}}{=} \sup_{\mu \geq \mu_{\text{réf}}} \mathbb{P}_{\mu, \sigma_{\text{réf}}} (M < \tilde{\mu}_{\text{seuil}}) = \sup_{\mu \geq \mu_{\text{réf}}} \mathbb{P}(\text{Normale}(\mu, \sigma_{\text{réf}}) < \tilde{\mu}_{\text{seuil}}). \quad (\text{BS}')$$

La remarque (BK') nous suggère alors que le supremum ci-dessus sera atteint pour $\mu = \mu_{\text{réf}}$: et de fait, ce sera bien le cas. Mais comment le démontrer ? Ici, plutôt que de chercher à expliciter la façon dont $\mathbb{P}(\text{Normale}(\mu, \sigma_{\text{réf}}) < \tilde{\mu}_{\text{seuil}})$ dépend de μ , nous allons juste *comparer* les valeurs entre $\mu_{\text{réf}}$ et les autres valeurs de μ . Soit donc $\mu \geq \mu_{\text{réf}}$; notre but va être de montrer que

$$\mathbb{P}(\text{Normale}(\mu, \sigma_{\text{réf}}) < \tilde{\mu}_{\text{seuil}}) \leq \mathbb{P}(\text{Normale}(\mu_{\text{réf}}, \sigma_{\text{réf}}) < \tilde{\mu}_{\text{seuil}}). \quad (\text{BT}')$$

Pour ce faire, on observe que, d'après les propriétés des lois normales, on a

$$\text{Normale}(\mu, \sigma_{\text{réf}}) = \text{Normale}(\mu, \sigma_{\text{réf}}) + \mu - \mu_{\text{réf}} : \quad (\text{BU}')$$

en particulier, si M est une v.a. de loi Normale($\mu_{\text{réf}}, \sigma_{\text{réf}}^2 / n$), alors la v.a. $M' := M + \mu - \mu_{\text{réf}}$ sera de loi Normale($\mu, \sigma_{\text{réf}}^2 / n$) ; et en outre, puisque $\mu \geq \mu_{\text{réf}}$, on aura presque-surement que $M' \geq M$ ^[§]. Mais, au vu de cette comparaison, l'évènement $\{M' < \tilde{\mu}_{\text{seuil}}\}$ implique l'évènement $\{M < \tilde{\mu}_{\text{seuil}}\}$, de sorte que

$$\mathbb{P}(M' < \tilde{\mu}_{\text{seuil}}) \geq \mathbb{P}(M < \tilde{\mu}_{\text{seuil}}) : \quad (\text{BV}')$$

[‡]. En l'occurrence, l'arrondi par-dessous correspond à l'arrondi au plus proche ; mais si nous avions plutôt voulu arrondir point près, ou au centième de près, cette fois-ci notre règle aurait conduit à prendre des arrondis à resp. 94 et 94,73, qui n'auraient pas coïncidé avec l'arrondi au plus proche !

[§]. Cette idée de trouver des variables aléatoires instanciant les lois Normale($\mu_{\text{réf}}, \sigma_{\text{réf}}^2 / n$) et Normale($\mu, \sigma_{\text{réf}}^2 / n$) sur un même espace probabilisé, d'une façon qu'on ait une comparaison pertinente entre les deux v.a. qui soit de type *presque-sure* (ou, plus généralement, qui soit une comparaison en probabilité), s'appelle « réaliser un *couplage* entre les deux lois ».

ce qui, quand on se rappelle quelles lois suivent resp. M' et M , signifie bien que, conformément à ce qu'on souhaitait prouver, on a que $\mathbb{P}(\text{Normale}(\mu, \sigma_{\text{réf}}^2 / n) < \tilde{\mu}_{\text{seuil}}) \geq \mathbb{P}(\text{Normale}(\mu_{\text{réf}}, \sigma_{\text{réf}}^2 / n) < \tilde{\mu}_{\text{seuil}})!$

Ainsi, on a que

$$\tilde{\alpha} = \text{répartNormale}(\mu_{\text{réf}}, \sigma_{\text{réf}}^2 / n; \tilde{\mu}_{\text{seuil}}-), \quad (\text{BW}')$$

ce qu'il ne reste plus qu'à calculer numériquement [¶] :

```
> pnorm(94.7, muref, sqrt(sigma ^ 2 / n))
[1] 0.04873142
```

On trouve ainsi que le niveau véritable de notre test est de 4,88 % [¶] (en arrondissant ici *par-dessus*, afin qu'il soit bien légitime de dire que « c'est un test de niveau 4,88 % »). ☺

12.4 Recherche de statistiques de test pertinentes

Pour appliquer la théorie des tests de façon utile, il va falloir construire des tests qui auront à la fois un bon niveau (typiquement, $\alpha \leq 10\%$) et une fonction de puissance aussi élevée que possible. Comment construire de tels tests? Comme expliqué ci-dessus, en pratique on construit les tests à partir de statistiques de test (associées à un critère de suspicion) : la question devient donc de savoir comment on peut trouver une bonne statistique de test!

En tant qu'ingénieurs généralistes, en fait, il devrait être exceptionnel que vous soyez amené(e) à créer vous-mêmes des statistiques de test : le point essentiel sera surtout de s'assurer qu'une statistique de test donnée est pertinente! Et sur ce point, en général, on procédera plus souvent par des raisonnements informels que par des preuves rigoureuses, l'enjeu étant juste de s'assurer que le choix de la statistique de test est "raisonnable"...

Le principe fondamental pour vérifier qu'une statistique de test est pertinente est toujours le même. Ci-dessous nous allons légèrement anticiper sur la section suivante du chapitre, en parlant d'un critère de suspicion que vous n'avez pas encore rencontré : le « critère de suspicion bilatéral ». Ce qu'on entend par là, et la façon de construire les tests correspondants, sera développé un peu plus tard.

Principe (BX'). *Considérons un modèle statistique, où l'espace du paramètre caché est partitionné en deux hypothèses complémentaires, avec les notations standard. Une statistique de test $t(X)$ est appropriée dans les cas suivants :*

- (i) *Lorsque $\text{Loi}_\theta(t(X))$ a "tendance à prendre des valeurs plus grandes^[**] lorsque $\theta \in \Theta_1$ que lorsque $\theta \in \Theta_0$ ^[††]" : dans ce cas, le critère de suspicion approprié est un critère de suspicion à droite.*
- (ii) *Symétriquement, lorsque $\text{Loi}_\theta(t(X))$ a tendance à prendre des valeurs plus petites lorsque $\theta \in \Theta_1$ que lorsque $\theta \in \Theta_0$: dans ce cas, il convient de soupçonner les valeurs les plus à gauche.*

[¶]. Notez que, comme les lois normales sont diffuses, il n'est pas important de savoir si on considère la fonction de répartition en $\tilde{\mu}_{\text{seuil}}-$ ou en $\tilde{\mu}_{\text{seuil}}+$.

[¶]. Cette valeur est très proche de 5 %, ce qui montre que notre arrondi n'était effectivement pas bien méchant! ☺

!

- (iii) Lorsque $\text{Loi}_\theta(t(X))$ n'a jamais tendance à prendre des valeurs trop grandes aussi longtemps que $\theta \in \Theta_0$, qu'il existe des valeurs $\theta \in \Theta_1$ pour lesquelles $\text{Loi}_\theta(t(X))$ prend des valeurs (nettement) plus grandes que dans n'importe quel cas associé à Θ_0 , et qu'il n'existe aucune valeur $\theta \in \Theta_1$ pour lesquelles $\text{Loi}_\theta(t(X))$ prenne des valeurs (nettement) plus petites que tous les cas associés à Θ_0 : alors le critère de suspicion approprié est une suspicion à droite.
- (iv) Symétriquement, lorsque $\text{Loi}_\theta(t(X))$ n'a jamais tendance à prendre des valeurs trop grandes aussi longtemps que $\theta \in \Theta_0$, qu'il existe des valeurs $\theta \in \Theta_1$ pour lesquelles $\text{Loi}_\theta(t(X))$ prend des valeurs (nettement) plus petites que dans n'importe quel cas associé à Θ_0 , et qu'il n'existe aucune valeur $\theta \in \Theta_1$ pour lesquelles $\text{Loi}_\theta(t(X))$ prenne des valeurs (nettement) plus grandes que tous les cas associés à Θ_0 : alors le critère de suspicion approprié est une suspicion à gauche.
- (v) Lorsque, aussi longtemps que θ est dans Θ_0 , $\text{Loi}_\theta(t(X))$ n'a jamais tendance à prendre des valeurs ni trop petites, ni trop grandes ; qu'il existe en revanche des valeurs θ dans Θ_1 pour lesquelles $\text{Loi}_\theta(t(X))$ prend des valeurs nettement plus grandes que dans les cas associés à Θ_0 , et qu'il existe aussi des valeurs θ dans Θ_1 pour lesquelles $\text{Loi}_\theta(t(X))$ prend des valeurs nettement plus petites que dans les cas associés à Θ_0 : alors le critère de suspicion approprié est une suspicion bilatérale. \diamond

Remarque (BY'). En fait, à partir d'une statistique de test et d'un critère de suspicion, on peut *toujours* définir un test de niveau α : prendre le "mauvais" critère de suspicion ne signifie donc pas qu'on n'arrivera pas à construire de test, mais que ce test sera dénué d'intérêt, parce que sa puissance sera beaucoup trop faible ! (typiquement, quand le critère de suspicion est mal choisi, ce qui risque de se passer est que la puissance du test se retrouve uniformément majorée par α (qui, rappelons-le, est typiquement une valeur très petite) : ainsi, lorsque l'hypothèse alternative est vraie, on ne conclut alors quasiment jamais en sa faveur, et ce, peu importe la valeur θ_ν dans \mathcal{H}_1 ! \clubsuit

! *Remarque (BZ').* On notera que, lorsqu'on vérifie qu'une statistique de test est pertinente (et qu'on détermine du critère de suspicion approprié), c'est à cet endroit, et *seulement* à cet endroit, qu'on est amené à regarder *aussi* quelle est l'hypothèse alternative \mathcal{H}_1 ! (Alors que, pour déterminer le seuil correspondant à un niveau donné, ou pour déterminer le niveau correspondant à un seuil donné, seule l'hypothèse nulle intervient...).

C'est tout à fait logique, puisque l'idée qu'une statistique de test soit pertinente va se traduire par la *puissance* qu'elle permet d'obtenir pour un niveau de risque donné. Or pour avoir une bonne puissance, il faut que les valeurs prises de la statistique de test sous l'hypothèse alternative soient (en général) suffisamment

[††]. Dans le cadre de ce cours, nous ne chercherons pas à définir précisément ce qu'on entend par là. De manière générale, l'idée serait de dire qu'il y a *domination stochastique* (au moins de manière approximative) des lois $\text{Loi}_\theta(t(X))$ pour $\theta \in \Theta_0$ par les lois $\text{Loi}_\theta(t(X))$ pour $\theta \in \Theta_1$.

[‡‡]. Ici, il faut comprendre qu'on exige que, pour *n'importe quelles* v valeurs $\theta_0 \in \Theta_0$ et $\theta_1 \in \Theta_1$, $\text{Loi}_{\theta_1}(t(X))$ a tendance à prendre des valeurs plus grandes (au moins approximativement) que $\text{Loi}_{\theta_0}(t(X))$. De manière générale, juste dire, par exemple « pour tout $\theta_0 \in \Theta_0$, on peut trouver un $\theta_1 \in \Theta_1$ tel que $\text{Loi}_{\theta_1}(t(X))$ ait tendance à prendre des valeurs plus grandes que $\text{Loi}_{\theta_0}(t(X))$; et pour tout $\theta_1 \in \Theta_1$, on peut trouver un $\theta_0 \in \Theta_0$ tel que $\text{Loi}_{\theta_0}(t(X))$ ait tendance à prendre des valeurs plus petites que $\text{Loi}_{\theta_1}(t(X))$ » ne serait pas suffisant ! Néanmoins, s'il s'agit juste de trouver le critère de suspicion approprié pour une statistique de test donnée, on pourra se contenter d'un tel argument.

différentes des valeurs pouvant être typiquement prises sous l'hypothèse nulle pour qu'on puisse se rendre compte que l'hypothèse nulle n'est pas plausible : on est donc amené à *comparer* le comportement de la statistique de test entre les hypothèses nulle et alternative pour comprendre (au moins de façon grossière) si la puissance va bien se comporter. ♣

Nous allons maintenant voir un exemple où on peut vérifier la pertinence d'une statistique de test à l'aide de nos critères qualitatifs. En fait, nous allons reprendre l'exemple trouvé ci-dessus, et la statistique de test de Neyman-Pearson que nous avons calculée : il s'agira juste de vérifier que cela conduit bien à quelque chose satisfaisant les critères qualitatifs du principe présenté ci-dessus ! ☺

Exemple (CA'). Nous sommes donc dans le modèle du pédagogue, en testant l'« hypothèse du vieux grincheux^[*] » $\{\mu = \mu_{\text{réf}} \text{ et } \sigma = \sigma_{\text{réf}}\}$ contre l'hypothèse alternative complémentaire (dans Θ^{orig}). Et nous voulons utiliser pour ce faire la statistique de test^[†]

$$T := (M - \mu_{\text{réf}})^2 + S^2 - 2\sigma_{\text{réf}}^2 \ln S, \quad (\text{CB}')$$

dont nous allons feindre de ne pas connaître le critère de suspicion approprié.

Commençons par observer qu'ici, puisque l'hypothèse nulle est simple, il y a une seule loi $\text{Loi}_\theta(T)$ possible lorsque $\theta \in \Theta_0$: dès lors, on peut dire automatiquement que cette loi ne pourra jamais avoir tendance à prendre des valeurs arbitrairement grandes, ni arbitrairement petites, dans un tel cas.

Maintenant, que se passe-t-il sous l'hypothèse alternative ? Il y a quatre façons (non exclusives) dont le vieux grincheux peut avoir tort, selon ce que vaut le vrai paramètre $\theta_{\checkmark} = (\mu_{\checkmark}, \sigma_{\checkmark})$:

- Soit parce que $\mu_{\checkmark} < \mu_{\text{réf}}$, c.-à-d. que la nouvelle méthode fait en réalité globalement régresser le niveau des élèves : dans ce cas, sous réserve que l'écart entre $\mu_{\text{réf}}$ et μ_{\checkmark} soit suffisamment grand, on s'attend à ce que les notes soient très mauvaises, autrement dit que M soit (sauf coup de chance) nettement plus petit que $\mu_{\text{réf}}$, ce qui va rendre la statistique de test très élevée à cause de son premier terme.
- Soit, à l'inverse, parce que $\mu_{\checkmark} > \mu_{\text{réf}}$, c.-à-d. que la nouvelle méthode fait globalement progresser les élèves : auquel cas M sera nettement plus grand que $\mu_{\text{réf}}$, ce qui fera à nouveau exploser le premier terme de la statistique de test.
- Soit parce que $\sigma_{\checkmark} < \sigma_{\text{réf}}$, c.-à-d. que la nouvelle méthode donne des résultats plus homogènes que l'ancienne : dans ce cas, on s'attend à avoir un S très proche de 0, auquel cas la statistique de test explosera à cause de son troisième terme.
- Soit parce que $\sigma_{\checkmark} > \sigma_{\text{réf}}$, c.-à-d. que la nouvelle méthode crée plus d'hétérogénéité entre les élèves que l'ancienne : dans ce cas, sous réserve que le ratio entre $\sigma_{\text{réf}}$ et σ_{\checkmark} soit suffisamment grand, on s'attend à ce que les notes soient très dispersées, autrement dit que S soit grand, ce qui fera exploser la statistique de test à cause de son second terme (certes, dans ce cas le troisième

[*]. Concernant cette appellation, confer exemple (CF') un peu plus loin.

[†]. Comme nous le verrons dans l'exemple (CF'), cette statistique de test est en fait celle du rapport des vraisemblances, ce qui implique en particulier que son critère de suspicion sera à droite : mais nous allons feindre d'ignorer ce dernier point... Au demeurant, le raisonnement que nous allons présenter ci-dessous n'a rien de spécifique au rapport des vraisemblances, et pourrait fonctionner de même avec plein d'autres statistiques de test ! ☺

terme sera négatif, mais la décroissance logarithmique ne compense pas la croissance quadratique! ^[‡]).

On voit donc que, quelle que soit la façon dont le vieux grincheux a tort, la statistique de test tendra à prendre des valeurs plus grandes que sous l'hypothèse nulle : ainsi la statistique de test est effectivement pertinente, et le critère de suspicion approprié consiste effectivement à suspecter ses grandes valeurs! \heartsuit

Rapport de vraisemblances

Au-delà des principes généraux énoncé dans les sous-sections précédente et suivante, il existe une procédure systématique pour construire une statistique de test pertinente : c'est de considérer le *rapport des vraisemblances* entre l'hypothèse alternative et l'hypothèse nulle (avec suspicion du côté où l'hypothèse alternative prend de plus grandes valeurs relativement à l'hypothèse nulle). Je n'ai pas marqué ce principe comme à retenir, dans le mesure où on l'utilise assez rarement en pratique ; néanmoins, c'est une "ligne directrice" qu'il est utile de conserver à l'esprit :

Principe (CC '). *De manière générale, une statistique de test pertinente pour tester une hypothèse alternative \mathcal{H}_1 contre une hypothèse nulle \mathcal{H}_0 consiste en la v.a. dont la réalisation est le rapport de vraisemblances ^[§]*

$$\frac{\mathcal{L}(\mathcal{H}_1 \mid X = x_{\mathcal{J}})}{\mathcal{L}(\mathcal{H}_0 \mid X = x_{\mathcal{J}})}, \quad (\text{CD '})$$

le critère de suspicion consistant alors à trouver suspectes les grandes valeurs de ce rapport. \diamond

Remarque (CE '). Il existe un théorème, appelé « théorème de Neyman-Pearson » (voir § ??), qui dit que cette statistique de test est *rigoureusement* optimale (dans un sens bien précis) dans le cas où Θ_0 et Θ_1 sont toutes les deux des singletons. Lorsque ce n'est pas le cas, on n'a malheureusement plus de tel résultat d'optimalité ; néanmoins la stratégie de Neyman-Pearson reste pertinente malgré tout. \heartsuit

Exemple (CF '). Voyons un exemple pour lequel le statistique du rapport des vraisemblances peut être calculée. On se place toujours dans le modèle du pédagogue, avec pour espace du paramètre caché la version originelle $\Theta^{\text{orig}} = \mathbb{R} \times \mathbb{R}_+^*$; mais cette fois-ci, plutôt que de chercher à savoir si la nouvelle méthode va faire progresser les élèves comme prévu ou non, on s'intéresse à l'affirmation d'un vieux professeur grincheux qui affirme que « de toutes façons, le choix de la méthode pédagogique n'a aucune influence sur les résultats des élèves ! ». Ainsi, le vieux professeur défend l'hypothèse que $\{\theta = \theta_{\text{réf}}\}$, avec $\theta_{\text{réf}} := (\mu_{\text{réf}}, \sigma_{\text{réf}})$: nous prendrons cette hypothèse comme l'hypothèse nulle. Notre hypothèse alternative, ici, sera que « le professeur grincheux a tort », à savoir, que la nouvelle méthode pédagogique donne des résultats différents de l'ancienne, que ce soit parce que la moyenne change et/ou parce que la variance change : autrement dit, on prend $\Theta = \mathbb{R} \times \mathbb{R}_+^*$ tout

[‡]. Plus précisément, on vérifie sans difficulté que la fonction $s \mapsto s^2 - 2\sigma_{\text{réf}}^2 \ln s$ a les variations suivantes sur \mathbb{R}_+^* : tendant vers l'infini près de 0, elle décroît jusqu'à l'abscisse $s = \sigma_{\text{réf}}$, avant de croître à nouveau jusque vers l'infini lorsque $s \rightarrow \infty$.

[§]. Rappelons que la notion de vraisemblance d'une hypothèse composite a été définie dans la définition (UJ).

entier et $\Theta_1 := \Theta \setminus \{\theta_{\text{réf}}\}$. Calculons le rapport des vraisemblances dans ce cas. Nous avons vu (cf. § 8.4) que

$$\ln \mathcal{L}((\mu, \sigma) = (\mu, \sigma) \mid X = x_{\mathcal{J}}) = -n \left(\ln \sigma - \frac{1}{2} \sigma^{-2} ((\mu - m_{\mathcal{J}})^2 + s_{\mathcal{J}}^2) \right). \quad (\text{CG}')$$

Pour calculer la vraisemblance de l'hypothèse alternative, il faut passer au supremum (en observant que le supremum sur Θ_1 est égal au supremum sur Θ tout entier par continuité) : vu que, à σ fixé, le supremum est atteint pour $\mu = m_{\mathcal{J}}$, il ne reste plus qu'à optimiser en σ sous la contrainte $\{\mu = m_{\mathcal{J}}\}$:

$$\begin{aligned} \ln \mathcal{L}(\theta \in \Theta_1 \mid X = x_{\mathcal{J}}) &= \sup_{\sigma \in \mathbb{R}_+^*} \left(-n \ln \sigma - \frac{1}{2} n s_{\mathcal{J}}^2 \sigma^{-2} \right) \\ &= -n \ln s_{\mathcal{J}} - \frac{1}{2} n s_{\mathcal{J}}^2 s_{\mathcal{J}}^{-2} = -n \left(\ln s_{\mathcal{J}} + \frac{1}{2} \right), \end{aligned}$$

où l'avant-dernière égalité vient de ce que le supremum est atteint pour $\sigma = s_{\mathcal{J}}$, comme nous l'avons vu lorsque nous avons calculé l'estimateur du maximum de vraisemblance. On obtient donc le rapport de vraisemblances suivant :

$$\ln \frac{\mathcal{L}(\theta \neq \theta_{\text{réf}} \mid X = x_{\mathcal{J}})}{\mathcal{L}(\theta = \theta_{\text{réf}} \mid X = x_{\mathcal{J}})} = -n \left(\ln(s_{\mathcal{J}}/\sigma_{\text{réf}}) + \frac{1}{2} \sigma_{\text{réf}}^{-2} (m_{\mathcal{J}} - \mu_{\text{réf}})^2 + \frac{1}{2} (\sigma_{\text{réf}}^{-2} s_{\mathcal{J}}^2 - 1) \right), \quad (\text{CH}')$$

qui peut être mis en bijection croissante avec la statistique de test suivante (dont on suspectera donc les grandes valeurs), plus simple :

$$(M - \mu_{\text{réf}})^2 + S^2 - 2\sigma_{\text{réf}}^2 \ln S. \quad (\text{CI}')$$

♣

Remarque (CJ'). Comme les calculs ci-dessus le montrent, même dans le cas particulièrement simple du modèle du pédagogue avec hypothèse nulle simple, il peut être assez laborieux de calculer le rapport des vraisemblances : et calculer le seuil approprié pour le test correspondant peut s'avérer encore plus compliqué ! C'est pourquoi, malgré son caractère extrêmement général, la statistique de Neyman-Pearson n'est utilisée que dans une minorité de cas : le plus souvent, il s'avèrera bien plus commode de bricoler une statistique de test ad hoc, facile à calculer et pour laquelle on sache déterminer efficacement des seuils, en s'appuyant simplement sur les principes qualitatifs présentés plus haut pour vérifier la pertinence de cette statistique de test. Et dans les cas où on souhaitera utiliser la stratégie de Neyman-Pearson, il faudra en général procéder à un certain nombre d'approximations pour se ramener à des choses calculables... ♣

Statistique à loi uniforme

Un autre principe utile à connaître pour trouver des statistiques de test pertinentes est le suivant, qui s'appuie directement sur le principe (BX') :

Principe (CK'). Si on arrive à trouver une statistique $t(X)$ ayant la même loi ^[¶] sous tous les $\mathbb{P}_{\theta(\bullet)}$ pour $\theta \in \Theta_0$, mais une loi (en général) différente pour $\theta \in \Theta_1$, alors c'est très probablement une statistique de test pertinente. En fonction de la façon dont les lois pour $\theta \in \Theta_1$ sont susceptibles de différer de la loi (commune) pour $\theta \in \Theta_0$, le critère de suspicion approprié pourra être à gauche, à droite, ou bilatéral. ◇

!

Remarque (CL'). En fait, il n'y a pas besoin que la statistique $t(X)$ ait *exactement* la même loi sous tous les $\mathbb{P}_\theta(\bullet)$ de $\theta \in \Theta_0$ pour que ce principe reste valable : si la loi de $t(X)$ sous tous ces $\mathbb{P}_\theta(\bullet)$ est toujours *quasiment* la même, au sens où on est capable de déterminer un contrôle uniforme sur la ou les queue(s) appropriées des $\text{Loi}_\theta(t(X))$ (par exemple, si toutes ces lois ont la même moyenne et la même variance, ce qui permet alors de les contrôler via l'inégalité de Bienaymé-Tchebychev), alors cela marche aussi! \smile Cependant, le cas qu'on rencontre le plus souvent est où les $\text{Loi}_\theta(t(X))$ sont rigoureusement identiques : cas qui présente en outre l'avantage d'être plus facile à comprendre et à traiter! \smile \clubsuit

Remarque (CM'). Le principe (CK') sera très utilisé dans le chapitre ?? pour construire des *intervalles de confiance*. \clubsuit

Il y a également une variante du principe ci-dessus où l'idée de « avoir toujours la même loi » ne vaut que sur le “bord” de l'hypothèse nulle :

! **Principe (CN')**. Lorsque l'hypothèse nulle Θ_0 est définie par une inégalité (disons, $\gamma(\theta) \leq \gamma_0$), si on arrive à trouver une statistique $t(X)$ qui :

- (i) A toujours la même loi (ou quasiment la même loi) dans le cas d'égalité (autrement dit, lorsque θ vérifie $\gamma(\theta) = \gamma_0$, ce qui correspond donc au “bord” de l'hypothèse nulle);
- (ii) A tendance à prendre des valeurs plus petites que dans le cas d'égalité quand on est à l'intérieur de l'hypothèse nulle (autrement dit, lorsqu'on considère des θ pour lesquels $\gamma(\theta) < \gamma_0$);
- (iii) A tendance, à l'inverse, à prendre des valeurs plus grandes que le cas d'égalité quand on est dans l'hypothèse alternative (i.e. que $\gamma(\theta) > \gamma_0$),

alors il s'agit (normalement) d'une statistique de test pertinente, avec critère de suspicion à droite.

On a bien entendu un résultat similaire (avec critère de suspicion à gauche) lorsque les valeurs de $t(X)$ deviennent plus grandes à l'intérieur de Θ_0 , mais plus petites dans Θ_1 . \diamond

12.5 Tests bilatéraux

J'ai évoqué un peu plus haut l'existence d'un critère de suspicion « bilatéral » : c'est à ce point, et aux développements afférents, que nous allons nous intéresser dans cette section.

À titre d'introduction, reprenons l'exemple du vieux professeur grincheux évoqué ci-dessus (cf. exemple (CF')), à ceci près que cette fois-ci nous considérons que tout le monde s'accorde à dire que les différentes méthodes pédagogiques ne changent rien à l'hétérogénéité des résultats des élèves ; autrement dit, nous nous restreignons à $\Theta = \mathbb{R} \times \{\sigma_{\text{réf}}\}$. Dans ce cadre, le vieux professeur défend l'hypothèse nulle $\{\theta = \theta_{\text{réf}}\}$, avec $\theta_{\text{réf}} := (\mu_{\text{réf}}, \sigma_{\text{réf}})$. On voit qu'il y a qualitativement *deux* façons dont l'hypothèse nulle peut être fautive : soit l'hypothèse nulle est fautive parce que la vraie valeur de μ est $\mu_{\mathcal{V}} < \mu_{\text{réf}}$, soit l'hypothèse nulle est fautive parce qu'on aurait $\mu_{\mathcal{V}} > \mu_{\text{réf}}$.

Si nous reprenons ce que nous avons vu dans les sections précédentes, nous voyons que :

[¶]. Ou « quasiment la même loi », confer remarque (CL').

- S'il s'agissait de tester l'hypothèse nulle $\{\mu = \mu_{\text{réf}}\}$ contre l'hypothèse alternative $\{\mu > \mu_{\text{réf}}\}$, une statistique de test appropriée serait la moyenne M obtenue par les élèves, avec critère de suspicion à droite ;
- Si, en revanche, il s'agissait de tester l'hypothèse nulle $\{\mu = \mu_{\text{réf}}\}$ contre l'hypothèse alternative $\{\mu < \mu_{\text{réf}}\}$, on trouverait *aussi* que M est une statistique appropriée, mais cette fois-ci avec critère de suspicion à *gauche* !

On a donc toujours la même statistique de test M ; mais cette fois-ci il y a deux types de valeurs extrêmes qu'on a envie de rejeter : d'un côté, les valeurs trop grandes de M nous semblent suspectes du point de vue de l'hypothèse nulle parce qu'elles nous inciteraient plutôt à penser que $\mu > \mu_{\text{réf}}$; et de l'autre, les valeurs trop petites de M sont *aussi* suspectes, cette fois-ci parce qu'elles nous inciteraient plutôt à penser que $\mu < \mu_{\text{réf}}$! En résumé, on a envie de dire que plus les valeurs de M sont "extrêmes" (que ce soit vers la droite ou vers la gauche), plus elles sont suspectes... Eh bien, c'est précisément cela qu'on appellera « critère de suspicion bilatéral » :

Définition (CO'). Lorsqu'on associe un *critère de suspicion bilatéral* à une statistique de test T , cela signifie qu'on considère que les valeurs les plus extrêmes sont les plus suspectes :

- D'une part, asymptotiquement (au-delà d'une certaine valeur, je veux dire), on considère que plus la réalisation t_{\checkmark} de la statistique de test est grande, plus cela nous incite à pencher vers l'hypothèse alternative au détriment de l'hypothèse nulle ;
- D'autre part, toujours asymptotiquement (mais cette fois-ci *en deçà* d'une certaine valeur), on considère que plus la réalisation t_{\checkmark} de la statistique de test est petite, plus cela nous incite à pencher vers l'hypothèse alternative au détriment de l'hypothèse nulle !

♡

Remarque (CP'). Attention ! Il pourrait être tentant de croire que le critère ci-dessus équivaudrait à quelque chose comme « plus la valeur absolue $|T|$ est grande, plus c'est suspect ». Sauf que cela sous-entendrait qu'il devrait être *exactement aussi suspect* d'observer, mettons, $t_{\checkmark} = -177$ que $t_{\checkmark} = +177$... Or, cela n'a aucune raison d'être le cas ! (et ce ne le sera pas en général, confer plus loin). Tout ce que dit la définition ci-dessus, c'est que le niveau de suspicion associé à l'observation de la valeur t est maximal pour $t \rightarrow -\infty$, puis décroît lorsqu'on se rapproche des valeurs "centrales", avant d'augmenter à nouveau pour devenir, à nouveau, maximal lorsque $t \rightarrow +\infty$. ♣

Dans le cas d'un critère de suspicion bilatéral, puisque nous devons rejeter à la fois les valeurs de T trop petites et les valeurs trop grandes, il ne va pas y avoir un, mais deux seuils. Mais attention : c'est le risque *total* d'avoir un test positif qui doit demeurer $\leq \alpha$, par juste le risque associé à chaque seuil ! Cela nous amène donc à définir notre test à partir des intervalles de fluctuation :

Procédure (CQ'). Une partition de l'espace du paramètre caché en hypothèse nulle et hypothèse alternative étant donnée, un niveau $\alpha \in]0, 1/2]$ étant fixé, et une statistique de test *bilatérale* $t(X)$ étant choisie, on construit canoniquement, à partir de ces éléments, « le » test de niveau α de la façon suivante. On construit l'*intervalle d'acceptation* de notre test comme l'union des intervalles de fluctuation

!

des $\text{Loi}_\theta(t(X))$ pour θ dans l'hypothèse nulle :

$$I_{\text{acc}} := \text{Adh}\left(\bigcup_{\theta \in \Theta_0} [\text{Qtile}(\text{Loi}_\theta(t(X)); \alpha/2), \text{Qtile}(\text{Loi}_\theta(t(X)); 1 - \alpha/2)]\right). \quad (\text{CR}')$$

Quitte à remplacer δ_{acc} par son adhérence topologique (ce qu'on peut faire sans scrupule, car cela ne change rien du point de vue des applications pratiques), cela s'écrit aussi

$$I_{\text{acc}} = \left[\inf_{\theta \in \Theta_0} \text{Qtile}(\text{Loi}_\theta(t(X)); \alpha/2), \sup_{\theta \in \Theta_0} \text{Qtile}(\text{Loi}_\theta(t(X)); 1 - \alpha/2) \right]. \quad (\text{CS}')$$

Ensuite, on définit notre test comme $\{t(X) \notin I_{\text{acc}}\}$, autrement dit, on accepte d'hypothèse nulle si et seulement si la statistique de test tombe dans l'intervalle d'acceptation. \heartsuit

12.6 Régime asymptotique des tests

Intérêt de l'étude asymptotique des tests

Le problème avec les mathématiques, c'est que bien souvent la beauté des constructions théoriques se heurte à la dure réalité de calculs inextricables... Les tests n'échappent malheureusement pas à cette règle : l'information cruciale sur un test, son niveau (dans le cas booléen) ou sa p -valeur, est bien souvent trop compliquée à calculer exactement. En revanche, on peut parfois obtenir des approximations intéressantes... C'est ce qui conduit à l'objet de cette section : les tests *asymptotiques*. Dans le cadre asymptotique, on a souvent la possibilité de prouver des résultats mathématiques intéressants qui nous rassurent quant à la pertinence de tel ou tel test : autrement dit, on va mettre au point des méthodes statistiques qui sont censées bien marcher sous réserve que le nombre d'observations soit suffisamment grand.

Les tests asymptotiques, comme leur nom l'indique, appartiennent au domaine de la statistique... asymptotique : autrement dit, on va se placer dans un modèle statistique avec un (ou plusieurs) paramètres du modèle pour lequel on considère un certain régime asymptotique (cf. § 7.4), et pour toute valeur λ du paramètre du modèle, on considèrera un test booléen $T^{(\lambda)} : \mathcal{X}^{(\lambda)} \rightarrow \{\text{VRAI}, \text{FAUX}\}$; nous chercherons alors à dire comment se comporte le test $T^{(\lambda)}$ pour un certain régime asymptotique de λ : typiquement, on aura $\lambda \in \mathbb{N}$ et on étudiera la limite $\lambda \rightarrow \infty$.

Remarque (CT'). Attention!! Les mathématiciens adorent faire de la statistique asymptotique, parce que cela leur permet d'avoir de jolis théorèmes. Mais les ingénieurs, eux, doivent s'en méfier comme de la peste! En effet, il est bien beau de savoir que le résultat est valable si le nombre d'observations est « suffisamment grand », mais comment savoir si le nombre d'observations dont on dispose *en pratique* est *réellement* suffisamment grand?... C'est pourquoi un ingénieur raisonnable ne devrait appliquer une méthode asymptotique que s'il dispose au moins d'une méthode (idéalement rigoureuse, à défaut heuristique) pour savoir si l'approximation asymptotique est raisonnable ou non! \clubsuit

Tests asymptotiques

! **Définition (CU')** (Niveau asymptotique d'un test). Dans un modèle ayant un paramètre pour lequel on considère une asymptotique, soit $\text{Test}^{(\lambda)}$ un test booléen

de l'hypothèse $\{\theta \in \Theta_0\}$. On dit que le test est *asymptotiquement de niveau α* lorsque :

$$\forall \theta \in \Theta_0 \quad \overline{\lim}_{\lambda \rightarrow \infty} \mathbb{P}(\text{Test}^{(\lambda)} \mid \theta = \theta) \leq \alpha. \quad (\text{CV}') \quad \heartsuit$$

Remarque (CW'). Notez que la définition pour un test asymptotique requiert seulement qu'il y ait convergence pour chaque θ , mais qu'il est tout à fait possible que chaque θ converge à son propre rythme, et qu'on ne puisse trouver aucune borne de convergence valable pour tous les θ à la fois! \clubsuit

Proposition (CX'). Si un test $\text{Test}^{(\lambda)}$ est de niveau α (pour tout λ s'entend), alors il est asymptotiquement de niveau α (Par contre, la réciproque est fautive). \diamond

Démonstration. Immédiat. \diamond

Remarque (CY'). Il est tout-à-fait possible d'avoir un niveau asymptotique pour un test égal à... 0 (même si en pratique on préfère fixer un niveau α et s'y tenir, ou mieux, raisonner en termes de p -valeurs, confer chapitre suivant) : cela n'a en fait rien de bien extraordinaire, pas plus que la notion de consistance que nous allons voir ci-après. \clubsuit

Consistance

Toujours dans le cadre du régime asymptotique, un critère de qualité très important, lorsqu'on parle de tests, est la notion de *consistance* :

Définition (CZ'). Un test $\text{Test}^{(\lambda)}$ de l'hypothèse $\{\theta \in \Theta_0\}$ est dit *consistant* lorsque, pour tout $\theta \notin \Theta_0$, $\mathbb{P}(\text{Test}^{(\lambda)} \mid \theta = \theta) \rightarrow 1$ lorsque $\lambda \rightarrow \infty$. \heartsuit !

Remarque (DA'). On a l'impression avec la définition ci-dessus que la philosophie du test est inversée : asymptotiquement, on n'a plus aucun risque d'erreur dans le cas où $\theta_{\checkmark} \in \Theta_1$, alors qu'il reste toujours un risque α de faux positif si $\theta_{\checkmark} \in \Theta_0$... Mais, comme dit dans la remarque (CT'), la « vraie » statistique est en fait la statistique *non* asymptotique (ou au moins *uniformément* asymptotique) : les résultats asymptotiques sont juste là pour nous rassurer mathématiquement sur ce qu'on fait! \smile \clubsuit

Remarque (DB'). Tous les tests que nous serons amenés à voir dans ce cours, lorsqu'il y a un régime asymptotique, seront en fait consistants. On peut en fait montrer que, au moins en ce qui concerne les modèles d'échantillonnage, sous des hypothèses générales assez souples sur le modèle probabiliste et la structure de l'hypothèse nulle, il est toujours possible de trouver un test consistant pour l'hypothèse nulle. \clubsuit

Exemple (DC') (Test du χ^2). Le test du χ^2 [III] est un test extrêmement classique pour tester l'indépendance de deux variables qualitatives. Soient \mathcal{X} et \mathcal{Y} des ensembles finis non vides et μ_{\checkmark} une loi portée par $\mathcal{X} \times \mathcal{Y}$. Le paramètre caché de notre modèle probabiliste est μ_{\checkmark} elle-même, et l'expérience statistique consiste à observer un nombre n librement fixé de réalisations i. i. d. de la loi μ_{\checkmark} . Ces réalisations

[III]. 'χ' est la lettre grecque « chi » — attention, on prononce /ki/ —; à l'oral, on lit « test du chi-deux » — prononcé /kidø/.

sont notées resp. $(x_{1\checkmark}, y_{1\checkmark}), \dots, (x_{n\checkmark}, y_{n\checkmark})$, vues comme les réalisations des variables aléatoires $(X_1, Y_1), \dots, (X_n, Y_n)$.

Les test du χ^2 se propose de tester l'hypothèse nulle que, sous la véritable loi μ_{\checkmark} , les variables X et Y sont indépendantes. Autrement dit, notant $\mu(x, y)$ la probabilité sous μ de tirer le couple (x, y) , resp. $\mu(x) := \sum_{y' \in \mathcal{Y}} \mu(x, y')$ la probabilité que la première variable soit x , resp. $\mu(y) := \sum_{x' \in \mathcal{X}} \mu(x', y)$ la probabilité que la seconde variable soit y , on veut tester l'hypothèse que

$$\forall (x, y) \in \mathcal{X} \times \mathcal{Y} \quad \mu_{\checkmark}(x, y) = \mu_{\checkmark}(x)\mu_{\checkmark}(y). \quad (\text{DD}')$$

Pour ce faire, on introduit une statistique de test appelée *statistique du χ^2* , définie de la façon suivante. Notons $N(x, y)$ la variable aléatoire comptant le nombre d'observations i pour lesquelles on a $\{(X_i, Y_i) = (x, y)\}$, resp. $N(x) = \sum_{y' \in \mathcal{Y}} N(x, y')$ le nombre de fois qu'on a $\{X_i = x\}$, resp. $N(y) = \sum_{x' \in \mathcal{X}} N(x', y)$ le nombre de fois qu'on a $\{Y_i = y\}$. Alors la statistique du χ^2 est

$$T := n \sum_{(x, y) \in \mathcal{X} \times \mathcal{Y}} \frac{N(x, y)^2}{N(x)N(y)} - n. \quad (\text{DE}')$$

(En prenant la convention $0/0 = 0$ si nécessaire). (On peut montrer que T est toujours positive).

Le résultat fondamental sur lequel s'appuie le test du χ^2 est le suivant :

Proposition (DF'). Si μ_{\checkmark} vérifie l'hypothèse d'indépendance, alors quand $n \rightarrow \infty$, on a

$$\text{Loi}_{\checkmark}(T^{(n)}) \rightarrow \text{ChiDeux}((|\mathcal{X}| - 1) \times (|\mathcal{Y}| - 1)) \quad (\text{DG}')$$

lorsque μ_{\checkmark} associe une masse non nulle à tous les éléments de \mathcal{X} et \mathcal{Y} , et plus généralement

$$\text{Loi}(T^{(n)}) \rightarrow \text{ChiDeux}((|\text{supp}_{\mathcal{X}} \mu_{\checkmark}| - 1) \times (|\text{supp}_{\mathcal{Y}} \mu_{\checkmark}| - 1)), \quad (\text{DH}')$$

où $|\text{supp}_{\mathcal{X}} \mu_{\checkmark}|$ est le nombre de $x \in \mathcal{X}$ pour lesquels $\mu_{\checkmark}(x) > 0$, resp. $|\text{supp}_{\mathcal{Y}} \mu_{\checkmark}|$ est le nombre de $y \in \mathcal{Y}$ pour lesquels $\mu_{\checkmark}(y) > 0$. \diamond

On en déduit en particulier que pour un niveau t fixé, lorsque n tend vers l'infini,

$$\overline{\lim} \mathbb{P}(T^{(n)} > t) \leq \mathbb{P}(\text{ChiDeux}((|\mathcal{X}| - 1) \times (|\mathcal{Y}| - 1)) > t) \quad (\text{DI}')$$

(où l'on a admis le fait que $\mathbb{P}(\text{ChiDeux}(d) > t)$ est croissant lorsque d augmente — ce qui est très facile à vérifier dès qu'on connaît la définition de la loi du ChiDeux).

Par conséquent, cela suggère le protocole de test suivant : étant fixé un niveau α , chercher le plus grand t tel que $\mathbb{P}(\text{ChiDeux}((|\mathcal{X}| - 1) \times (|\mathcal{Y}| - 1)) \geq t) \geq \alpha$, puis rejeter l'hypothèse nulle lorsque $T > t$.

La proposition (DF') nous dit alors précisément que ce test sera asymptotiquement de niveau α . Par ailleurs, il est aisé de voir que le test proposé ci-dessus est consistant. \clubsuit

Exemple (DJ') (Test du χ^2 , suite). Dans le test du χ^2 proposé ci-dessus, on peut cependant démontrer que le niveau $\alpha^{(n)}$ de ce test pour une valeur donnée de n ne converge pas, quand n tend vers l'infini, vers α (il n'y a donc pas « uniforme consistance » du test), mais vers une valeur un peu plus élevée.

Néanmoins, dans le cas de ce modèle, on peut en fait trouver une règle qui permet d'obtenir une estimation quantitative sur la convergence du test vers son asymptotique. Pour ce faire, un certain seuil N_{\min} étant introduit (on prend généralement $N_{\min} = 5$), on considère la statistique modifiée

$$\check{T} := \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} n \frac{\check{N}(x,y)^2}{N(x)N(y)} - n, \quad (\text{DK}')$$

où l'on pose

$$\check{N}(x,y) = \begin{cases} N(x,y) & \text{pour } N(x,y) \geq N_{\min} \\ 0 & \text{pour } N(x,y) < N_{\min} : \end{cases} \quad (\text{DL}')$$

cette statistique a la même asymptotique que T sous l'hypothèse nulle, et conduit donc elle aussi à un test de niveau α . En outre, quoique le test fondé sur \check{T} soit un peu moins puissant que celui fondé sur T , il est lui aussi consistant. Et surtout, cette fois-ci notre test a un niveau borné *uniformément* sur Θ_0 de façon *non* asymptotique, uniformément en n , le niveau en question pouvant être fixé arbitrairement proche de α pour peu qu'on ait choisi N_{\min} suffisamment grand.

En pratique, pour $N_{\min} = 5$ et pour des valeurs de n de quelques dizaines ou quelques centaines, le niveau réel ainsi obtenu pour le test modifié est raisonnablement proche de α (à condition qu'on n'ait pas choisi α trop petit). \clubsuit

Chapitre 13

La p -valeur

13.1 p -valeur associée à une statistique de test

Dans le chapitre précédent, vous venez de voir comment fonctionne un test booléen. Maintenant “oubliez tout ce que vous venez d’apprendre^[*]” : en pratique, nous allons utiliser une autre façon de formuler les tests, plus précise que la notion “binaire” de test booléen. À savoir : la p -valeur !

Pour comprendre la raison derrière l’introduction du concept de p -valeur, revenons sur la procédure que nous avons suivie lorsqu’il y a une statistique de test T avec critère de suspicion unilatéral. (Pour fixer les idées, disons ici qu’on suspecte les *petites* valeurs de T). Notre procédure fonctionnait ainsi : nous étant fixé un niveau de tolérance α pour le risque de première espèce, nous avons choisi un seuil $t_{\text{seuil}}(\alpha)$ tel que

$$\forall \theta \in \Theta_0 \quad \mathbb{P}_\theta(T < t_{\text{seuil}}(\alpha)) \leq \alpha; \quad (\text{DM}')$$

et, tant qu’à faire, nous avons choisi ce seuil aussi grand que possible : autrement dit, nous avons pris

$$t_{\text{seuil}}(\alpha) = \inf\{\text{Qtile}(\text{Loi}_\theta(T); \alpha) \mid \theta \in \Theta_0\}, \quad (\text{DN}')$$

où $\text{Qtile}(\bullet; \alpha)$ désigne le quantile de niveau α (dans le cas d’un critère de suspicion à droite, nous aurions pris cette fois-ci le supremum des quantiles de niveau $1 - \alpha$). La réponse de notre test (booléen) est alors VRAI si $t_{\checkmark} < t_{\text{seuil}}(\alpha)$, resp. FAUX si $t_{\checkmark} \geq t_{\text{seuil}}(\alpha)$.

Par conséquent, on voit que la réponse de notre test booléen résulte de la conjonction de deux facteurs :

- La valeur effectivement observée t_{\checkmark} de la statistique de test ;
- Le niveau de risque α (via la détermination de $t_{\text{seuil}}(\alpha)$).

[*]. Façon de parler...! ☹ Vous devez *quand même* maîtriser la notion de test booléen : d’une part, parce qu’il arrive parfois qu’on le rencontre en pratique malgré tout ; et d’autre part, parce que comprendre les tests booléens est une étape extrêmement utile pour comprendre les p -valeurs (l’objet du présent chapitre) ainsi que les intervalles de confiance (chap. 14). Ce que je veux dire par « oubliez les tests booléens », c’est que, aussi loin qu’il s’agit d’étudier la plausibilité d’une hypothèse nulle, l’approche par p -valeur est à la fois plus rapide, plus précise et plus riche : c’est pourquoi, dans la pratique industrielle, cette pratique méritera d’être systématiquement préférée dans la mesure du possible ! Au demeurant, c’est bien ce que font les logiciels de statistique : ils expriment toujours les résultats des tests fréquentistes d’hypothèse en termes de p -valeur, jamais en termes booléens ! ☺

Il se trouve que la façon de conjuguer ces deux facteurs peut être faite sans jamais avoir à expliciter $t_{\text{seuil}}(\alpha)$. En effet, pour tout $t \in \mathbb{R}$, introduisons la quantité

$$p(t) := \sup\{\mathbb{P}_\theta(T \leq t) \mid \theta \in \Theta_0\}. \quad (\text{DO}')$$

On a alors le résultat suivant, pas très difficile à démontrer :

Lemme (DP').

$$t_{\text{seuil}}(\alpha) = \inf\{t \in \mathbb{R} \mid p(t) > \alpha\}. \quad (\text{DQ}')$$

◇

Comme, en outre, la fonction $t \mapsto p(t)$ est évidemment croissante, il découle du lemme (DP') que, sauf cas particulier^[†] (que nous négligerons dans la suite), la condition $\{t_{\mathcal{V}} < t_{\text{seuil}}(\alpha)\}$ est équivalente à la condition $\{p(t_{\mathcal{V}}) \leq \alpha\}$, laquelle ne demande cette fois-ci que l'évaluation de $p(t_{\mathcal{V}})$, sans avoir à calculer explicitement $t_{\text{seuil}}(\alpha)$!

La valeur $p(t_{\mathcal{V}})$, que dorénavant nous noterons $p_{\mathcal{V}}$, permet donc de savoir si le test booléen de niveau α associé à notre statistique de test (pour le critère de suspicion considéré) est positif ou négatif, ce qui montre clairement son intérêt ! C'est pourquoi on lui donne un nom : on l'appelle la *p-valeur* de notre test. Comme la définition ci-dessous l'explique, on peut comprendre la *p-valeur* comme « la probabilité (supremale), lorsque l'hypothèse nulle est vraie, d'obtenir un résultat au moins aussi suspect que celui qui a été effectivement observé » :

! **Définition (DR')** (*p-valeur pour une statistique de test, cas unilatéral*). Soit T une statistique destinée à servir de statistique de test pour l'hypothèse nulle $\{\theta \in \Theta_0\}$, et soit $t_{\mathcal{V}}$ la réalisation de cette statistique. Alors la *p-valeur*^[‡] de notre test est la (réalisation de) statistique $p_{\mathcal{V}}$ définie ainsi : c'est le supremum, sur toutes les valeurs θ appartenant à l'hypothèse nulle, de la probabilité que, sous \mathbb{P}_θ , la statistique de test T prenne une valeur au moins aussi suspecte que la valeur $t_{\mathcal{V}}$ effectivement observée. Autrement dit, c'est la quantité suivante :

— Si ce sont les petites valeurs de T qui sont suspectes (« test à gauche ») :

$$p_{\mathcal{V}} := \sup_{\theta \in \Theta_0} \mathbb{P}(T \leq t_{\mathcal{V}} \mid \theta = \theta). \quad (\text{DS}')$$

[†]. Le « cas particulier » en question concerne la situation où on a *exactement* $p(t_{\mathcal{V}}) = \alpha$. Néanmoins, on peut considérer comme assuré que cela n'arrivera jamais en pratique, car les valeurs α susceptibles d'intervenir dans vos analyses seront toujours des seuils de références fixés indépendamment du modèle (p. ex. 5 %, 1 %, ...,), seuils auxquels la *p-valeur*, qui est en général un nombre réel tout ce qu'il y a de plus quelconque, n'aura aucune raison d'être rigoureusement égale... ! D'où ma volonté de ne pas s'ennuyer avec ce cas d'égalité. En fait, il serait *possible* de traiter rigoureusement l'équivalence entre théorie des tests booléens et approche par *p-valeurs* d'une façon qui ne présente aucune exception au niveau des cas d'égalité ; mais cela obligerait alors à revenir sur la théorie des tests booléens, et même sur celle des intervalles de fluctuation, pour savoir dans quel cadre on doit ouvrir ou pas les bornes de tel ou tel intervalle : ce qui représenterait beaucoup de complexité supplémentaire pour un apport dérisoire... ∴ Cet apport serait en outre d'autant plus dérisoire que nous savons très bien (merci George Box ☺) que les nombres réels intervenant dans les modèles statistiques et les observations sont eux-mêmes inévitablement entachés d'une certaine imprécision : il serait donc passablement stupide de compliquer le cours pour la seule satisfaction de savoir traiter rigoureusement le cas d'une égalité parfaite qui, non seulement n'a aucune raison de se produire, mais n'a de toutes façons aucune réalité concrète ! ☺

[‡]. La *p-valeur* est traditionnellement notée par la lettre '*p*', d'où son nom. Dans le cadre de ce cours, nous utiliserons plutôt la notation ' $p_{\mathcal{V}}$ ', reflétant ainsi le fait que la *p-valeur* est en fait la réalisation d'une certaine statistique — statistique que nous pourrions alors noter ' P ' lorsqu'on la verra en tant que variable aléatoire).

— Si ce sont les grandes valeurs de T qui sont suspectes (« test à droite ») :

$$p_{\mathcal{J}} := \sup_{\theta \in \Theta_0} \mathbb{P}(T \geq t_{\mathcal{J}} \mid \theta = \theta). \tag{DT'}$$

♡

Remarque (DU'). Attention! Lorsqu'on définit la p valeur, on considère (le supremum de) la probabilité d'être *au moins aussi* suspect que l'observation; alors que lorsqu'on définissait le test booléen, on considèrerait le fait d'être *strictement plus suspect* que (l'infimum de) le quantile de la loi de T ! On a donc une inégalité *large* dans le cadre de la définition (DR'), contrairement à la procédure (BD') qui requerrait une inégalité *stricte*... [§] ♣

On peut adapter la définition de la p -valeur aux tests bilatéraux :

Définition (DV') (p -valeur pour les tests bilatéraux). Dans le cas d'une statistique de test T avec critère de suspicion bilatéral, la p -valeur déduite de la statistique de test T est définie de la façon suivante, où $t_{\mathcal{J}}$ est la réalisation de T (en pratique, on utilisera plutôt de la majoration donnée par le membre de droite [¶]) :

$$p_{\mathcal{J}} := 2 \sup_{\theta \in \Theta_0} (\mathbb{P}_{\theta}(T \leq t_{\mathcal{J}}) \wedge \mathbb{P}_{\theta}(T \geq t_{\mathcal{J}})) \leq 2 \left(\sup_{\theta \in \Theta_0} \mathbb{P}_{\theta}(T \leq t_{\mathcal{J}}) \wedge \sup_{\theta \in \Theta_0} \mathbb{P}_{\theta}(T \geq t_{\mathcal{J}}) \right). \tag{DW'}$$

♡

Remarque (DX'). La formule ci-dessus peut avoir l'air passablement obscure; mais en réalité, elle s'interprète de façon tout à fait logique, comme je vais vous l'expliquer à présent.

De manière générale, l'idée de la p -valeur, c'est qu'on considère la probabilité d'être *au moins aussi suspect* que $t_{\mathcal{J}}$. Dans le cas d'un critère de suspicion unilatéral, disons à droite, c'était facile à interpréter : cela signifiait juste « probabilité d'être au moins aussi *grand* que $t_{\mathcal{J}}$ ». Maintenant, si on est *bilatéral*, comment faut-il traduire les choses?...

Prenons un exemple pour comprendre. Mettons ici que la probabilité d'être strictement inférieur à $t_{\mathcal{J}}$ vaut 11 %, que celle d'être exactement égal à $t_{\mathcal{J}}$ vaut 2 %, et que celle d'être strictement supérieur à $t_{\mathcal{J}}$ vaut 87 %. En l'occurrence, puisqu'il y a moins de probabilité à gauche de $t_{\mathcal{J}}$ qu'à droite de $t_{\mathcal{J}}$, ce qui rend ce $t_{\mathcal{J}}$ -là suspect, c'est le fait qu'il soit particulièrement *petit*! Ainsi, la probabilité d'être au moins aussi suspect que $t_{\mathcal{J}}$ *en étant suspect dans le même sens que lui* vaut ici 13 %. Cette

[§]. En fait, ce n'est pas si paradoxal qu'il n'y paraît à première vue. Certes, pour les tests booléens, on avait à considérer des événements du type $\{T < t_{\text{seuil}}\}$, alors que pour la p valeur, on considère des événements du type $\{T \leq t_{\mathcal{J}}\}$. Mais, dans le premier cas, la variable aléatoire T était destinée à être remplacée *in fine* par $t_{\mathcal{J}}$; alors que dans le second, ' T ' fait en réalité référence à un "autre" tirage de la statistique de test, si on *relançait* l'expérience, tirage qu'on pourrait appeler ' T' '... Dès lors, si on se place au niveau des réalisations, on s'aperçoit qu'on considère dans le premier cas l'affirmation « $t_{\mathcal{J}} < t_{\text{seuil}}$ », et dans le second cas, l'affirmation « $t'_{\mathcal{J}} \leq t_{\mathcal{J}}$ ». Le membre commun entre ces deux affirmations étant $t_{\mathcal{J}}$, si on veut réellement rendre les deux affirmations comparables, on va placer " $t_{\mathcal{J}}$ " du côté gauche à chaque fois : et on se rend compte alors que les expressions deviennent resp. « $t_{\mathcal{J}} < t_{\text{seuil}}$ » et « $t_{\mathcal{J}} \not< t'_{\mathcal{J}}$ » : ainsi, en fait, on s'intéressait bien, dans les deux cas, à des questions sur le fait que $t_{\mathcal{J}}$ soit strictement inférieur, ou pas, à une valeur de comparaison!

[¶]. Dans les situations rencontrées en pratique, en effet, la comparaison entre les valeurs $\mathbb{P}_{\theta}(T \leq t_{\mathcal{J}})$ et $\mathbb{P}_{\theta}(T \geq t_{\mathcal{J}})$ a lieu dans le même sens pour toutes les valeurs de $\theta \in \Theta_0$, de sorte que la majoration ne perd alors aucunement en précision.

valeur correspond au minimum « $\mathbb{P}(T \leq t_{\checkmark}) \wedge \mathbb{P}(T \geq t_{\checkmark})$ » qui apparaît dans notre formule : selon que t_{\checkmark} est situé du côté gauche de la distribution de T ou du côté droit de la probabilité, la « probabilité d'être au moins suspect que lui en penchant du même côté » sera soit $\mathbb{P}(T \leq t_{\checkmark})$, soit $\mathbb{P}(T \geq t_{\checkmark})$: et dans tous les cas, ce sera toujours *la plus petite* des deux probabilités en question : étant donné que dire que t_{\checkmark} est suspect *du fait de sa petitesse* (plutôt que du fait de sa grandeur) devient précisément à dire que « être plus petit que t_{\checkmark} » est quelque chose de plus rare que « être plus grand que t_{\checkmark} » !

Soit ; nous avons donc compris ce fameux symbole de minimum. Maintenant, d'où vient le facteur 2 ? Eh bien, si nous reprenons l'exemple ci-dessus, nous avons dit que la probabilité d'être au moins aussi suspect que t_{\checkmark} *en étant suspect du fait qu'on est petit* valait 13 % : d'accord ; mais on pourrait tout aussi bien être au moins aussi suspect que t_{\checkmark} *en étant suspect du fait qu'on est grand* ! (Puisqu'une valeur "très très" grande sera bel et bien plus suspect qu'une valeur juste "très" petite). Mais que vaut cette seconde probabilité ? En fait, comme nous avons décidé de répartir le risque de façon symétrique entre les deux queues de la distribution, « être plus suspect qu'un niveau donné du fait qu'on est grand » doit avoir la même probabilité que « être plus suspect que ce même niveau donné du fait qu'on est petit », d'où la multiplication par 2 ^[III].

Donc en fait, on devrait plutôt écrire le membre du milieu dans la formule comme

$$\sup_{\theta \in \Theta_0} (2 \times (\mathbb{P}_{\theta}(T \leq t_{\checkmark}) \wedge \mathbb{P}_{\theta}(T \geq t_{\checkmark}))), \quad (\text{DY}')$$

pour bien rendre visible le fait qu'on est toujours en train de considérer la probabilité d'être au moins aussi suspect : mais cela rend la formule difficilement lisible du fait du grand nombre de parenthèses... ☹

Un mot, enfin, en ce qui concerne le membre de droite de la formule. Une fois qu'on a compris l'idée que « la probabilité d'être au moins aussi suspect, dans le cas bilatéral, c'est le double de celle des deux valeurs qui est la plus faible entre, d'une part, la probabilité d'être au moins aussi petit, et d'autre part, la probabilité d'être au moins aussi grand », on interprète le membre de droite en disant : « la (majoration de la) p -valeur dans le cas bilatéral peut se définir ainsi : on regarde la p -valeur qu'on aurait eue dans le cas d'un critère de suspicion à gauche, d'une part ; puis la p -valeur qu'on aurait eue dans le cadre d'un critère de suspicion à droite, d'autre part ; on prend *la plus petite* de ces deux p -valeurs (puisque c'est celle qui nous dit quel est le critère de suspicion qui rendrait notre observation effective la plus suspecte) ; et on multiplie par 2 pour tenir compte du fait qu'il y a *deux* façons possibles, considérées de façon symétrique, d'être suspect ». Ouf ! ☺ ☼

Remarque (DZ'). On remarquera que la formule donnée ci-dessus pour déterminer la p -valeur (ainsi que sa majoration) dans le cas d'un test bilatéral peut produire, dans certains cas, des valeurs *strictement supérieures* à 1 : ce qui n'est cohérent, ni avec l'interprétation de la p valeur comme un supremum de probabilités, ni avec le reste du chapitre (où il est toujours entendu que les p -valeurs sont ≤ 1), ni avec

[III]. Si on veut vraiment rentrer dans le détail, il s'avère malheureusement que la phrase que je viens d'écrire requiert, pour être parfaitement correcte, la prise en compte d'un certain nombre de subtilités lorsqu'on considère des distributions de probabilité discrètes... Je ne rentrerai pas dans ces subtilités ici, vous demandant simplement de me croire sur parole quand je dis que le raisonnement ci-dessus justifiant la multiplication par 2 est, malgré son allure un peu bancal, bel et bien solide dans le cadre de la façon dont ce cours traite les statistiques de test ! ☺

l'usage fait dans les logiciels (qui ne renvoient jamais de p -valeur supérieure à 1). Pour autant, la formule n'est pas *techniquement* fautive pour autant, car rien n'empêche de construire la théorie des p -valeurs en autorisant les valeurs > 1 ; et tous les théorèmes que nous avons énoncés restent alors encore valables... En fait, dans le paradigme où on autorise les p -valeurs à dépasser 1, on s'aperçoit que, si P est une p -valeur, alors $P \wedge 1$ est *aussi* une p -valeur, ce qui explique pourquoi il est en fait inutile de s'encombrer de valeurs > 1 . Un corolaire de ce point est que, dans la formule ci-dessus définissant la p -valeur, on pourrait rajouter « $\wedge 1$ » sans rien changer à la validité de ce qui est écrit. (Du coup, si à un moment vous calculez par cette formule une p -valeur de, mettons, 1,04, vous pouvez tout à fait légitimement préciser « que l'on peut tronquer à 1 », et considérer ensuite que la *véritable* p -valeur est plutôt 1). ♣

Comme annoncé, à partir de la p -valeur, on peut avoir la réponse à n'importe quel test booléen :

Théorème (EA') (Lien entre p -valeur et résultats des tests booléens). *Une statistique de test (et un critère de suspicion (unilatéral) associé) étant fixée, le test de niveau α défini à partir de cette statistique sera positif si, et seulement si, la p -valeur associée à cette statistique est inférieure ou égale à α ^[**] ^[††]* ◇ !

Remarque (EB'). Ainsi, la p -valeur nous permet d'un seul coup de donner la réponse des tests booléens pour *n'importe quelle* valeur de α : soit α est plus petit que la p -valeur et alors on accepte l'hypothèse nulle ; soit α est plus grand, et alors on la rejette ! Du coup, plutôt que de fixer arbitrairement un α dès le départ, il est beaucoup plus malin de calculer la p -valeur, et de considérer que celle-ci donne une réponse *graduée* à notre test. ♣ !

13.2 Notion intrinsèque de p -valeur

Cette courte section est entièrement “culturelle”, et peut donc être sautée sans encombre ☹

Lorsque nous avons défini les tests booléens, nous avons d'abord dit rigoureusement ce que signifiait « être un test booléen de niveau α » *dans l'absolu*, avant de voir, *en pratique*, comment on construisait de tels tests booléens à partir de statistiques de test. Dans le traitement ci-dessus que nous avons fait des p -valeurs, par contre, nous nous sommes appuyés de façon *fondamentale* sur l'idée de statistique de test : notre définition n'a en effet pas de sens en l'absence de statistique de test ! Est-ce à dire qu'il n'y a pas de définition « absolue » de la notion de p -valeur ?...

[**]. En fait, de la façon dont ont été construits les tests booléens dans ce cours (à savoir, en incluant systématiquement les valeurs seuil dans la zone de négativité), ce théorème n'est pas cent pour cent rigoureux... Un théorème exact serait le suivant : d'une part, si le test booléen de niveau α est positif, alors la p -valeur sera inférieure ou égale à α ; d'autre part, le test {la p -valeur est $\leq \alpha$ } est lui-même un test booléen de niveau α . (Mais il est possible que, lorsque la statistique de test est précisément égale à la valeur seuil, le test construit à partir de la p -valeur soit positif, alors que par convention le test construit à partir de la statistique de test est négatif au niveau du seuil). Néanmoins, les cas où l'équivalence écrite dans l'énoncé ne tient pas sont extrêmement particuliers, donc il vaut mieux retenir cet énoncé-là, qui de toutes façons est “quasiment vrai”, seuls les cas où la p -valeur est rigoureusement égale à α pouvant poser problème.

[††]. On remarquera que, pour le coup, le test consistant à regarder si la p -valeur est $\leq \alpha$ est un test booléen qui sera *positif* au niveau de la valeur seuil, contrairement à ce qui a été fait dans le chapitre précédent. Il y a des raisons techniques profondes à ce que j'aie préféré écrire les choses ainsi, mais sur lesquelles je préfère ne pas m'apesantir... En pratique, soyez rassurés : si jamais vous rencontrez une p -valeur égale à *exactement* 6 %, et que vous écrivez « puisqu'on est pile au niveau du seuil, le test est négatif, donc on accepte l'hypothèse nulle au risque 6 % », je considérerai cela comme correct aussi ! ☹ Et de toutes façons, ce genre d'égalité exactes n'arrive normalement *jamais* dans la vie réelle... ☹

En fait, si : on peut bel et bien définir la notion de p -valeur sans passer par des statistiques de test ; mais comme c'est un usage assez rare en pratique, il m'a semblé dispensable de préciser cette définition dans le cadre de ce cours. En fait, il y a même *deux* définitions possibles de la notion de p -valeur : la première, la plus usuelle dans les ouvrages de statistique, où on généralise légèrement l'idée de « test reposant sur une statistique donnée » en celle de « famille de tests imbriqués » ; et la seconde, très générale, où la notion de « statistique de p -valeur » est définie de façon complètement directe, qui est ma préférée $\hat{~}$. Voici ces deux définitions :

Définition (EC') (p -valeur associée à une famille de tests imbriqués). Soit un modèle statistique avec les notations standards, et une division de l'espace du paramètre caché en des parties nulle et alternatives Θ_0 et Θ_1 . On dit qu'un ensemble $(Test_\alpha)_{\alpha \in]0, 1[}$ est une *famille de tests imbriqués* pour l'hypothèse nulle $\{\theta \in \Theta_0\}$ lorsque :

- D'une part, pour tout $\alpha \in]0, 1[$, $Test_\alpha$ est un test de notre hypothèse nulle au niveau de risque $^{\{\ddagger\}}$ α ;
- D'autre part, pour tous $\alpha < \alpha'$, $Test_\alpha \implies Test_{\alpha'}$: autrement dit, si une observation x donne lieu à une réponse positive pour le test $Test_\alpha$, elle devra aussi donner une réponse positive pour tous les $Test_{\alpha'}$ avec $\alpha' > \alpha$.

Dans ces conditions, la p -valeur associée à l'observation $x_{\mathcal{V}}$ pour notre famille de tests est définie comme :

$$p_{\mathcal{V}} := \inf\{\alpha \in]0, 1[\mid test_\alpha(x_{\mathcal{V}})\}; \tag{ED'}$$

autrement dit, c'est la « première » valeur de α pour laquelle, lorsqu'on relève progressivement le niveau de risque du test, on se met à avoir un test positif pour l'observation $x_{\mathcal{V}}$. \heartsuit

Définition (EE') (Notion générale de statistique de p -valeur). De manière générale, un modèle statistique et une hypothèse nulle étant donnés (avec les notations standard), une *statistique de p -valeur* est une statistique P à valeurs dans $[0, 1]$ telle que

$$\forall \alpha \in [0, 1] \quad \forall \theta \in \Theta_0 \quad \mathbb{P}_\theta(P < \alpha) \leq \alpha^{[*]}. \tag{EF'}$$

\heartsuit

Proposition (EG') (Test par p -valeur et test booléen, version sans statistique de test). Une *statistique de p -valeur* P étant donnée pour un test de l'hypothèse nulle \mathcal{H}_0 , pour $\alpha \in]0, 1[$ un niveau de risque, le test booléen $\{P \leq \alpha\}$ est un test de niveau α de l'hypothèse \mathcal{H}_0 . \diamond

13.3 Interprétation d'une p -valeur

!! **Principe (EH')**. La p -valeur étant la probabilité d'obtenir un résultat au moins aussi suspect dans le cas où l'hypothèse nulle est vraie, plus la p -valeur est petite et plus on sera enclin à rejeter l'hypothèse nulle. Voir aussi le paragraphe « Interprétation d'une p -valeur » un peu plus loin. \diamond

Remarque (EI'). Lorsqu'on arrondit une p -valeur, il est préférable de le faire toujours par excès, afin d'être conservatif concernant le niveau du test. (Confer remarque (BJ')). \clubsuit

!

$^{\{\ddagger\}}$. À noter que, exceptionnellement, on est amené ici à considérer des niveau de risque plus grands que 1/2.

$^{\{*\}}$. Dans le jargon technique, deux lois Q_0 et Q_1 sur \mathbb{R} étant données, on dit que Q_1 domine stochastiquement Q_0 lorsque la fonction de répartition de Q_0 est partout supérieure ou égale à la fonction de répartition de Q_1 (ce qui signifie, intuitivement, que Q_1 prend ses valeurs « plus à droite » que Q_0). Avec ce vocabulaire, la condition de la formule peut se reformuler en disant que, pour tout θ participant de l'hypothèse nulle, la loi de la P sous \mathbb{P}_θ doit être stochastiquement dominée par $\text{Unif}^{\text{me}}(0, 1)$.

Définition (EJ'). En analyse statistique, lorsqu'on souhaite exprimer l'idée qu'une observation semble suffisamment incompatible avec une hypothèse nulle (autrement dit, que la p -valeur du test soit suffisamment petite) pour qu'on puisse conclure au rejet de cette hypothèse nulle, on parle de résultat « significatif » : par exemple, une phrase comme « les résultats de la promotion cobaye sont significativement meilleurs que ce à quoi on aurait pu s'attendre avec l'ancienne méthode » signifie, en termes plus techniques : « le test de l'hypothèse alternative selon laquelle la nouvelle méthode pédagogique donne des résultats strictement meilleurs que l'ancienne, lorsqu'on l'applique aux données de la promotion cobaye, est positif pour un niveau de risque très faible (ou, de façon équivalente, a une très petite p -valeur) ».

Attention ! Dans le vocabulaire courant, « significativement meilleur » n'exprime en général pas le *niveau de confiance* qu'on peut avoir sur le fait qu'on est effectivement meilleur, mais à *quel point* on est meilleur (cela correspond au concept de « taille d'effet », cf. § 17.5) : l'usage statistique de ce mot est donc bien particulier, et on doit être très prudent avant d'utiliser le mot « significatif » devant des non-statisticiens ! (et, à l'inverse, on veillera à ce que les statisticiens, eux, pourront être induits en erreur si on a voulu employer l'adjectif « significatif » au sens du langage courant...). ♡

Remarque (EK'). La p -valeur étant une grandeur quantitative, son interprétation qualitative en termes de « résultat plus ou moins significatif » ne saurait être définie de façon canonique. Cependant, l'usage a consacré certains seuils et appellations qu'il est préférable d'avoir vu au moins une fois afin de ne pas trop s'en éloigner : voir la table 13.1. Vous n'êtes pas obligés de les suivre à la lettre (de sorte que je ne vous demande pas de les apprendre par cœur), mais cela vous donnera au moins une idée de ce qui peut être considéré comme un usage correct : par exemple, dire qu'une p -valeur à 6,2 % est « significative » ne peut guère être considéré comme fautif ; par contre, dire qu'elle est « très significative » est assez clairement inacceptable... ! ♣

Remarque (EL'). Vous serez peut-être surpris, dans la table 13.1, de voir qu'on considère qu'il n'y a *aucune* significativité pour une p -valeur supérieure à 20 %, alors qu'une p -valeur peut potentiellement prendre n'importe quelle valeur entre 0 et 1 : dès lors, ne semblerait-il pas logique de commencer à trouver un peu de significativité dès lors que la p -valeur est inférieure à 1/2 ?...

Eh bien non, ce serait là une conclusion parfaitement inepte ! Parler de significativité pour une p -valeur supérieure à 20 %, ce serait considérer que, lorsqu'on voit un événement qui avait 21 % de chances de se produire sous l'hypothèse nulle, ce serait suffisant pour dire « ouh là, 21 % de chances, c'est vraiment improbable : le fait que l'événement se soit produit est donc un indice substantiel que l'hypothèse nulle est en fait fautive » ! Pour prendre un exemple concret, imaginons que vous vous demandiez si un dé donné est équilibré ou s'il a tendance à privilégier les grandes valeurs ; et que pour trancher, vous lanciez le dé une fois, lancer qui vous donne un '6' (ce qui correspond à une p -valeur de 1/6 pour l'hypothèse nulle que le dé n'est pas truqué, si on considère que comme statistique de test le résultat du dé, avec suspicion des grandes valeurs). Vous semblerait-il raisonnable d'y voir là un début de preuve que le dé est truqué ? C'est tout au plus un indice *très, très faible* ! Et c'est bien ce que dit notre tableau, qui qualifie une p -valeur de 17 % de « très faiblement significative ». ♣

p -valeur	Interprétation
$20 \% \leq p_{\checkmark}$	Non significatif
$10 \% \leq p_{\checkmark} < 20 \%$	Très faiblement significatif
$5 \% \leq p_{\checkmark} < 10 \%$	Faiblement significatif
$2 \% \leq p_{\checkmark} < 5 \%$	Significatif
$1 \% \leq p_{\checkmark} < 2 \%$	Assez fortement significatif
$1 \% \leq p_{\checkmark} < 1 \%$	Fortement significatif
$10^{-6} \leq p_{\checkmark} < 1 \%$	Très fortement significatif
$p_{\checkmark} < 10^{-6}$	Preuve considérable comme “définitive”

TABLE 13.1 – Interprétations conventionnelles de la p -valeur en termes de significativité

13.4 Un exemple de test par p -valeur

Dans cette section, nous allons développer un exemple montrant comment un test par p -valeur se construit et se calcule en pratique. Nous allons nous placer, à nouveau, dans le cas du pédagogue ; mais contrairement au chapitre précédent, nous prenons cette fois-ci $\Theta = \mathbb{R} \times \mathbb{R}_+^*$ tout entier : autrement dit, on envisage ici que la variance des notes des élèves puisse éventuellement changer suite à l'introduction de la nouvelle méthode. L'hypothèse nulle que l'on souhaite tester est l'affirmation de l'enseignant que la nouvelle méthode va faire progresser ses élèves d'*au moins* un point en moyenne, autrement dit c'est l'hypothèse $\{\mu \geq \mu_{\min}\}$, avec $\mu_{\min} := 76$.

Pour une telle situation, en consultant la littérature^[†]^[‡], on trouve que la statistique de test appropriée est la suivante, avec test à gauche (autrement dit, plus on ira vers $-\infty$ et plus l'hypothèse nulle sera suspecte) :

$$T := \frac{M - \mu_{\min}}{S}, \quad (\text{EM}')$$

où nous rappelons que nous avons convenu de noter resp. par M et S les moyenne et écart-type empiriques des observations :

$$M := n^{-1} \sum_{i=0}^{n-1} X_i; \quad S := \left(n^{-1} \sum_{i=0}^{n-1} (X_i - M)^2 \right)^{1/2}.$$

Ayant pris connaissance de la réalisation t_{\checkmark} de notre statistique de test, nous souhaitons déterminer la p -valeur associée à celle-ci. Puisque notre critère de suspicion est à gauche, et au vu de l'hypothèse nulle que nous souhaitons tester, ladite p -valeur est définie comme

$$\sup_{\substack{\mu \geq \mu_{\min} \\ \sigma \in \mathbb{R}_+^*}} \mathbb{P} \left(\frac{M - \mu_{\min}}{S} \leq t_{\checkmark} \mid \theta = (\mu, \sigma) \right). \quad (\text{EN}')$$

Pour calculer cela, il nous faudrait connaître, pour $(\mu, \sigma) \in [\mu_{\min}[\times \mathbb{R}_+^*$, la loi de la variable aléatoire $(M - \mu_{\min}) / S$ sous le contexte $\mathbb{P}_{\mu, \sigma}$. De manière générale, cette loi est susceptible dépendre de la valeur de (μ, σ) au sein de l'hypothèse nulle ; néanmoins, nous espérons être dans un cas où un des principes (CK') ou (CN')

[†]. Concernant la façon d'exploiter la littérature, confier la § 13.5 ci-après.

[‡]. En l'occurrence, il s'avère d'ailleurs que cette statistique des test est aussi celle suggérée par le rapport des vraisemblances ! ☺

s'applique. En l'occurrence, vu que notre hypothèse nulle fait intervenir une inégalité, nous serions le cas échéant dans le cadre du principe (CN') ; et le "bord" de notre hypothèse nulle correspondrait aux cas où $\mu = \mu_{\min}$.

Est-il donc vrai que, sur ce bord, la loi de $(M - \mu_{\min})/S$ ne dépend pas de σ ? De fait, oui ! C'est même assez facile à démontrer : en effet, lorsque μ vaut μ_{\min} , alors sous le contexte $\mathbb{P}_{\mu_{\min}, \sigma}$, en introduisant les variables "réduites" $\check{X}_i := (X_i - \mu_{\min})/\sigma$, on a que

$$\frac{M - \mu_{\min}}{S} \stackrel{\text{déf}}{=} \frac{n^{-1} \sum_{i=1}^n (\mu_{\min} + \sigma \check{X}_i) - \mu_{\min}}{\text{var}_{\text{emp}}^{1/2}((\mu_{\min} + \sigma \check{X}_i)_i)} = \frac{n^{-1} \sigma \sum_{i=1}^n \check{X}_i}{\sigma \text{var}_{\text{emp}}^{1/2}((\check{X}_i)_i)} = \frac{\sum_{i=1}^n \check{X}_i}{n \text{var}_{\text{emp}}^{1/2}((\check{X}_i)_i)}, \quad (\text{EO}')$$

dont la loi ne dépend donc que de la loi jointe des \check{X}_i ; or, peu importe la valeur de σ , les variables réduites sont, sous $\mathbb{P}_{\mu_{\min}, \sigma(\bullet)}$, indépendantes [en tant que fonctions de v.a. indépendantes] et chacune de loi normale standard [d'après les propriétés d'invariance des lois normales], de sorte que leur loi jointe est toujours la même.

En fait, dans les situations pratiques, le fait que la statistique de test ait toujours la même loi dès lors que μ vaut μ_{\min} nous serait donné par la littérature E . D'autre part, la littérature nous dirait aussi — ce qui est très important ! — *quelle* est cette loi que suit la statistique de test. En l'occurrence, il s'agit du résultat suivant :

Théorème (EP'). *Dans le modèle du pédagogue $[\text{S}]$, lorsque le paramètre caché est de la forme (μ_{\min}, σ) , la loi de la statistique de test introduite ci-dessus vaut*

$$n^{-1/2} \times T_{\text{St}}(n-1), \quad (\text{EQ}')$$

où " $T_{\text{St}}(n-1)$ " est ce qu'on appelle la loi t de Student à $(n-1)$ degrés de liberté, dont les tables sont implémentées dans tous les bons logiciels utilisables pour traiter des données statistiques. \diamond

Grâce à ce théorème, nous avons que, lorsque (μ, σ) est au bord de notre hypothèse nulle, la probabilité que la statistique de test tombe en-dessous de $t_{\check{v}}$ est toujours la même, et vaut

$$\mathbb{P}(n^{-1/2} \times T_{\text{St}}(n-1) \leq t_{\check{v}}). \quad (\text{ER}')$$

Cependant, dans le supremum (EN') définissant notre p -valeur, il faut considérer non seulement les cas où $\mu = \mu_{\min}$, mais *aussi* ceux où $\mu > \mu_{\min}$... ! Pour ces cas-là, en vertu du principe (CN'), notre espoir est que la statistique de test ait tendance à prendre des valeurs *plus à droite* que dans le cas du "bord", de sorte que sa probabilité d'être $\leq t_{\check{v}}$ sera plus petit que pour les cas du bord, et n'interviendra donc pas dans la détermination du supremum. Cet espoir semble parfaitement raisonnable, puisqu'on comprend bien que T , qui a tendance à être d'autant plus grande que les X_i sont élevés (à cause du terme en M au numérateur), devrait donc tendre à prendre de plus grandes valeurs sous $\mathbb{P}_{\mu, \sigma(\bullet)}$ pour $\mu > \mu_{\min}$ que sous $\mathbb{P}_{\mu_{\min}, \sigma}$.

Est-il possible de rendre rigoureux cet argument intuitif ? Oui ; par exemple en considérant à nouveau les variables aléatoires "réduites" de façon à suivre la loi Normale(0, 1). En l'occurrence, notre paramètre (μ, σ) étant fixé avec $\mu > \mu_{\min}$,

$[\text{S}]$. Le « modèle du pédagogue » est bien connu de la littérature — quoique pas sous ce nom, bien entendu ! \check{v} — : il s'agit en effet du modèle d'échantillonnage associé à des lois normales de paramètres inconnus, qui apparaît naturellement dans nombre de situations pratiques E .

nous définissons les \check{X}_i comme $(X_i - \mu) / \sigma$: on a alors, comme précédemment, que les \check{X}_i i.i.d. normaux standards ; tandis que, en ce qui concerne la façon de ré-écrire la statistique de test à partir des variables réduites, on voit apparaître un terme supplémentaire par rapport à (EO')

$$T = \frac{\sum_{i=1}^n \check{X}_i}{n \operatorname{var}_{\text{emp}}^{1/2}((\check{X}_i)_i)} + \frac{\mu - \mu_{\min}}{\sigma \operatorname{var}_{\text{emp}}^{1/2}((\check{X}_i)_i)}. \quad (\text{ES}')$$

Or, la variable aléatoire que constitue ce terme supplémentaire — et que nous noterons ‘ Ξ ’ dans les lignes qui suivent — ne prend que des valeurs positives ! On a donc donné un sens rigoureux à l’idée que la loi de T , sous $\mathbb{P}_{\mu, \sigma}(\bullet)$, “prend des valeurs plus grandes que” $(\sum_{i=1}^n \check{X}_i) / n \operatorname{var}_{\text{emp}}^{1/2}((\check{X}_i)_i)$; et en particulier, on est maintenant en mesure de démontrer qu’elle a moins de risque de tomber en-dessous de $t_{\check{V}}$, par le raisonnement suivant :

$$\begin{aligned} \mathbb{P}_{\mu, \sigma}(T \leq t_{\check{V}}) &= \mathbb{P}_{\mu, \sigma}\left(\frac{\sum_{i=1}^n \check{X}_i}{n \operatorname{var}_{\text{emp}}^{1/2}((\check{X}_i)_i)} + \underbrace{\leq t_{\check{V}} - \Xi}_{\leq t_{\check{V}} \text{ p.-s.}}\right) \\ &\leq \mathbb{P}_{\mu, \sigma}\left(\underbrace{\frac{\sum_{i=1}^n \check{X}_i}{n \operatorname{var}_{\text{emp}}^{1/2}((\check{X}_i)_i)}}_{\text{de loi } n^{-1/2} \times T_{\text{St}}(n-1)} \leq t_{\check{V}}\right) = \mathbb{P}(n^{-1/2} \times T_{\text{St}}(n-1) \leq t_{\check{V}}). \quad (\text{ET}')$$

Par conséquent, dans le supremum défini par (EN'), les cas correspondant au bord de l’hypothèse nulle valent tous $\operatorname{répart} T_{\text{St}}(n-1; n^{1/2}t_{\check{V}}+)$, tandis que les autres cas ne sont jamais plus grand que cette valeur : *in fine*, on a donc trouvé que (la réalisation de) notre p -valeur vaut

$$\operatorname{répart} T_{\text{St}}(n-1; n^{1/2}t_{\check{V}}+) \quad \ddot{\circ} \quad (\text{EU}')$$

Remarque (EV'). Dans le calcul ci-dessus, nous avons *démontré* rigoureusement que seuls les cas au bord de l’hypothèse nulle nécessitaient d’être pris en compte pour le calcul du supremum. Néanmoins, comme nous l’avons signalé, c’était en réalité très intuitif... Dans un contexte d’exercice ou d’examen, vous auriez vous contenter [¶], dans un tel cas, d’une justification intuitive, du type : « Il est évident, au vu de sa forme, que la statistique de test prend (à σ fixé) des valeurs d’autant plus grandes que μ est grand : en effet, quand μ augmente, le numérateur tend à augmenter, tandis que les valeurs prises par le dénominateur restent les mêmes. Dès lors, pour déterminer le supremum définissant la p -valeur, il suffit de considérer les cas où μ est minimal (à σ fixé), autrement dit, ceux où $\mu = \mu_{\min}$ ». ♣

Application numérique. La formule pour la p -valeur que nous avons donnée s’implémente ainsi sous R :

```
> mumin = 76
> notes = c(91.7, 69.6, 81.5, 45.0, 69.7, 77.3, 42.1, 70.7, 66.9, 93.6,
+ 75.7, 93.7, 59.8, 75.2, 49.8, 88.2, 63.7, 90.0, 54.2, 57.2, 72.8, 58.3)
> (n = length(notes))
[1] 22
> (m = mean(notes)) # Moyenne empirique
```

[¶]. À moins, bien sûr, qu’il ne fût explicitement demandé une démonstration formelle! ☹

```
[1] 70.30455
> (s = sqrt(mean(notes ^ 2) - m ^ 2)) # Écart-type empirique
[1] 15.20698
> pt(sqrt(n) * (m - mumin) / s, n - 1)
[1] 0.04677142
```

Une telle p -valeur de 4,7 % est généralement considérée comme significative (quoique pas *très* significative non plus) : c'est un indice fort pour rejeter l'hypothèse nulle, même si gros coup de manchanse reste vaguement envisageable sans faire preuve de trop de mauvaise foi. Mais si l'enseignant est raisonnable, il devrait reconnaître qu'il s'est trompé : sa méthode, manifestement, échoue à faire progresser les élèves autant qu'il ne l'aurait espéré ! (voire, les fait carrément régresser).

Il peut être intéressant de regarder ce qu'aurait donné le résultat de notre test pour d'autres valeurs de l'observation. Par exemple, si toutes les notes avaient été deux points plus élevées (avec donc une moyenne empirique de 72,30, l'écart-type empirique restant le même), la p -valeur serait montée à 13,4 %, ce qui pour le coup n'est que très peu significatif : cela montre donc qu'il est tout à fait possible qu'on ait une moyenne de promotion plus de deux points et demi *en dessous* du seuil de référence tout en jugeant largement envisageable l'hypothèse que, *en réalité*, la méthode passe, quand on moyenne l'effet du hasard, *progresser* les élèves d'au moins un point !

Imaginons aussi le cas où les notes des élèves auraient eu la même moyenne que dans la réalité (soit 70,30), mais avec une dispersion *plus grande* : mettons, par exemple, qu'on aurait dans ce scénario $s_{\checkmark} = 21$. Dans ce cas aussi, la p -valeur serait sensiblement plus élevée (à 10,9 %) : des notes plus hétérogènes pour la promotion tendant en effet à indiquer que la nouvelle méthode présente manifestement une valeur assez élevée de σ , ce qui crée un fort "bruit" sur les notes faisant de la moyenne de la première promotion une moins bonne approximation de la véritable valeur μ_{\checkmark} ... À l'inverse, une dispersion plus petite, disons $s_{\checkmark} = 11$, (toujours avec la même moyenne) ferait plonger la p -valeur à 1,3 % : de tels résultats permettraient alors de réfuter l'hypothèse de l'enseignant avec une confiance très élevée !

13.5 Utilisation pratique d'un test par l'ingénieur

Ci-dessus, le traitement des statistiques de test a été mené assez en détail. Nous avons notamment vu comment utiliser les rapports de vraisemblances pour déterminer des statistiques de test pertinentes, comment on pouvait introduire des variables réduites pour ramener la loi de la statistique de test (sous l'hypothèse nulle) à un objet mathématique bien spécifié, et comment démontrer, si nécessaire, que seuls certains cas particuliers de l'hypothèse nulle sont à prendre en compte. Voir comment tous ces éléments marchent est en effet important, me semble-t-il, pour bien comprendre la démarche du test.

Dans la pratique néanmoins, l'ingénieur généraliste (non statisticien) n'aura pas à ré-inventer la roue ! Concrètement, que fait-il face à un problème donné ? Eh bien, il ouvre son ouvrage de statistique favori pour voir s'il y a un test bien adapté à cette situation ! On tombe alors sur une description dans le style suivant :

Tester la moyenne d'une population normale de variance inconnue

Lorsque X_1, \dots, X_n sont des v.a. normales i.i.d. Normale($\mu_{\mathcal{J}}, \sigma_{\mathcal{J}}$) avec $\mu_{\mathcal{J}}$ et $\sigma_{\mathcal{J}}$ inconnus, pour tester l'hypothèse $\{\mu_{\mathcal{J}} = \mu_{\text{réf}}\}$ (ou $\{\mu_{\mathcal{J}} \leq \mu_{\text{réf}}\}$, ou $\{\mu_{\mathcal{J}} \geq \mu_{\text{réf}}\}$), il faut utiliser la statistique de test

$$T := n^{1/2} \frac{\text{moy}(X_1, \dots, X_n) - \mu_{\text{réf}}}{\text{var}_{\text{emp}}^{1/2}(X_1, \dots, X_n)}, \quad (\text{EW}')$$

sachant que, lorsque $\mu_{\mathcal{J}} = \mu_{\text{réf}}$, cette variable aléatoire suit la loi $T_{\text{St}}(n-1)$.

Comment faut-il utiliser une description ? C'est ce que je détaille ci-dessous, la procédure étant ici illustrée pour le cas « tester la moyenne d'une population normale de variance inconnue », mais étant parfaitement générale :

! **Procédure (EX')**. Déjà, le texte de l'ouvrage de statistique nous fournit déjà la statistique de test et la loi de référence à laquelle il faut se ramener !

Il vous appartient néanmoins de vous répondre à la question « dois-je faire un test à gauche, à droite ou bilatéral ? ». Sur ce point, on pourrait bien sûr détailler le raisonnement mathématique ; mais le simple bon sens montre que, pour tester l'hypothèse $\{\mu \leq \mu_{\text{réf}}\}$, c'est un test à droite qu'il faut faire : en effet, quand l'hypothèse nulle est (largement) fautive, cela signifie que les X_i tendent à être (beaucoup) plus grands que $\mu_{\text{réf}}$, mais alors leur moyenne aura aussi tendance à être (beaucoup) plus grande, et donc la statistique de test aura tendance à être (très) positive : ce sont donc les plus petites grandes (au sens de « plus positives ») de la statistique de test qui seront les plus suspectes, alors qu'à l'inverse, de petites valeurs semblent indiquer que $\mu_{\mathcal{J}}$ est plus *petite* que $\mu_{\text{réf}}$, et ne doivent donc pas conduire au rejet de l'hypothèse nulle ! De même, si on avait eu à tester l'hypothèse nulle $\{\mu \geq \mu_{\text{réf}}\}$, il aurait convenu de faire un test à gauche, et pour tester l'hypothèse nulle $\{\mu = \mu_{\text{réf}}\}$, un test bilatéral.

Le point suivant consiste à voir comment seuls certains cas méritant d'être pris en compte. Ainsi, si l'hypothèse nulle que nous souhaitons tester est $\{\mu \leq \mu_{\text{réf}}\}$ (on aurait un argument similaire pour $\{\mu \geq \mu_{\text{réf}}\}$), vous allez argüer (ou démontrer si vous avez un doute) que, pour $\theta \in \Theta_0$, le cas où la valeur de $\mathbb{P}(T \geq t_{\mathcal{J}} \mid \theta = \theta)$ est le plus élevé est forcément lorsque $\mu = \mu_{\text{réf}}$, de sorte que vous pouvez vous limiter à ce cas.

Mais alors, dans le cas en question, T suit la loi $T_{\text{St}}(n-1)$, dont vous disposez de la fonction de répartition, ce qui vous permet de calculer la p -valeur : celle-ci correspond en effet à $\mathbb{P}(T_{\text{St}}(n-1) \geq t_{\mathcal{J}})$, autrement dit à $1 - \text{répart}T_{\text{St}}(n-1; t_{\mathcal{J}})$ (ici peu importe si on considère la fonction de répartition en $t_{\mathcal{J}}^-$ ou en $t_{\mathcal{J}}^+$, car les lois de Student sont diffuses), où la fonction $\text{répart}T_{\text{St}}$ sera implémentée dans votre logiciel de statistique préféré et $t_{\mathcal{J}}$ sera calculé à partir de la formule pour T en prenant les valeurs des observations effectives.

Dans le cas d'un test de l'hypothèse nulle $\{\mu \geq \mu_{\text{réf}}\}$, où le critère de suspicion est à gauche, la formule pour la p -valeur aurait de même été donnée par $\text{répart}T_{\text{St}}(n-1; t_{\mathcal{J}})$; et dans le cas d'un test de l'hypothèse nulle $\{\mu = \mu_{\text{réf}}\}$, où le critère de suspicion est bilatéral, la formule aurait été $2 \times ((1 - \text{répart}T_{\text{St}}(n-1; t_{\mathcal{J}})) \wedge \text{répart}T_{\text{St}}(n-1; t_{\mathcal{J}}))$. ♡

Chapitre 14

Intervalles de confiance et de prédiction fréquentistes

14.1 Intervalles de confiance

Définition et interprétation

Les tests d'hypothèses nous permettent d'obtenir de l'information sur une quantité d'intérêt $\varphi(\theta) =: \varphi$, dans la mesure où ils peuvent nous permettre de rejeter avec une certitude assez grande des hypothèses nulles sur la valeur de φ qui ne seraient pas réalistes. Cependant, comme on sait pas à l'avance quelles valeurs de φ sont réalistes ou non, cela implique que, avant d'avoir une idée honnête de ce que peut valoir $\varphi(\theta)$, il va falloir tester un grand nombre d'hypothèses nulles distinctes, ce qui n'est guère pratique... Nous, ce qu'on voudrait plutôt, c'est tout simplement connaître directement quelle est la plage des valeurs plausibles pour $\varphi(\theta)$!

C'est à ce problème que répondent les intervalles de confiance. Il s'agit en fait, du point de vue mathématique, de mener *simultanément* des tests pour toutes les valeurs envisageables de $\varphi(\theta)$; mais nous allons voir que les résultats se présentent sous une apparence fort différente.

Définition (EY') (Intervalle de confiance). Soit $\Theta \ni \theta \mapsto \text{Loi}(X \mid \theta = \theta)$ un modèle statistique, $\varphi(\theta)$ notre quantité d'intérêt; et soit $I(X)$ une statistique à valeurs dans les intervalles de \mathbb{R} [*]. On dit que $I(X)$ est un *intervalle de confiance au niveau α* pour $\varphi(\theta)$ lorsque, pour tout $\theta \in \Theta$, $\mathbb{P}(I(X) \not\supseteq \varphi(\theta) \mid \theta = \theta) \leq \alpha$; autrement dit, pour tout $\theta \in \Theta$, « $I(X) \not\supseteq \varphi(\theta)$ » est un test de niveau α de l'hypothèse nulle $\{\theta = \theta\}$. !

Définition (EZ') (Intervalles de confiance asymptotiques). Dans le cas d'un modèle statistique ayant un paramètre du modèle pour lequel on considère une asymptotique, on définit de même le concepts d'intervalle de confiance *asymptotiquement de niveau α* lorsque pour tout $\theta \in \Theta$, $\overline{\lim}_{n \rightarrow \infty} \mathbb{P}_\theta(I(X) \not\supseteq \varphi(\theta)) \leq \alpha$. !

Remarque (FA'). En fait, rien dans les définitions ci-dessus n'oblige conceptuellement $I(X)$ à être un intervalle : cela pourrait être n'importe quelle partie de

[*]. N'oubliez pas qu'une « statistique » peut être à valeurs dans *n'importe quel type* d'objets mathématiques : ce n'est pas forcément un nombre réel! (Nous avons d'ailleurs vu, dans le chapitre 11, le cas des lois empiriques, qui sont des statistique à valeurs « distribution de probabilité »...). En l'occurrence, $I(X)$ est à valeurs « intervalles » : c'est donc un *intervalle* défini à partir de la variable aléatoire X .

\mathbb{R} (et il faudrait alors plutôt parler de « zone de confiance »). D'ailleurs, de manière générale, $\varphi(\theta)$ pourrait être à valeurs dans n'importe quel ensemble \mathcal{G} , et alors $I(X)$ serait une partie de \mathcal{G} . Cependant, essentiellement toutes les situations pratiques concernent des intervalles de \mathbb{R} ; d'autant que, comme nous allons le voir ci-dessous, on dispose alors d'une procédure (partiellement) standardisée pour construire des intervalles de confiance. C'est pourquoi nous ne nous intéresserons qu'à ce cas ici. ♣

Remarque (FB'). Contrairement à l'intervalle de confiance bayésien (alias « intervalle de croyance »), qui correspond à un intervalle de fluctuation pour la loi à postériori et peut donc être défini “ x par x ” (puisque la loi à postériori a bien un sens x par x), un intervalle de confiance fréquentiste ne peut être défini que comme une *application* de \mathcal{X} dans l'ensemble *Ivl* des intervalles de \mathbb{R} . ♣

Remarque (FC'). En fait, on peut montrer que sous des hypothèses très peu restrictives, pour tout $x \in \mathcal{X}$, on peut trouver un intervalle de confiance $ICF: \mathcal{X} \rightarrow \text{Ivl}$ tel quel $ICF(x) = \emptyset$! Alors que dans le cadre bayésien (une fois, s'entend, qu'on a choisi notre priore), cela n'est pas possible, puisque la réalisation $ICB(x_{\mathcal{J}})$ d'un intervalle de croyance $ICB: \mathcal{X} \rightarrow \text{Ivl}$ doit nécessairement être un intervalle de fluctuation pour la loi à postériori de $\varphi(\theta)$.

Au passage, le fait qu'on puisse toujours trouver un intervalle de croyance fréquentiste ICF pour lequel $ICF(x) = \emptyset$ souligne un point particulièrement important en statistique fréquentiste (dont nous reparlerons dans la chapitre 17) : à savoir, qu'il faut décider de la procédure d'analyse statistique *avant* de connaître (ou du moins de regarder) les données : car en “bricolant” notre technique d'analyse au vu de nos données, on pourrait arriver à “conclure” à n'importe quel résultat qui nous arrange, aussi absurde soit-il...! ♣

!

Remarque (FD'). À cause de ce que nous venons de dire, il ne serait pas correct de conclure une analyse statistique par intervalle de confiance fréquentiste par une phrase comme « il y a une probabilité $(1 - \alpha)$ que $\varphi \in ICF(x_{\mathcal{J}})$ » : en fait, cette formulation correspondrait à l'interprétation d'un intervalle de confiance *bayésien*...! Tout ce que nous pouvons dire lors du calcul d'un intervalle de confiance fréquentiste, c'est qu'*avant* de lancer l'analyse statistique, donc *avant* de tirer $x_{\mathcal{J}}$, on savait qu'il y aurait une probabilité $(1 - \alpha)$ que notre intervalle de confiance aille contenir la véritable valeur de $\varphi_{\mathcal{J}}$. Mais peut-être qu'il y a certains x pour lesquels notre intervalle de confiance ne contiendra pas aussi souvent (voire jamais, dans le cas d'un intervalle vide!) la véritable valeur de φ , et d'autres pour lesquels au contraire $ICF(x)$ contiendra encore plus souvent (voire toujours, dans le cas d'un intervalle total!) la véritable valeur... Seule l'analyse bayésienne pourrait nous en dire plus long sur ce que vaut la probabilité de contenir la vraie valeur *une fois qu'on a observé $x_{\mathcal{J}}$* !

Pour cette raison, je vous demande de toujours suivre une des deux formulations suivantes (avec une préférence pour la première, plus rigoureuse) :

- « Notre procédure conduit à un intervalle de confiance au risque α pour φ dont la réalisation (pour nos données effectives) est $ICF(x_{\mathcal{J}})$ » ;
- « Au risque α , on conclut que $\varphi_{\mathcal{J}} \in ICF(x_{\mathcal{J}})$ ». (La locution « au risque α » ne veut à vrai dire pas dire grand-chose ; mais elle a le mérite d'être simple et de distinguer la phrase de conclusion pour un intervalle de confiance fréquentiste de celle pour un intervalle de croyance! ☺).

♣

Cette définition étant posée, il y a deux méthodes principales pour fabriquer des intervalles de confiance, que nous allons maintenant présenter.

Première méthode : Test avec paramètre

L'idée de cette méthode est de calculer le résultat du test de l'hypothèse nulle $\{\varphi = \varphi\}$ pour chaque valeur possible de φ , et à renvoyer comme intervalle de confiance l'ensemble des valeurs compatibles. Cela donne la procédure suivante, dans laquelle l'adjectif « asymptotique » correspond au cas d'un intervalle de confiance asymptotique :

Procédure (FE').

- 1° Pour tout $\varphi \in \mathbb{R}$, on considère un test (asymptotique) de niveau α $Test_\varphi$ de l'hypothèse $\{\varphi = \varphi\}$, ce qui nous donne donc une famille de tests $(Test_\varphi)_{\varphi \in \mathbb{R}}$: pour chacun de ces tests, on regarde quelle est la zone d'acceptation pour l'observation. On s'arrangera évidemment pour que tous ces tests s'obtiennent par le même calcul, où seule la valeur de φ changera. Il peut également être souhaitable que tous ces tests utilisent la même statistique de test T (auquel cas la loi de la statistique de test, et donc la région d'acceptation, dépendra bien entendu de la valeur de φ).
- 2° Pour une réalisation x_\vee de l'observation (ou pour une réalisation t_\vee de la statistique de test), la réalisation de l'intervalle de confiance au niveau α pour φ_\vee sera alors donnée par l'ensemble des $\varphi \in \mathbb{R}$ tels que le test $Test_\varphi$ soit négatif au vu de cette observation (autrement dit, tel que x_\vee (ou t_\vee) soit dans la région d'acceptation de $Test_\varphi$). ♡

Remarque (FF'). Une variante de cette méthode consiste à tester toutes les hypothèses élémentaires $\{\theta = \theta\}$ pour $\theta \in \Theta$, à regarder la "zone de confiance" correspondante dans Θ , puis à prendre l'image par φ de cette zone de confiance. Par exemple, si on note $T^{(\theta)}$ la statistique utilisée pour tester l'hypothèse $\{\theta = \theta\}$, et que tous nos tests reposent sur un critère de suspicion bilatéral, on pose

$$\Theta_{\text{conf}\vee} := \{\theta \in \Theta \mid t_\vee^{(\theta)} \text{ est dans l'intervalle de fluctuation au risque } \alpha \text{ de } \text{Loi}_\theta(T^{(\theta)})\},$$

(FG')

qui est l'ensemble des valeurs de θ jugées crédibles (au risque α) au vu de notre observation ; et la réalisation de notre intervalle de confiance est alors donnée par $\varphi(\Theta_{\text{conf}\vee})$. ♣

Je n'ai pas encore écrit d'exemple pour la procédure (FE') ; mais la mise en pratique est assez directe... Cette méthode a l'avantage d'être (théoriquement) applicable à essentiellement tout modèle, surtout une fois prise en compte la remarque (FF') ; néanmoins, les calculs peuvent se révéler particulièrement lourds, voire impraticables... D'où l'intérêt de la seconde méthode que nous allons à présent exposer, qui ne s'applique que dans certains cas, mais qui est beaucoup plus rapide.

Seconde méthode : Combiner intelligemment la quantité d'intérêt et l'observation

Nous allons d'abord expliquer le raisonnement conduisant à cette seconde méthode, avant d'expliquer, dans la procédure (FH'), comment on applique cette méthode en pratique.

L'idée est de partir d'une variable aléatoire $t(\varphi, X)$ ^[†] dépendant de l'observation et de la quantité d'intérêt φ , dont la loi est la même sous tous les contextes probabilistes $\mathbb{P}_{\theta(\bullet)}$ ^[‡] (ou, dans le cas d'un test asymptotique, dont la loi sous $\mathbb{P}_{\theta}^{(n)}$ tend, lorsque $n \rightarrow \infty$, vers une limite indépendante de θ). (Attention : rien ne garantit que trouver une telle v.a. sera facile, ni même possible... C'est pourquoi cette méthode est moins générale que la précédente. Néanmoins ici nous supposons disposer d'une telle variable, pertinente pour notre objectif ^[§]).

Ensuite, en remplaçant φ par φ_{\vee} dans l'expression de $t(\varphi, X)$, on obtient une statistique de test pour l'hypothèse $\{\varphi = \varphi\}$. Il faut en toute rigueur vérifier que cette statistique est intelligente, et en particulier que si φ_{\vee} n'est pas réellement égal à φ , on s'attend à ce que la statistique de test en question prenne des valeurs bien différentes de ce que serait sa loi sous l'hypothèse nulle. Normalement, la situation doit être la suivante : notre statistique de test (qui est à valeurs dans \mathbb{R} ou un des ses sous-ensembles) aura tendance à s'écarter dans une direction si $\varphi_{\vee} < \varphi$, et dans l'autre direction si $\varphi_{\vee} > \varphi$; de sorte que le critère de suspicion consistera à soupçonner les valeurs les plus extrêmes, les valeurs les moins suspectes étant celles les plus centrales. Par conséquent, l'intervalle d'acceptation pour la statistique de test sera toujours $[q_{\alpha/2}, q_{1-\alpha/2}]$, où $q_{\alpha/2}$ et $q_{1-\alpha/2}$ sont les quantiles aux niveaux respectifs $\alpha/2$ et $1 - \alpha/2$ de la loi censée être suivie par la statistique de test sous l'hypothèse nulle, autrement dit, de la loi $\text{Loi}(t(\varphi, X) \mid \varphi = \varphi)$, laquelle loi est connue ^[¶] : c'est en effet la loi $\text{Loi}_{\vee}(t(\varphi, X))$, vu que celle-ci est indépendante de la valeur de θ_{\vee} .

On regarde alors, pour une observation donnée, quel est l'ensemble des valeurs φ pour lesquelles la statistique construite ci-dessus aurait donné un test négatif : cet ensemble de valeurs est la réalisation de notre intervalle de confiance pour la quantité d'intérêt φ .

La procédure que nous venons de décrire, où nous faisons intervenir explicitement des tests bilatéraux, est en fait inutilement compliquée : nous ne l'avons donnée qu'à titre d'éclaircissement pédagogique. En réalité, le résultat de cette procédure est rigoureusement identique à celui de la procédure suivante, bien plus simple :

! Procédure (FH').

1° On part d'une variable aléatoire $t(\varphi, X)$ bien choisie ^[||], dépendant de l'observation et de la quantité d'intérêt φ , dont la loi (sous \mathbb{P}_{\vee}) est connue (ou, dans le

[†]. Attention : Même si cette variable aléatoire joue un peu le même rôle qu'une statistique de test, on ne peut pas la qualifier de « statistique », puisqu'elle dépend de φ , qui n'est pas une fonction de l'observation...!

[‡]. En pratique, plutôt que de dire qu'on a la même loi sous tous les contextes probabilistes $\mathbb{P}_{\theta(\bullet)}$, je dirai fréquemment « la loi de $t(\varphi, X)$ est connue, et ce, indépendamment de la valeur θ_{\vee} ».

[§]. Dans l'absolu, trouver une variable aléatoire vérifiant la propriété requise pourrait toujours être fait en prenant, par exemple, $t(\varphi, X) \equiv 0$: cependant, ce ne serait pas un choix *pertinent*, parce que la fonction $t(\varphi, x)$ ne dépendrait pas *réellement* de φ , de sorte que la connaissance de sa valeur n'apporterait aucune information sur φ — sans compter qu'ici la loi de $t(\varphi, X)$ serait une loi de Dirac, pour laquelle on ne peut pas partitionner le support en des zones de confiance et de rejet non triviales...

[¶]. Comme nous sommes dans le cadre fréquentiste, il n'est *prima facie* pas clair qu'on sache déterminer la loi conditionnée par $\{\varphi = \varphi\}$: c'est uniquement à cause des propriétés particulières de la v.a. $t(\varphi, X)$ qu'un tel conditionnement peut être calculé ici.

[||]. Dans les exercices et problème qui vous seront soumis, cette variable aléatoire sera donnée directement par l'énoncé, ou du moins fortement suggérée.

cas d'un test asymptotique, tend vers une loi connue lorsque $n \rightarrow \infty$), indépendamment de la valeur de $\theta_{\mathcal{J}}$.

- 2° On détermine l'intervalle de fluctuation IF pour la loi de la v.a. $t(\varphi, X)$. Ici la notion d'« intervalle de fluctuation » a bien du sens, même dans le cadre fréquentiste, car nous avons supposé avoir construit notre v.a. de sorte que $\text{Loi}_{\theta}(t(\varphi, X))$ ne dépende pas de θ . Cet intervalle de fluctuation correspond, informellement, aux valeurs « crédibles » pour $t(\varphi, X)$.
- 3° La réalisation de notre intervalle de confiance correspond alors aux valeurs de φ pour lesquelles $t(\varphi, x_{\mathcal{J}}) \in IF$.

♥

Remarque (FI'). Ci-dessus, j'ai supposé que la loi de $t(\varphi, X)$ devait être *exactement* la même sous toutes les lois \mathbb{P}_{θ} . En réalité, ce n'est pas vraiment nécessaire, et tout ce qui compte est que cette loi ne varie « pas trop », au sens où on peut trouver un intervalle IF qui, pour tout $\theta \in \Theta$, soit un intervalle de fluctuation au risque α de $\text{Loi}_{\theta}(t(\varphi(\theta), X))$: concrètement, il suffit pour ce faire de considérer tous les intervalles de fluctuation de toutes les $\text{Loi}_{\theta}(t(\varphi(\theta), X))$, puis de prendre IF comme étant leur réunion. Cependant, je ne connais pas vraiment de cas où cette construction donne lieu à quelque chose d'intéressant bien que les lois $\text{Loi}_{\theta}(t(\varphi(\theta), X))$ soient distinctes : c'est pourquoi j'ai seulement mis en valeur le cas où on a une identité parfaite entre toutes ces lois. ♣

Exemples de construction d'intervalles de confiance

Nous allons maintenant donner deux exemples de construction d'intervalles de confiance. Tous les deux s'appuient sur la seconde méthode présentée ci-dessus, mais dans les deux cas nous soulignons le lien avec l'approche par test avec paramètre plutôt que de recourir directement à la procédure (FH'). Le premier exemple est non asymptotique, tandis que le second est asymptotique.

Exemple (FJ'). Considérons le modèle du pédagogue (p. ??), pour lequel nous cherchons un intervalle de confiance pour μ . Il se trouve que dans ce modèle, on peut montrer que, quelle que soit la valeur de $\theta_{\mathcal{J}}$, la variable aléatoire

$$\frac{n^{-1} \sum_{i=1}^n X_i - \mu_{\mathcal{J}}}{\text{var}_{\mathbb{B}}^{1/2}(X_1, \dots, X_n) / n^{1/2}} \quad (\text{FK}')$$

suit (sous $\mathbb{P}_{\mathcal{J}}$) une loi qu'on appelle « loi de Student à $(n - 1)$ degrés de liberté », notée $T_{\text{St}}(n - 1)$. Nous ne chercherons pas à démontrer ce résultat ici : ce que nous allons faire, en revanche, est de montrer comment on peut l'utiliser pour obtenir un intervalle de confiance sur la valeur de μ . (Plus généralement, dans ce cours, le but de vos exercices sera de savoir utiliser des résultats mathématiques comme le précédent pour construire des intervalles de confiance, pas de démontrer les résultats en question).

Pour commencer, que cherchons-nous ? Réponse : à savoir quelles valeurs de μ sont plausibles ou pas. En quoi la propriété ci-dessus nous y aide-t-elle ? Réponse : parce qu'elle spécifie la loi que doit suivre une quantité qui ne dépend que de μ et de l'observation. Par conséquent, pour $\mu \in \mathbb{R}$, si la quantité

$$\frac{n^{-1} \sum_{i=1}^n X_i - \mu}{\text{var}_{\mathbb{B}}^{1/2}(X_1, \dots, X_n) / n^{1/2}} \quad (\text{FL}')$$

a une valeur particulièrement peu plausible pour une loi $T_{\text{St}}(n-1)$, ce sera vraisemblablement que μ n'était pas la vraie valeur de μ_{\checkmark} ! Maintenant, il s'agit de dire quelles sont les valeurs pour la quantité ci-dessus qui nous laissent le plus soupçonner que le véritable μ_{\checkmark} n'est pas égal à μ . Si par exemple μ est plus petit que la véritable valeur μ_{\checkmark} , on s'attend à ce que les valeurs trouvées pour notre statistique de test dévient vers $+\infty$; tandis que si μ est plus grand que μ_{\checkmark} , on s'attend à ce qu'elles dévient vers $-\infty$. Par conséquent, nous allons considérer que les valeurs les plus suspectes pour notre statistique de test sont celles qui sont les plus "extrêmes", et donc que les valeurs les moins suspectes sont celles les plus "centrales". (Comme expliqué plus haut, normalement lorsque notre but est de construire un intervalle de confiance, on doit toujours se retrouver dans ce cas).

Sous l'hypothèse que μ_{\checkmark} est effectivement égal à μ , la statistique de test doit suivre la loi $T_{\text{St}}(n-1)$, en vertu de ce que nous avons dit ci-dessus. Nous voulons une procédure qui rejette alors l'hypothèse que $\{\mu = \mu\}$ dans une proportion au plus α des cas, et par ailleurs nous avons dit qu'il nous paraissait intelligent de rejeter les valeurs trop petites ou trop grandes : on va par conséquent se placer dans une situation de test *bilatéral*. Soient donc $q_{\alpha/2}^{(n-1)}$ et $q_{1-\alpha/2}^{(n-1)}$ les quantiles aux niveaux resp. $\alpha/2$ et $(1-\alpha/2)$ de la loi $T_{\text{St}}(n-1)$: nous décidons de rejeter l'hypothèse nulle lorsque la statistique de test est $< q_{\alpha/2}$ ou $> q_{1-\alpha/2}$. Maintenant, la condition sous laquelle nous acceptons l'hypothèse nulle, à savoir

$$q_{\alpha/2}^{(n-1)} \leq \frac{n^{-1} \sum_{i=1}^n X_i - \mu}{\text{var}_B^{1/2}(X_1, \dots, X_n) / n^{1/2}} \leq q_{1-\alpha/2}^{(n-1)} \quad (\text{FM}')$$

peut se réécrire

$$\mu \in \left[n^{-1} \sum_{i=1}^n X_i - q_{1-\alpha/2}^{(n-1)} \text{var}_B^{1/2}(X_1, \dots, X_n) / n^{1/2}, \right. \\ \left. n^{-1} \sum_{i=1}^n X_i - q_{\alpha/2}^{(n-1)} \text{var}_B^{1/2}(X_1, \dots, X_n) / n^{1/2} \right] :$$

ça y est, nous l'avons, notre intervalle de confiance (au niveau α), c'est-à-dire l'ensemble des valeurs de μ pour lesquelles nous acceptons l'hypothèse que $\{\mu = \mu\}$ (au risque α) ! ♣

Exemple (FN'). Voyons maintenant un exemple plus compliqué. Dans cet exemple, on considère le modèle suivant de sondage électoral (idéalisé). Il y a un très grand nombre d'électeurs, tellement grand qu'on peut le considérer comme infini, dont une proportion $\pi_{A\checkmark}$ souhaite voter pour la candidate A , une proportion $\pi_{B\checkmark}$ pour le candidat B , et une proportion $\pi_{C\checkmark}$ pour la candidate C (donc dans ce modèle, l'espace du paramètre caché est donc $\Theta = \{(\pi_A, \pi_B, \pi_C) \mid \pi_A + \pi_B + \pi_C = 1\}$; et on note $(\pi_A, \pi_B, \pi_C) =: \theta$). On sonde n électeurs, qu'on suppose tirés uniformément et indépendamment dans la population; parmi ceux-là, on obtient X_A électeurs qui souhaitent voter pour A , X_B électeurs qui souhaitent voter pour B et X_C électeurs qui souhaitent voter pour C (l'espace de l'observation est donc $\mathcal{X}^{(n)} = \{(x_A, x_B, x_C) \in \mathbb{N}^3 \mid x_A + x_B + x_C = n\}$ — notez que notre modèle possède n comme paramètre du modèle).

Supposons que nous nous intéressons à la quantité $\varphi(\theta) := \pi_A - \pi_B$, pour laquelle nous cherchons à établir un intervalle de confiance. Notons $\hat{\pi}_A := X_A / n$, resp. $\hat{\pi}_B := X_B / n$, la proportion parmi les sondés d'électeurs souhaitant voter pour A ,

resp. pour B . Le théorème-limite central nous assure alors que, lorsque n tend vers l'infini,

$$\text{Loi}_{\checkmark} \left(\frac{(\hat{\pi}_A - \hat{\pi}_B) - (\pi_{A\checkmark} - \pi_{B\checkmark})}{n^{1/2}} \right) \rightarrow \text{Normale}(0, \pi_{A\checkmark}(1 - \pi_{A\checkmark}) + \pi_{B\checkmark}(1 - \pi_{B\checkmark}) + 2\pi_{A\checkmark}\pi_{B\checkmark}).$$

La variable aléatoire de notre membre de gauche a bien l'allure que nous souhaitons, car elle ne dépend que de l'observation et de la quantité d'intérêt $\varphi(\theta_{\checkmark})$. Par contre, le membre de droite ne va pas du tout, car il dépend du paramètre caché θ_{\checkmark} ! Pour pallier cela, on peut essayer de diviser par l'écart-type de la variable normale :

$$\text{Loi}_{\checkmark} \left(\frac{(\hat{\pi}_A - \hat{\pi}_B) - (\pi_{A\checkmark} - \pi_{B\checkmark})}{(\pi_{A\checkmark}(1 - \pi_{A\checkmark}) + \pi_{B\checkmark}(1 - \pi_{B\checkmark}) + 2\pi_{A\checkmark}\pi_{B\checkmark})^{1/2} n^{1/2}} \right) \xrightarrow{n \rightarrow \infty} \text{Normale}(0, 1). \quad (\text{FO}')$$

Cela règle le problème du membre de droite, mais alors nous avons à nouveau un souci avec le membre de gauche, car maintenant il ne dépend pas uniquement (outre les observations) de $(\pi_A - \pi_B)$ mais aussi de $(\pi_A(1 - \pi_A) + \pi_B(1 - \pi_B) + 2\pi_A\pi_B)^{1/2}$, qui n'est pas une fonction de $\varphi(\theta)$! L'idée est alors de remplacer le facteur $(\pi_A(1 - \pi_A) + \pi_B(1 - \pi_B) + 2\pi_A\pi_B)^{-1/2}$ par quelque chose qui en est très proche, mais qui ne dépend que des observations : typiquement, en remplaçant π_A et π_B par des *estimateurs* de ces quantités (cf. chap. ??) : le choix le plus naturel est, en l'occurrence, de remplacer π_A par $\hat{\pi}_A$ et π_B par $\hat{\pi}_B$. Pour notre modèle, on peut alors montrer que, dans ce cas, on a effectivement

$$\text{Loi}_{\checkmark} \left(\frac{(\hat{\pi}_A - \hat{\pi}_B) - (\pi_{A\checkmark} - \pi_{B\checkmark})}{(\hat{\pi}_A(1 - \hat{\pi}_A) + \hat{\pi}_B(1 - \hat{\pi}_B) + 2\hat{\pi}_A\hat{\pi}_B)^{1/2} n^{1/2}} \right) \xrightarrow{n \rightarrow \infty} \text{Normale}(0, 1). \quad (\text{FP}')$$

Maintenant, nous pouvons reprendre le raisonnement de l'exemple précédent. Considérons, pour $\delta \in \mathbb{R}$, la statistique de test

$$\frac{(\hat{\pi}_A - \hat{\pi}_B) - \delta}{(\hat{\pi}_A(1 - \hat{\pi}_A) + \hat{\pi}_B(1 - \hat{\pi}_B) + 2\hat{\pi}_A\hat{\pi}_B)^{1/2} n^{1/2}} \quad (\text{FQ}')$$

pour tester l'hypothèse $\{\pi_A - \pi_B = \delta\}$. Si l'hypothèse est vraie, cette statistique de test suivra asymptotiquement la loi Normale(0, 1). Si maintenant l'hypothèse est fautive, dans le cas où $\pi_{A\checkmark} - \pi_{B\checkmark} < \delta$, la statistique de test aura tendance à s'écartier vers $-\infty$, tandis que dans le cas où $\pi_{A\checkmark} - \pi_{B\checkmark} > \delta$, elle aura tendance à s'écartier vers $+\infty$. Nous allons donc considérer comme suspectes les valeurs les extrêmes prises par la loi Normale(0, 1) et comme plausibles ses valeurs les plus centrales. (Comme dit précédemment, ce critère de suspicion était en fait attendu dans le contexte de la construction d'un intervalle de confiance). Soit α le niveau (asymptotique) de l'intervalle de confiance que nous souhaitons mettre en place, et soient $q_{\alpha/2}$ et $q_{1-\alpha/2}$ les quantiles aux niveaux $\alpha/2$ et $(1 - \alpha/2)$ de la loi Normale(0, 1). Regarder si

$$\frac{(\hat{\pi}_A - \hat{\pi}_B) - \delta}{(\hat{\pi}_A(1 - \hat{\pi}_A) + \hat{\pi}_B(1 - \hat{\pi}_B) + 2\hat{\pi}_A\hat{\pi}_B)^{1/2} / n^{1/2}} \notin [q_{\alpha/2}, q_{1-\alpha/2}] \quad (\text{FR}')$$

sera alors bien un test asymptotiquement de niveau α de l'hypothèse nulle $\{\pi_A - \pi_B = \delta\}$. Et l'ensemble des valeurs pour lesquelles l'hypothèse nulle est acceptée est le

suisant : (où nous avons utilisé, pour simplifier la formule, qu'on a $q_{\alpha/2} = -q_{1-\alpha/2}$ par symétrie de la loi normale standard),

$$[\hat{\pi}_A - \hat{\pi}_B \pm q_{1-\alpha/2}(\hat{\pi}_A(1 - \hat{\pi}_A) + \hat{\pi}_B(1 - \hat{\pi}_B) + 2\hat{\pi}_A\hat{\pi}_B)^{1/2} / n^{1/2}] : \quad (\text{FS}')$$

cela est donc un intervalle de confiance asymptotiquement de niveau α pour la quantité $\pi_A - \pi_B$. \clubsuit

14.2 Intervalles de prédiction

Définition

! **Définition (FT')** (Intervalle de prédiction (fréquentiste)). Soit un modèle statistique de prédiction ayant pour espace du paramètre caché Θ et pour espaces des observations resp. passée et future \mathcal{X} et \mathcal{Y} , la loi d'une observation pour la valeur θ du paramètre étant $\text{Loi}_\theta(X, Y)$. Soit $g := \mathcal{Y} \rightarrow \mathbb{R}$ une application définissant une quantité d'intérêt réelle $g(Y)$. Notant Ivl l'ensemble des intervalles de \mathbb{R} , on dit alors qu'une application $IPF : \mathcal{X} \rightarrow Ivl$ définit un *intervalle de prédiction (fréquentiste) au risque α pour $g(Y)$* lorsque, pour tout $\theta \in \Theta$,

$$\mathbb{P}_\theta(IPF(X) \not\ni g(Y)) \leq \alpha. \quad (\text{FU}')$$

\heartsuit

Remarque (FV'). Bien que l'intervalle de prédiction fréquentiste porte sur une fonction de l'observation future (et non du paramètre caché), la condition sur le fait que l'intervalle contienne la quantité d'intérêt est quantifiée, elle, par les valeurs possibles du paramètre caché ! C'est parfaitement naturel, quand on y réfléchit bien : en statistique, tout le problème est qu'on ignore la valeur du paramètre caché, mais qu'on sait en revanche ce qui se passe à valeur donnée du paramètre caché ; par conséquent, pour faire des prédictions qui fonctionnent bien, il faut avoir quelque chose qui fonctionne bien quelle que soit la valeur du paramètre caché : et cela n'a rien à voir avec le fait que notre prédiction elle-même concerne le paramètre caché ou l'observation future ! \clubsuit

Remarque (FW'). Toutes les remarques sur la différence entre l'interprétation bayésienne et l'interprétation fréquentiste que nous avons faites pour les intervalles de confiance s'appliquent aussi, *mutatis mutandis*, aux intervalles de prédiction. \clubsuit

Construction

Lorsqu'il s'agit de construire un intervalle de prédiction fréquentiste, l'idée de faire une famille de tests paramétrée par la valeur testée pour la quantité d'intérêt ne fonctionne plus : en effet, par nature, un test d'hypothèse ne peut porter que sur le paramètre caché, et pas sur l'observation future... En revanche, la technique exposée en § 14.1, elle, peut être adaptée sans difficulté. Cela donne la procédure suivante :

! **Principe (FX')**.

1° On détermine une variable aléatoire

$$t(X, \bullet g(Y)) \quad (\text{FY}')$$

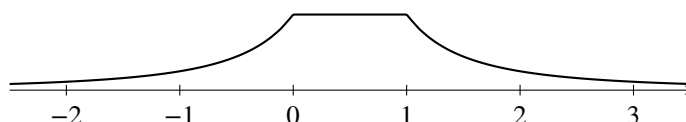


FIGURE 14.1 – Densité de la loi de $(Y_3 - \min(X_1, X_2)) / |X_2 - X_1|$ lorsque X_1, X_2, Y_3 sont tirées indépendamment selon une loi uniforme.

telle que $\text{Loi}_\theta(t(X, \bullet g(Y)))$ ne dépende pas de θ ^[**] Dans les exercices et problèmes qui vous seront soumis, cette variable aléatoire sera donnée directement par l'énoncé, ou du moins fortement suggérée.

2° On détermine ensuite l'intervalle de fluctuation J pour $\text{Loi}_\theta(t(X, \bullet g(Y)))$. (Rappelons que cette loi est connue indépendamment de θ). L'intervalle J est donc tel que

$$\forall \theta \in \Theta \quad \mathbb{P}_\theta(t(X, \bullet g(Y)) \notin J) \leq \alpha. \quad (\text{FZ}')$$

(De même que dans la remarque (FI'), on n'était en fait pas obligé que $\theta \mapsto \text{Loi}_\theta(t(X, \bullet g(Y)))$ fût constante : tout ce qui compte est d'avoir la propriété ci-dessus, qu'on peut obtenir par exemple en prenant pour J la réunion de tous les intervalles de fluctuation des $\text{Loi}_\theta(t(X, \bullet g(Y)))$. Mais en pratique, je ne connais pas de cas intéressant où la loi de $t(X, \bullet g(Y))$ soit variable...).

3° Un intervalle de prédiction pour $\bullet g(Y)$ est alors donné par

$$x \mapsto \{g \in \mathbb{R} \mid t(x, g) \in J\}. \quad (\text{GA}')$$

◇

[††].

Un exemple

Exemple (GB'). On considère un modèle où l'observation passée est un couple d'observations $X_1, X_2 \in \mathbb{R}$ et l'observation future est $Y_3 \in \mathbb{R}$, sachant que X_1, X_2, Y_3 sont i.i.d. selon une loi $\text{Unif}^{\text{me}}(\alpha_\nu, \beta_\nu)$ avec α_ν et β_ν inconnus ($\alpha_\nu < \beta_\nu$). On peut montrer que, quels que soient α_ν et β_ν , la loi (sous \mathbb{P}_ν) de $(Y_3 - \min(X_1, X_2)) / |X_2 - X_1|$ est toujours la même : et plus précisément, on a (voir la figure 14.1)

$$\mathbb{P}_\nu \left(\frac{Y_3 - \min(X_1, X_2)}{|X_2 - X_1|} \in dt \right) = \begin{cases} \frac{1}{3}(1-t)^{-2} \text{vol}(dt) & \text{pour } t \leq 0; \\ \frac{1}{3} \times \text{vol}(dt) & \text{pour } t \in]0, 1]; \\ \frac{1}{3}t^{-2} \text{vol}(dt) & \text{pour } t > 1. \end{cases} \quad (\text{GC}')$$

On calcule que l'intervalle de fluctuation à 90 % de cette loi est $[-3,48, 4,48]$: il y donc 90 % de chances que $(Y_3 - \min(X_1, X_2)) / |X_2 - X_1| \in [-3,48, 4,48]$, autrement dit que Y_3 tombe entre $\min(X_1, X_2) - 3,48 \times |X_2 - X_1|$ et $\min(X_1, X_2) + 4,48 \times |X_2 - X_1|$, qui sont donc les bornes de notre intervalle de prédiction. L'application numérique donne, pour $x_{1\nu} = 64$ et $x_{2\nu} = 61$, qu'on aura

$$Y_3 \in [50,58, 74,42] \quad (\text{GD}')$$

au risque 10 %.

♣

[††]. Attention : Ici il s'agit bien d'une dépendance en θ ! Confer remarque (FV').

Quatrième partie

Du bon usage de la statistique

Chapitre 15

Choix de la priore

15.1 Problématique

Principe général pour choisir la priore

Le choix de la priore est absolument crucial en analyse bayésienne : en partant de deux priores différentes, on arrive à deux postérieures différentes, et donc potentiellement à deux conclusions différentes pour notre analyse... Mais alors, comment faut-il choisir la priore ?! La réponse est paradoxale : il n'y a pas forcément *une* façon bien définie de choisir la priore... mais on ne peut pas la choisir arbitrairement pour autant !

Principe (GE'). *La priore a un sens bien précis : elle exprime la probabilité à priori que le paramètre caché vaille telle ou telle valeur. Mais comment donner un sens à la « probabilité » que, par exemple, la nouvelle méthode de l'enseignant donne un meilleur résultat moyen ?... Nous avons en effet l'habitude de dire « la probabilité d'un évènement, c'est la proportion de fois où cet évènement se produirait si on répétait un très grand nombre de fois l'expérience à l'identique » : c'est ce qu'on appelle parfois la définition fréquentiste de la notion de probabilité. Or ici, il n'y a pas d'expérience répétable à l'identique : l'enseignant a eu l'idée d'une nouvelle méthode bien particulière, qui est peut-être efficace et peut-être pas, mais cela ne dépend en rien du hasard : simplement du fait que l'enseignant n'arrive pas à deviner à l'avance les valeurs de μ et σ ... On entre alors dans la vision épistémique de la notion de probabilité. Si, par exemple, l'enseignant dit « Il y a 30 % de probabilité que σ soit plus petit que $\sigma_{\text{réf}}/2$ », cela veut en fait dire qu'il serait exactement aussi surpris que s'apercevoir que $\{\sigma < \sigma_{\text{réf}}/2\}$ qu'il le serait de voir qu'une variable aléatoire uniformément tirée entre 0 et 1 soit inférieure à 0,30 : par exemple, si on lui proposait un pari rapportant 5 € (moins déduction de la mise) dans le cas où $\{\sigma < \sigma_{\text{réf}}/2\}$ (après expérimentation de la nouvelle méthode, s'entend, sur un nombre d'élèves tellement grand qu'on n'aurait plus de doute sur la vraie valeur de σ), il serait prêt à entrer dans un tel pari pour une mise de 1 €, mais pas pour une mise de 2 €^[*]. ◇*

!!

[*]. En supposant ici que, pour des montants inférieurs à 5 € en valeur absolue, l'utilité de « recevoir un gain net de x € » (une perte étant comptée comme un gain négatif) est pratiquement une fonction affine de x : cette approximation est très raisonnable lorsque les montants en jeu sont petits, car dans ce cas il n'y a guère d'aversion (ni de propension) au risque, conférer Remarque (HQ').

Probabilités épistémiques

La notion de probabilité épistémique suggère deux difficultés conceptuelles, que nous discutons dans cette sous-section.

Point (GF') (Subjectivité des probabilités épistémiques). Une première difficulté avec la notion de probabilités épistémiques est que ces quantités sont fondamentalement *subjective* : dans le cas du modèle du pédagogue, par exemple, une enseignante expérimentée arrivera certainement à deviner avec beaucoup plus de précision l'effet de la nouvelle méthode qu'un enseignant novice, et attribuera donc une loi différente à θ pour la *même* situation... Mais est-ce à dire que l'enseignante expérimentée a *raison* pour sa loi et que le novice a tort ? Non : le point est qu'une probabilité épistémique n'est définie que *par rapport à la personne qui pose cette probabilité*, et est donc correcte dès lors que la personne a bien pris en compte, de façon appropriée, toute l'information dont elle disposait ! ♣

Remarque (GG'). Il convient de noter une subtilité dans le point précédent : « prendre en compte de façon appropriée toute l'information dont on dispose » inclut, le cas échéant, l'information concernant le fait que d'autres personnes (ayant accès à des informations que nous n'avons pas) attribuent telle ou telle probabilité épistémique au phénomène concerné ! Ainsi, dans le cas de l'enseignant novice, cela signifie qu'il ne devra pas donner la même loi à priori à θ selon qu'il a eu, ou non, l'occasion de discuter avec sa collègue expérimentée : en cas de discussion, et pour peu qu'il accorde un minimum de foi à l'expertise de sa collègue, il devra modifier sa croyance à priori pour tenir compte de l'avis de celle-ci... [†] ♣

Point (GH'). La seconde difficulté est que, bien souvent, nous avons tendance à nous tromper sur notre évaluation des probabilités épistémiques ! Lequel (ou laquelle) d'entre vous n'a-t-il jamais dit à des amis, par exemple, « oui, je viendrai à votre soirée ; c'est sûr à 99,99 % », et pourtant a eu un contretemps de dernière

[†]. En mathématiques et en philosophie de la connaissance, la problématique de la prise en compte des probabilités épistémiques émises par des tiers donne lieu à des développements que je trouve absolument fascinants... Imaginons ainsi deux individus supposés parfaitement rationnels, qui ont recueilli chacun de son côté diverses informations concernant une affirmation donnée. S'ils n'ont pas communiqué entre eux, les probabilités épistémiques qu'ils auront concernant cette affirmation seront évidemment différentes : rien d'étonnant jusque-là. À l'inverse, si chacun va transmettre à l'autre la totalité des informations dont il dispose, les deux individus se retrouvent alors avec exactement les mêmes informations, et (puisque'ils sont parfaitement rationnels) vont en déduire la même probabilité épistémique. Mais bien entendu, communiquer la totalité des informations dont on dispose est en pratique impossible, tant le détail de ces informations est encodé de façon subtile dans notre cerveau... Mais c'est alors qu'un théorème fascinant intervient : supposons que le premier individu (appelons-le *A*) envoie au second (*B*) un message indiquant simplement la probabilité épistémique qu'il attribue au phénomène. *B* répond ensuite en informant *A* de la probabilité épistémique qu'il attribue au phénomène *après prise en compte du premier message de A*. Puis *A* réplique à son tour en indiquant sa probabilité épistémique *après prise en compte du message de B*, et ainsi de suite... Dans ce cas, on peut montrer que les probabilités épistémiques de *A* et *B* doivent converger (en les supposant parfaitement rationnels) vers une valeur commune au fil des échanges, et assez rapidement qui plus est ! Or, on observe en pratique que ce n'est pas ce qui se passe : des gens très intelligents et très rationnels peuvent camper sur des évaluations épistémiques fortement différentes concernant le même phénomène, même en ayant connaissance de leur discordance... En philosophie, on appelle cela la question du *désaccord entre pairs épistémiques*. Les philosophes essayent de comprendre si l'omniprésence de tels désaccords est l'indication que les humains sont irrémédiablement profondément irrationnels, ou s'il n'existe pas des mécanismes subtils qui font que le théorème mathématique sur la convergence des probabilités épistémiques ne s'applique pas dans la vraie vie... ?

minute...?! Pourtant, lorsque nous annonçons un tel niveau de certitude, l'annulation ne devrait survenir qu'une fois sur 10 000 : beaucoup plus que le nombre de soirées auxquelles vous êtes jamais allés... De manière générale, notre intuition est très mauvaise pour évaluer les probabilités extrêmement petites (ou, par complémentarité, extrêmement proches de 1) : par exemple, nous avons beaucoup de mal à admettre que la probabilité que deux personnes prises au hasard dans le monde soient nées le même jour à moins de 5 minutes d'intervalle est beaucoup plus grande que la probabilité qu'une grille de loto donnée remporte le jackpot ! (ce qui consiste à trouver 5 numéros corrects parmi 49, plus 1 parmi 10). Et que cette probabilité elle-même est beaucoup plus grande que la probabilité qu'une météorite qui tombe sur Terre se trouve atterrir à moins de 2 km de chez vous...

Or, autant la probabilité épistémique que nous pouvons attribuer à un événement dépend de l'*information* dont nous disposons, autant on ne peut pas lui donner *n'importe quelle* valeur selon notre convenance ! En effet, il faut dans tous les cas que, par exemple, les événements dont nous prédisons à priori que la probabilité qu'ils arrivent vaut 10 % se produisent effectivement environ 10 fois pour cent prédictions^[‡] (même s'il ne s'agit pas là à proprement parler de la répétition d'une expérience, puisque ce sont des événements complètement différents à chaque fois : c'est ce qu'on appelle, dans le langage technique, la *calibration* correcte des probabilités épistémiques^[§]. Il faut donc une méthode permettant de quantifier de façon raisonnablement précise la façon de convertir nos informations et nos croyances en probabilités... ♣

Nous allons présenter ci-après les façons appropriées pour choisir une priore. Comme il faut s'y attendre, la méthode à suivre dépendra du type d'informations dont nous disposons (ou pas) sur le paramètre caché. Nous allons commencer par présenter les cas où la loi du paramètre caché peut être déterminée sans ambiguïté, en allant progressivement vers des situations où l'information sur le paramètre caché sera décrite de façon de plus en plus vague.

Remarque (GI'). Dans le cadre de ce cours, on ne vous demandera normalement pas de choisir vous-mêmes la priore d'un modèle statistique ; cependant, il est important que vous connaissiez les règles pour le choix d'une priore, et que vous soyez

[‡]. Ici j'ai supposé implicitement que les différentes prédictions concernées portaient sur des événements pouvant être considérés comme indépendants les uns des autres. En effet, il va de soi que, si je m'intéresse aux cent événements « il pleuvra sur la place Stanislas le 22 septembre 2009 à 17 h 00 », « il pleuvra sur la place Stanislas le 22 septembre 2009 à 17 h 01 », ..., « il pleuvra sur la place Stanislas le 22 septembre 2009 à 18 h 39 », que j'attribue à chacun de ces cent événements une probabilité de 10 %, et que pourtant aucun d'entre eux ne se réalise, ce n'est pas le signe d'une mauvaise calibration de ma part : simplement que, pour des événements aussi corrélés entre eux, il n'y a aucune raison de s'attendre à observer un phénomène de type « loi des grands nombres » !...

[§]. Nous avons vu que deux individus rationnels peuvent très bien attribuer des probabilités épistémiques différentes au même événement, pour peu qu'ils s'appuient sur des informations différentes (et qu'ils n'aient pas pu en échanger entre eux) : ainsi, le fait de savoir qu'un individu donné attribue une probabilité épistémique de tant ou tant à tel ou tel phénomène ne peut jamais constituer une preuve que cet individu est irrationnel (à moins, bien sûr, de connaître précisément toutes les informations dont il dispose). En revanche, les probabilités épistémiques émises par un individu rationnel *doivent* être correctement calibrées : ainsi, si les pronostics de votre belle-sœur concernant les matchs de football de Ligue 1 sont tels que, quand elle dit qu'une équipe a au moins 90 % de chances de gagner, on constate en pratique, sur le long terme, que l'équipe désignée comme favorite ne gagne effectivement que dans 70 % des situations ainsi pronostiquées par votre belle-sœur, le problème n'est pas que votre belle-sœur manque d'information : cela prouve que, forcément, les évaluations qu'elle fait des probabilités épistémiques ne traitent (en général) pas les informations dont elle dispose selon un raisonnement correct !

capables de justifier le choix d'une priore donnée par l'énoncé, voire de proposer une priore en vertu de certains critères spécifiés. ♣

15.2 Cas où la priore est non ambiguë

Il s'agit là du cas le plus facile : c'est lorsque la situation confrontée correspond *réellement* à tirer au sort la valeur θ_{\checkmark} de θ (quoique sans la révéler), puis à tirer x_{\checkmark} [¶] selon la loi $\text{Loi}(X \mid \theta = \theta_{\checkmark})!$

Exemple (GJ'). Plaçons-nous dans la situation du chasseur, mais avec une petite nuance dans l'interprétation du problème mathématique : ici, on suppose que les qualités du chasseur sont parfaitement connues ; mais par contre, qu'on ignore la qualité de l'*arme* avec laquelle il va tirer : notre chasseur possède en effet trois fusils indistinguables à l'œil nu, mais dont l'un a un léger défaut de fabrication qui fait chuter sa précision : si le chasseur tire avec un bon fusil, il aura un taux de réussite de 30 % ; tandis que s'il tire avec le fusil défectueux, il n'aura un taux de réussite que de 15 %... La question correspond ici à identifier l'état du fusil avec lequel le chasseur tire : comme c'est équivalent à déterminer son taux de réussite, mathématiquement parlant, nous sommes bien dans la situation de base.

Quelle est la priore correspondante ? La réponse ne souffre d'aucune ambiguïté : puisque les fusils sont indistinguables à l'œil nu, l'arme qu'utilise notre chasseur est prise au hasard uniformément parmi ses trois fusils ; et il y a donc deux chances sur trois d'avoir un bon fusil (auquel cas θ_{\checkmark} vaut 30 %), et une chance sur trois d'avoir le fusil défectueux (auquel cas θ_{\checkmark} vaut 15 %). Ainsi, il *faut* prendre

$$\text{Loi}_{\text{pr}}(\theta) = \frac{2}{3}\delta_{30\%} + \frac{1}{3}\delta_{15\%}. \quad (\text{GK}') \quad \clubsuit$$

Ce genre de situations où la priore est déterminée sans ambiguïté sont celles que vous avez dû voir dans vos exercices de lycée. Mais toute la force de la statistique bayésienne est, justement, que sa démarche peut aussi s'appliquer dans des cas où la notion de « probabilité à priori » est définie de façon beaucoup plus floue ! C'est ce que nous verrons dans les § 15.4 et 15.5. Mais auparavant, il importe de préciser une « méta-règle » qui doit s'appliquer *dans tous les cas* pour déterminer une priore, quelle que soit la partie de ce chapitre à laquelle on se réfère : c'est ce que nous allons voir dans la section suivante !

15.3 Postérieure d'hier, priore d'aujourd'hui

!! Principe (GL'). Si on a déjà des informations sur θ issues d'une expérience précédente, alors il faut inclure ces informations dans la loi à priori qu'on attribue à θ !

En particulier, dans le cas où l'« expérience précédente » en question est formalisée en termes de statistique bayésienne, la loi à priori pour θ qu'il convient de prendre pour l'expérience présente est celle que l'expérimentateur de l'expérience précédente (peu importe qu'il s'agisse de nous-même ou quelqu'un d'autre) aurait trouvée comme loi à postériori. ◇

[¶]. Ou éventuellement $(x_{\checkmark}, y_{\checkmark})$ si on est dans une situation de prévision, mais, le cas échéant, en ne révélant que x_{\checkmark} dans un premier temps.

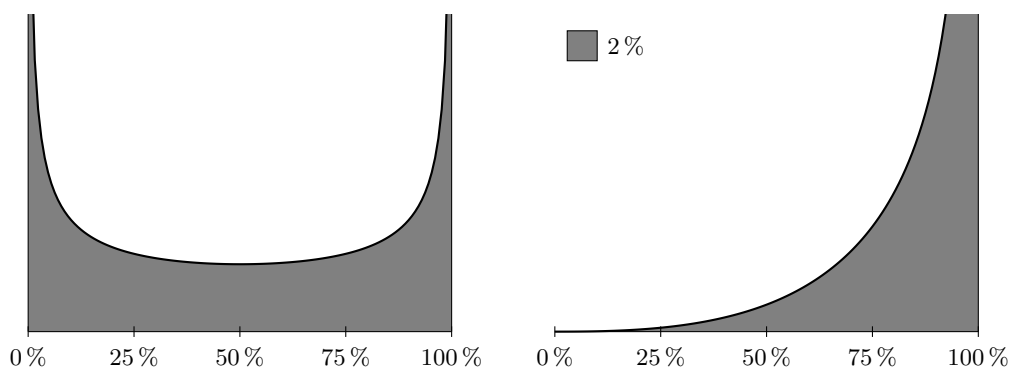


FIGURE 15.1 – Tracés, à la même échelle, de la loi arcsinus (à gauche) et de la loi Bêta($3\frac{1}{2}$, $1/2$) (à droite).

Exemple (GM'). Prenons à nouveau le cas du chasseur. Normalement, quand un nouveau candidat arrive, les statisticiens du club du Bouchonnois choisissent comme priore pour θ la loi arcsinus. Cependant, ici le cas est un peu particulier : car le candidat qui est examiné aujourd'hui a procédé à trois tirs d'essai hier (les conditions de ces tirs d'essais étaient les mêmes que pour le test), qu'il a réussis tous les trois ! Quelle serait la postérieure attribuée à θ pour un chasseur ayant réussi trois tirs ? Ce serait dans ce cas la loi Bêta($3\frac{1}{2}$, $1/2$), qui est très nettement différente de la loi arcsinus ! (voir figure 15.1). Et cela est parfaitement logique : puisqu'ici, au vu de ses trois tirs d'essais, avant même que notre candidat commence le test officiel, nous sommes déjà fortement convaincus qu'il a une bonne fiabilité...! ♣

Remarque (GN'). Dans l'exemple précédent, le résultat qu'on obtiendra, en partant de la priore Bêta($3\frac{1}{2}$, $1/2$), après analyse bayésienne des 25 tirs de test du chasseur, sera exactement le même que celui qu'on aurait obtenu en partant de la priore arcsinus, *mais en intégrant les trois premiers tirs à nos données !* (et donc en considérant que l'expérience a consisté à tirer sur 28 plateaux en tout). En effet, dans les deux cas, il s'agit du traitement bayésien des *mêmes* informations (à ceci près que dans le premier cas elles ont été traitées en deux étapes, alors que dans le second cas on a procédé directement à un traitement global) : il est donc intuitivement évident que les deux calculs doivent conduire au même résultat !

Il s'agit là, bien entendu, d'un phénomène parfaitement général : dire « on prend la postérieure du début comme priore pour la fin » ou dire « on considère que l'expérience englobe le début et la fin », ce ne sont que deux façons différentes, et mathématiquement parfaitement équivalentes, de prendre en compte toute l'information disponible ! ☺ ♣

Remarque (GO'). En sciences expérimentales, l'idée de choisir la priore au vu des observations précédentes est particulièrement importante. En effet, il est fréquent que plusieurs chercheurs réalisent des expériences différentes sur le même phénomène (par exemple, l'effet du tabagisme sur le cancer du poumon), correspondant donc au même paramètre caché (en l'occurrence, la façon dont fumer augmente — ou pas — le risque de cancer). Dès lors, il serait absurde et même malhonnête de faire table rase des expériences passées lorsqu'on mène une nouvelle étude : au contraire, notre croyance à priori sur le paramètre caché *doit* être la croyance à postériori consécutive aux expériences précédentes ! (sauf, bien sûr, si on a des raisons de penser que ces expériences précédentes ont été truquées ou mal conçues statistiquement). Imaginez ainsi qu'un élève vienne me voir et me montre des données épidémiologiques sur sa famille, menées à partir d'une priore non informative (cf. § 15.5 *infra*), dont il déduise,

au prix d'une analyse statistique impeccable, qu'il n'y a que 40 % de chances qu'il existe un lien entre tabagisme et cancer. Vais-je *réellement* croire qu'il y a 60 % de chances que fumer soit sans aucun risque? Certainement pas! Car, au vu des nombreuses études scientifiques déjà menées préalablement sur le sujet, j'avais du fait des résultats de ces études de très fortes raisons de croire à priori (c.-à-d. avant de voir les résultats sur la famille de mon élève) que fumer augmente le risque de cancer : le choix d'une priore non informative par mon élève n'était donc pas pertinent! La façon correcte de traiter les résultats d'une telle étude serait plutôt celle d'une *révision de croyances*, du type suivant : avant d'avoir vu l'enquête de mon élève, j'attribuais une probabilité de 99,995 % à l'existence d'un lien entre tabagisme et cancer ; après avoir vu son enquête, je n'y attribue plus "qu'"une probabilité de 99,994 %.

Curieusement, la démarche que je viens d'exposer n'est *pas* celle de la plupart des travaux de sciences expérimentales : les chercheurs dans ces matières ont au contraire tendance à utiliser l'approche *fréquentiste*, et à analyser les résultats des différentes expériences *indépendamment*. Peut-être est-ce parce que l'approche bayésienne est trop subtile à manipuler de façon convaincante sur de tels cas concrets, à moins que ce ne soit simplement parce que cette approche n'est pas celle qui a été enseignée aux chercheurs dans leur jeunesse... $\ddot{\smile}$ Seules les *méta-études* utilisent parfois cette approche bayésienne pour, justement, mettre à bout les résultats des différentes études individuelles : elles arrivent alors à montrer de façon quantifiable comment beaucoup de petites preuves peuvent conduire à une certitude très grande, ou quelle est la conclusion générale qu'il convient de tirer de résultats individuels en apparence contradictoires ; et ce, grâce à la méthode bayésienne en particulier \smile \clubsuit

15.4 Traduction mathématique d'une expertise préalable

Je parle d'« expertise » dans le cas où on n'a pas d'information de l'énoncé qui nous donne une connaissance *parfaite* de la priore, mais où on a cependant une idée *honnête* de ce que vaut cette priore au vu de notre expertise. En fait, « utiliser notre expertise pour choisir la priore » s'apparente à l'idée de « utiliser les expériences précédentes pour faire de la postérieure de ces expériences la priore de notre modèle », mais dans un cadre où on n'est pas forcément à même de poser un modèle précis, comme nous allons le voir sur les exemples suivants.

Exemple (GP'). Dans les exemples précédents traitant du problème du chasseur, nous avons proposé une priore de type neutre (voir § 15.5 *infra*) pour définir la loi à priori de θ sur $]0, 1[$, en l'occurrence la loi arcsinus. Mais en fait, pour peu que le club du Bouchonnois ait un peu d'expérience dans le recrutement des candidats, il peut faire bien mieux...! Par exemple, à force de voir des candidats défiler, puis d'avoir pu juger de leurs performances sur le plus long terme pour ceux qui ont été recrutés \llbracket , le club s'est certainement fait une idée honnête de la proportion de bons chasseurs (disons ici qu'un « bon chasseur » est un chasseur capable de toucher au moins un plateau sur quatre, autrement dit pour lequel $\{\theta \geq 1/4\}$) et de mauvais chasseurs. Ainsi, si sur les 20 candidats précédents ils s'est avéré y avoir 8 bons chasseurs (soit 40 % des candidats passés) et 12 mauvais chasseurs (soit 60 %), plutôt que de poser une priore arcsinus pour θ (θ représentant le niveau

\llbracket . Ici nous ferons l'hypothèse que le club ne rejette un candidat que s'il n'y a aucun doute raisonnable sur le fait que ce soit un mauvais chasseur. Dès lors, le club connaît pour chaque candidat son statut de bon ou de mauvais chasseur : soit que son test n'ait laissé aucun doute raisonnable, soit qu'il a été recruté et qu'on a pu jauger finement de son niveau réel par la suite.

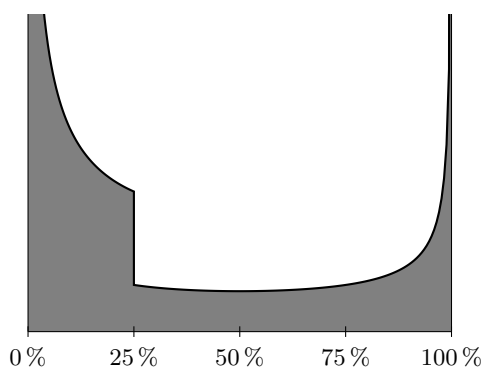


FIGURE 15.2 – Une priore donnant une masse de 60 % à l’hypothèse que le candidat ait une fiabilité inférieure à 1/4, resp. une masse de 40 % à l’hypothèse que sa fiabilité soit supérieure.

d’un nouveau candidat), il sera plus judicieux de poser

$$\mathbb{P}_{\text{pr}}(\theta \in d\theta) = \begin{cases} \alpha_{\text{bon}} \times \mathbb{P}(\text{LoiArcsin} \in d\theta) & \text{pour } \theta \geq 1/4; \\ \alpha_{\text{mal}} \times \mathbb{P}(\text{LoiArcsin} \in d\theta) & \text{pour } \theta < 1/4, \end{cases} \quad (\text{GQ}')$$

avec

$$\alpha_{\text{bon}} := \frac{40 \%}{\mathbb{P}(\text{LoiArcsin} \geq 1/4)}; \quad \alpha_{\text{mal}} := \frac{60 \%}{\mathbb{P}(\text{LoiArcsin} < 1/4)}; \quad (\text{GR}')$$

autrement dit, on *pondère* les différentes possibilités par rapport à la priore arcsinus de façon à avoir, après pondération, une probabilité à priori de resp. 40 % et 60 % pour être resp. un bon ou un mauvais chasseur.

D’ailleurs, même si le club ne dispose pas de statistiques précises sur ses candidats passés, il peut toujours se servir de son expérience pour proposer une estimation “au doigt mouillé” de la distribution des niveaux parmi les candidats potentiels : ainsi, si l’expérience semble montrer qu’il y a un peu plus de “mauvais” chasseurs que de bons, et que parmi ces bons chasseurs, rares sont les chasseurs qui soient carrément “très bons”, le club peut chercher, en guise de priore, une distribution de probabilité sur $]0, 1[$ qui donne une masse de 60 % à l’intervalle $]0, 1/4[$ (les mauvais chasseurs), 30 % à l’intervalle $[1/4, 3/4[$ (les chasseurs assez bons), et 10 % à l’intervalle $[3/4, 1[$ (les très bons chasseurs). En outre, pour éviter de créer des seuils trop nets au niveau des valeurs-seuils 1/4 et 3/4, il peut être judicieux de souhaiter aussi que la densité à priori soit C^1 , ce qui revient à dire qu’on veut que la fonction de répartition de la priore *interpole* de façon régulière les valeurs 0 % en 0, 60 % en 1/4, 90 % en 3/4 et 100 % en 1 : on peut résoudre un tel problème d’interpolation à l’aide de techniques d’analyse numérique^[**], ce que j’ai fait pour obtenir une

[**]. En analyse numérique, vous avez vu qu’une manière de résoudre ce genre de problème d’interpolation est de chercher, parmi les densités de probabilité respectant les contraintes sur la masse totale des intervalles respectifs $]0, 1/4[$, $[1/4, 3/4[$ et $[3/4, 1[$, celle qui minimise l’“énergie” $\int_0^1 f'(x)^2 dx$ ^[††], et que dans ce cas la fonction de répartition d’obtient à l’aide de « splines cubiques ». La méthode que j’ai suivie pour obtenir la figure 15.3 est apparentée à l’idée de minimisation de l’énergie, mais avec quelques fioritures supplémentaires pour tenir comptes des spécificités du problème.

[††]. Notez que c’est le carré de la dérivée *première* de la densité qu’on considère ici, puisque l’interpolation porte sur la fonction de répartition, de sorte que la dérivée seconde de la fonction de répartition est bien la densité première de la densité.

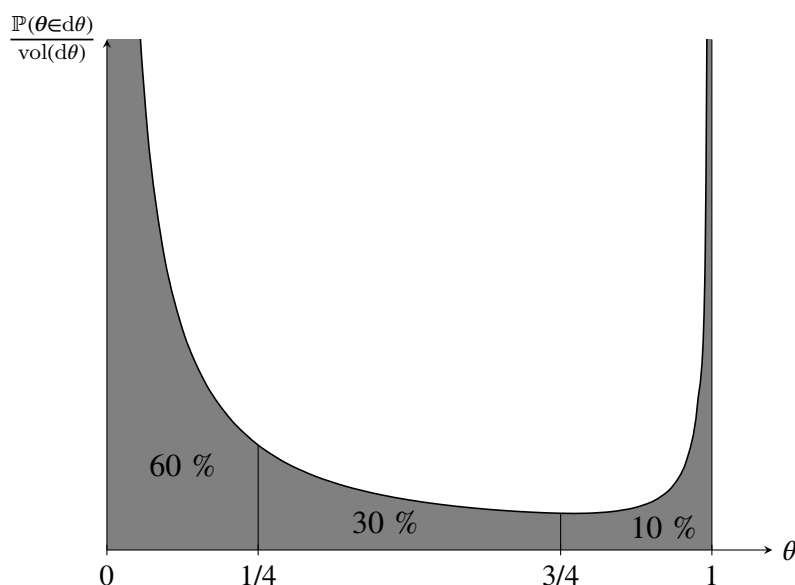


FIGURE 15.3 – La loi représentée ci-dessus propose une priore sur θ traduisant l’expertise informelle du club de chasse : il y a plus de mauvais chasseurs que de bons ; et parmi les bons, rares sont les très bons. Concrètement, on a fait en sorte que 60 % de la masse soit comprise entre 0 et $1/4$, 30 % de la masse entre $1/4$ et $3/4$, et 10 % de la masse entre $3/4$ et 1. On a en outre imposé que le comportement qualitatif de la priore près des valeurs extrêmes 0 et 1 soit le même que celui de la loi arcsinus. (Et, sous l’ensemble des contraintes ci-dessus, on a choisi une priore dont la densité soit “aussi régulière que possible”).

distribution de probabilité “aussi raisonnable que possible” vérifiant les contraintes requises, dont vous pouvez voir le tracé en figure 15.3.

La priore de la figure 15.3 est donc une traduction en termes numériques de l’expertise de notre club, que nous avons initialement formulée en termes assez vagues :

« Il y a un peu plus de mauvais chasseurs que de bons ; et parmi ces bons chasseurs, rares sont les chasseurs qui soient carrément très bons. »

Il importe néanmoins de noter que le fait de fixer les seuils à $1/4$ et $3/4$, de prendre des proportions de 60 %, 30 % et 10 %, et le choix de la méthode d’interpolation, tout cela a été fait “au doigt mouillé”^[††] : la façon dont le club a traduit ici son expertise sous forme d’une distribution de probabilité est parfaitement honnête ; néanmoins, son expertise n’étant pas quantifiable en des termes parfaitement précis, cette traduction garde une part d’arbitraire... !

♣

!! **Remarque (GS ')**. Lorsqu’on estime une priore par expertise, on sait que notre choix

[††]. Concernant le seuil de $1/4$ et la proportion de 60 %, vous aurez noté que, dans l’exemple précédant celui-ci, ils n’avaient pas été définis “au doigt mouillé”, mais respectivement en tant que *définition* de la notion de « bon chasseur » et que proportion *observée* historiquement. Néanmoins, même si j’ai repris les mêmes valeurs pour l’exemple actuel, on n’est pas obligé de considérer cette fois-ci que ces valeurs correspondent resp. à une définition et à une observation historique exactes : en fait, l’idée ici est plutôt que, sans jamais rien avoir formalisé rigoureusement ni compté précisément, le club considère qu’il considère les candidats comme mauvais lorsque leur fiabilité est inférieure à *environ* $1/4$, et que le « un peu plus » de mauvais chasseurs que de bons correspond à une proportion d’*environ* 60 %.

risque de ne refléter qu'imparfaitement la réalité... Dans ce cas, entre deux maux possibles, il est essentiel de choisir le moindre, à savoir : il vaut mieux risquer d'avoir une priore qui attribue une probabilité de 10^{-3} à quelque chose qui se produit en fait avec probabilité 10^{-6} , que l'inverse ! Dit en d'autres termes, il vaut mieux que notre priore soit trop bien répartie entre les différentes possibilités (en restant donc assez prudent sur la façon de traduire en priore les indices dont on dispose) que l'inverse. Ou encore : il vaut mieux garder un soupçon d'approche "neutre" (voir § 15.5) dans la composition de notre priore experte. Ou encore encore : il vaut mieux rester plutôt ouvert à la probabilité d'un phénomène "exceptionnel"... ☺

Remarque (GT'). La justification de la remarque ci-dessus provient de concepts de théorie de l'information : pour mesurer "de combien se trompe" une probabilité approchée Q par rapport à une probabilité véritable P en termes de manque d'information, l'outil adéquat est la *divergence de Kullback–Leibler* $\mathcal{D}_{\text{KL}}(P \parallel Q)$. On peut montrer que dans la situation évoquée ci-dessus, si on estime que la probabilité d'un événement rare est η au lieu de ε , cela donne une divergence de Kullback–Leibler proportionnelle à $\eta - \varepsilon - \ln(\eta/\varepsilon)\varepsilon$. L'application numérique indique alors que le fait de croire qu'un événement rare a une probabilité 10^{-6} au lieu de 10^{-3} en réalité a un cout, en termes d'information, environ 6 fois plus grand que de croire qu'un événement rare a une probabilité de 10^{-3} au lieu de 10^{-6} en réalité. ☺

Exemple (GU'). L'exemple ci-dessous montre pourquoi il ne faut pas suivre aveuglément les indices dont on dispose. Dans le cadre de l'exemple (GP'), nous avons déterminé la priore à partir du fait que, dans le passé, 16 candidats s'étaient révélés mauvais contre seulement 4 bons ; et nous en avons déduit que la probabilité à priori qu'un nouveau candidat soit un bon chasseur devait être de l'ordre de 20 %. Imaginons maintenant que *tous* les 20 candidats précédents rencontrés par le club aient été mauvais. Aurait-il fallu en déduire qu'il était (à priori) *certain* qu'un nouveau candidat soit un mauvais chasseur ? Évidemment non : après tout, si la proportion de bons chasseurs parmi les candidats potentiels vaut en réalité, disons, 2 %, il n'y a même rien de bien étonnant à n'avoir trouvé aucun bon chasseur parmi les candidats précédents ! (Puisque sur 20 candidats, on doit alors trouver seulement 0,4 bons chasseurs en moyenne). Il est donc plus raisonnable de supposer que la probabilité à priori que le nouveau candidat soit un bon chasseur est certes faible, mais pas nulle ; typiquement, de la prendre de l'ordre de 2,5 % (soit la moitié d'un vingtième), puisque c'est le seuil à partir duquel il n'y a plus rien d'étonnant à n'avoir trouvé aucun bon chasseur parmi les 20 candidats précédents. ☺

Remarque (GV'). En fait, une façon mathématiquement plus solide de traiter le cas de l'exemple ci-dessus (où tous les candidats passés étaient de mauvais chasseurs) aurait consisté à utiliser un argument « hyper-bayésien » : dans ce cas, on considère que la priore elle-même a été tirée aléatoirement selon une « hyper-priore », autrement dit une distribution de probabilité *sur les distributions de probabilité sur Θ* , qu'on aurait choisie en fonction d'un argument neutre (voir le paragraphe suivant), puis ajustée au vu des résultats des 20 candidats précédents en vertu de ce que nous avons vu à la § 8.2. (Il faut ensuite prendre l'espérance pour transformer la postérieure sur la priore en *une* priore effective). On trouverait alors, dans le cas où tous les chasseurs précédents étaient mauvais, une probabilité *à priori* pour que le chasseur suivant soit bon de 2,4 %. De manière intéressante, si on applique le même argument dans le cas où 4 candidats sur 20 se sont avérés bons dans le passé, on arrive par l'argument hyper-bayésien à trouver qu'une priore raisonnable doit donner à notre nouveau candidat une probabilité de 21,4 % : si cela ne change pas grand-chose au schmilblick (et montre donc que notre proposition de priore de l'exemple (GP') était pertinente), on voit cependant que nous avons effectivement *équilibré* la priore (au sens où nous avons considéré que le probabilité d'être un bon chasseur était un peu moins faible que ce suggéraient les données) par rapport à ce qu'aurait suggéré une confiance trop aveugle

en nos données passées : ainsi, l'argument hyper-bayésien rejoint les conclusions de la remarque (GS') !

Néanmoins, les approches hyper-bayésiennes dépassent le cadre de ce cours, et nous ne les aborderons donc pas davantage ici. ♣

! *Remarque (GW')*. En fait, l'exemple (GU') ci-dessus illustre une idée qui est un peu moins générale que celle exprimée par la remarque (GS'), mais qui a l'avantage d'être plus facile à appréhender et à vérifier :

La priore doit donner une masse non nulle à *toutes* les zones de l'espace du paramètre caché ! (À moins, s'entend, qu'il n'existe un argument pour éliminer certaines zones avec une certitude *absolue*).

En termes plus formels, cela revient à dire (en utilisant les notations génériques) que, lorsque Θ est muni d'une topologie naturelle, notre priore doit donner une masse strictement positive à n'importe quel ouvert non vide de Θ . ♣

15.5 Priors non informatives

Voyons enfin le cas où on n'a absolument *aucune* information sur ce que peut bien pouvoir valoir θ (pas même par expertise) ; auquel cas le choix le plus raisonnable est manifestement de choisir une priore qui donne « la même chance à chaque possibilité ». Typiquement, si Θ est un espace fini, on choisira comme loi à priori pour θ la loi uniforme sur Θ .

Exemple (GX'). Dans le cas du chasseur, où $\Theta =]0, 1[$, un choix raisonnable de priore non informative serait ainsi de prendre $\text{Loi}(\theta) = \text{Unif}^{\text{me}}(0, 1)$. ♣

Exemple (GY'). Dans le modèle du pédagogue, demandons-nous quelle loi à priori on souhaite mettre sur le paramètre de dispersion σ , à valeurs dans \mathbb{R}_+^* . Il pourrait être tentant de prendre la priore (impropre) proportionnelle à la mesure de Lebesgue sur \mathbb{R}_+^* ; néanmoins, quand on y réfléchit bien, ce n'est pas si satisfaisant que cela. En effet, une priore uniforme sur \mathbb{R}_+^* pour σ reviendrait à dire, par exemple, qu'on considère qu'il est à priori deux fois plus probable d'avoir σ entre 100 et 110 que d'avoir σ inférieur à 5 ; or, on a plutôt envie de dire que c'est le second évènement qui devrait être plus probable à priori : car il y a des façons très différentes les unes par rapport aux autres d'avoir une dispersion inférieure à 5 (une dispersion de 4 n'a rien à voir avec une dispersion de 0,4!), alors que le premier évènement nous donne sur σ une information qu'on perçoit intuitivement comme beaucoup plus spécifique... En fait, quand on y réfléchit mieux, on réalise qu'il est intuitivement plus satisfaisant de vouloir que ce soit l'*ordre de grandeur* de σ qui soit uniforme : autrement dit, de donner une densité à priori uniforme à $\log \sigma$, plutôt qu'à σ lui-même. Si on fait cela, on dira en substance qu'on considère comme à priori aussi probable que σ soit entre 2 et 4, qu'entre 10 et 20, qu'entre 50 et 100 : et cette fois-ci, cela colle beaucoup mieux avec notre idée intuitive de « n'avoir aucune information sur ce que vaut le paramètre de dispersion σ » ! ☺

Dès lors, si on prend pour loi à priori sur $\ln \sigma$ la mesure (impropre) proportionnelle à la mesure de Lebesgue sur \mathbb{R} , on en déduit la loi à priori sur σ lui-même en prenant (toujours à constante multiplicative impropre près) la mesure-image de ladite mesure de Lebesgue par la fonction exponentielle : et les calculs nous donnent alors, pour σ , une priore (impropre) caractérisée par

$$\mathbb{P}(\sigma \in d\sigma) \propto \sigma^{-1} \text{vol}_1(d\sigma). \quad (\text{GZ}')$$



Remarque (HA'). Comme le montre l'exemple ci-dessus, le problème de « prendre une priore uniforme » n'est pas si simple qu'il n'en a l'air, car lorsqu'on a affaire à des distributions de probabilité non discrètes, la notion d'« uniformité » n'est pas stable par changement de variables... Voyons une autre situation où c'est encore plus flagrant (et avec des priores propres). Imaginons que, dans le club de chasse des Canardeurs (concurrent de celui du Bouchonnois), on utilise une échelle pour mesurer la qualité d'un chasseur : la « qualité » (au sens des Canardeurs) d'un tireur, noté cette fois-ci τ , étant sa probabilité de réussite *après deux essais*. Il est alors évident que la « qualité » τ et la « fiabilité » θ sont reliées par la relation bijective $\tau := 2\theta - \theta^2$; de sorte que, du point de vue mathématique, décrire le niveau du tireur par θ ou par τ est rigoureusement équivalent (on parle de *reparamétrage*). Cependant, mettre une loi uniforme sur τ n'est *pas* la même chose que mettre une loi uniforme sur θ ...



Comme le montre la remarque ci-dessus, la simple idée d'« uniformité » n'est pas toujours suffisante pour déterminer une priore non informative de façon canonique. Alors, peut-on faire mieux... ? Nous n'entrerons pas dans le détail de cette question dans le cadre de ce cours ; néanmoins je présente ci-dessous en digression, à titre de complément, deux grandes stratégies qu'on peut suivre pour tenter d'être neutre indépendamment du paramétrage.

Définition (HB') (Priore de Haar). Une stratégie pertinente pour déterminer une priore non informative consiste à demander que la distribution à priori sur θ soit invariante par un certain nombre de transformations « naturelles » sur l'espace Θ : par exemple, dans le cas où Θ est le cercle unité (si par exemple θ représente la phase d'un phénomène sinusoïdal), la loi uniforme sur le cercle unité est la seule mesure qui soit invariante par toutes les rotations. Une distribution de probabilité caractérisée par de telles propriétés d'invariance est qualifiée de *mesure de Haar*.

Dans le cadre de ce cours, je ne vous demanderai pas de déterminer des priores en tant que mesures de Haar ; en revanche, il se peut que je vous demande de vérifier qu'une priore donnée par une certaine formule vérifie certaines propriétés d'invariance, et que le cas échéant, vous sachiez y voir un signe que cette priore vérifie une certaine forme de neutralité, puisqu'elle est indépendante du choix qu'on fait entre différentes façons naturelles de paramétrer Θ (la façon de passer d'un paramétrage à l'autre correspondant, justement, aux transformations « naturelles » évoquées ci-dessus).



Exemple (HC'). En fait, la priore sur σ proposée dans l'exemple (GY') *supra* correspond à une priore de Haar. En effet, un groupe de transformations agissant naturellement sur \mathbb{R}_+^* est celui des changements d'échelle : $\sigma \mapsto k\sigma$ pour $k \in \mathbb{R}_+^*$. Or, les mesures proportionnelles à $\sigma^{-1} \text{vol}(d\sigma)$ sont les seules qui soient invariantes par ces transformations, de sorte que la mesure de probabilité (impropre) qui s'en déduit est un choix naturel de priore non informative pour σ .



Définition (HD') (Priore de Jeffreys). La *priore de Jeffreys* a une définition plus subtile, liée à la théorie de l'information, et qu'il serait au-delà de vos connaissances actuelles de définir rigoureusement : il s'agit de dire, en substance, que, pour tout $\theta \in \Theta$, la probabilité que θ soit difficile à distinguer de θ doit être identique.

La priore de Jeffreys présente l'avantage de pouvoir toujours être définie, et ce, de façon dénuée d'ambiguïté, dès qu'on est face à un modèle de statistique paramétrique.



Exemple (HE'). Dans le modèle du chasseur, le calcul de la priore de Jeffreys montre qu'il s'agit de la loi de l'arcsinus ! C'est notamment pour cette raison que j'en ai fait mon choix de priore par défaut dans le cadre de ce polycopié. (Et aussi parce que ce choix présentait un certain nombre d'avantages pédagogiques \smile).



Remarque (HF'). C'est principalement lorsqu'on recourt à des priores^[*] non informatives qu'on se retrouve confronté à des priores impropres. (Peu importe la méthode non informative retenue, d'ailleurs : aussi bien la stratégie uniforme que celles de Haar et de Jeffreys produisent facilement des mesures de masse totale infinie!). Ce qui n'est pas un cas si rare que cela : car comme nous l'avons fait remarquer plus haut, il vaut mieux éviter de se fier trop aveuglément à son avis sur la valeur de θ , et donc utiliser (au moins partiellement) une forme d'approche non informative dans son raisonnement... ♣

Avant de clore ce chapitre, mentionnons un concept très intéressant émergeant naturellement en théorie de l'information, et qui permet de poser une priore "universelle" sur absolument tout type de données. [Nota : ci-dessous je ne définis pas certains des termes intervenant (comme « machine de Turing universelle »), dans la mesure où mon but est simplement de vous donner un aperçu des idées exposées, pas de vous faire un cours de théorie de l'information]. Voici ce que j'aimerais dire ici à ce sujet :

Définition (HG'). La *distribution de probabilité de Solomonoff* est la distribution de probabilité sur $\{0, 1\}^{\omega}$ (autrement dit l'ensemble des suites infinies de bits) qu'on obtient à partir d'une machine de Turing universelle qui n'efface jamais sa sortie, lorsque la bande d'entrée et la bande de sortie sont initialisées à une suite de bits i.u.d. Informellement, cela correspond à la loi de la sortie d'un programme informatique aléatoire (écrit dans un langage suffisamment général).

"La" distribution de Solomonoff est en fait dépendante du choix de machine de Turing universelle opérée; mais on peut montrer que tous les choix possibles amènent à des probabilités ayant des densités bornées les unes par rapport aux autres, de sorte qu'on peut considérer qu'on a bien défini « une » distribution de probabilité "à densité relative bornée près". ♡

Définition (HH'). L'*inférence de Solomonoff* consiste à encoder l'ensemble des données passées et futures sous la forme d'une suite de bits (les données passées étant autodélimitées et précédant les données futures), et à faire de l'inférence prédictive sur des données futures en considérant que notre jeu de données a été généré à priori selon la distribution de probabilité de Solomonoff. ♡

Point (HI'). Il existe des arguments philosophiques forts pour considérer l'inférence de Solomonoff comme la "bonne" façon de faire de la statistique bayésienne en l'absence d'information spécifique à priori. On peut montrer notamment montrer que, dans le cadre des modèles d'échantillonnage paramétriques, l'approche de Solomonoff re-donne les mêmes conclusions qu'on aurait obtenues à partir de techniques d'analyse statistique "classique". En outre, cela permet d'obtenir des réponses solidement justifiées à des questions comme « si, depuis $n \gg 1$ jours que l'humanité existe, le Soleil s'est levé tous les jours, quelle est la probabilité qu'il se lève demain ? » (en supposant ici qu'on n'ait aucune notion d'astronomie par ailleurs), où la théorie de Solomonoff nous dit en l'occurrence que la réponse correcte est de l'ordre de $1/n$. De même, dans le cas du modèle du chasseur, la théorie de Solomonoff fournit un argument en faveur de la priore impropre Bêta(0, 0) par rapport à d'autres choix possibles. Dans les problématiques de *sélection de modèle* que vous verrez notamment en analyse de données et en intelligence artificielle, c'est encore les raisonnements à la Solomonoff qui sont cachés derrière le critère de sélection BIC (pour « *Bayesian information criterion* »). ♣

Remarque (HJ'). On peut néanmoins démontrer que la probabilité de Solomonoff est rigoureusement impossible à calculer (et il s'agit là d'une impossibilité *théorique*, pas seulement pratique : *aucun* algorithme, si compliqué soit-il, ne permet de calculer la probabilité de Solomonoff! — C'est lié aux problématiques d'incomplétude en logique mathématique) : il s'agit donc d'un outil *philosophique* plus que pratique, qui permettra tout au plus d'avoir des heuristiques intelligentes (comme le critère BIC) en fonction de certaines options plutôt que d'autres. L'idée est néanmoins remarquable sur le plan conceptuel, car elle indique que, même en l'absence totale d'information, il existe malgré tout une priore sur nos données qui aboutira à des

[*]. Ou éventuellement à des hyper-priores, confer remarque (GV').

conclusions pertinentes : grâce à la priore de Solomonoff, on peut donc utiliser l'analyse bayésienne comme une théorie *complètement générale* de l'idée de connaissance, nous permettant (sur le papier!) de répondre à *toute* question du type : « quelle croyance devons-nous tirer de telles données au vu de telles informations dont nous disposions précédemment ^[†] ? » ! \clubsuit

[†]. Les « informations dont nous disposions précédemment étant vues », dans le cadre de la théorie de Solomonoff, comme des données passées encore plus anciennes que ce que le modèle appelle « observation passée » : en effet, grâce à son encodage binaire, la théorie de Solomonoff n'a aucune difficulté à agréger des informations de typologies à première vue différentes.

Chapitre 16

Fonctions d'utilité

16.1 Illustration du concept

Dans les parties II et III du cours, nous avons expliqué que, lorsqu'on cherche à évaluer quantitativement la pertinence d'une procédure d'estimation, de prévision ou plus généralement de décision, on le faisait en passant par une *fonction de perte*, qui nous dit à quel point une décision s'avère mauvaise dans une situation donnée (la description de la situation en question étant reflétée, typiquement, par la valeur d'une certaine quantité d'intérêt), et que le but était alors de minimiser l'espérance de valeur prise par ladite perte.

Dans ce chapitre, je voudrais développer plus avant cette notion de perte, dans la mesure où il me semble qu'il s'agit là d'une idée dont la compréhension est absolument essentielle *pour un ingénieur*. Du strict point de vue de la *théorie statistique*, cependant, comprendre précisément le sens des fonctions d'utilité et la façon de les employer en contexte industriel n'est pas central : c'est pourquoi, afin de ne pas trop se disperser dans les enjeux de ce cours, je n'ai indiqué aucun élément de cette annexe comme étant à retenir. Pour autant, je vous conseille vivement d'essayer de lire et de comprendre ce qui suit, parce que ce sera intéressant et important pour vous *de manière générale* ; mais en revanche, lors de vos révisions pré-examen, vous pourrez vous contenter de regarder le passage qui suit en diagonale! 😊

Pour commencer, rappelons la définition de la notion d'utilité (ou de perte) :

Définition (HK'). Lorsqu'on est amené à prendre une décision, représentée mathématiquement comme un élément d'un certain espace \mathcal{D} , et que l'impact de cette décision dépend de la valeur prise par une certaine quantité d'intérêt, représentée mathématiquement comme un élément d'un certain espace \mathcal{G} , pour comparer deux prises de décisions en tenant compte de l'incertitude sur la quantité d'intérêt, l'outil mathématique le plus classique consiste à introduire une *fonction d'utilité* $u : \mathcal{G} \times \mathcal{D} \rightarrow \mathbb{R}$, dont l'interprétation est la suivante :

- (i) La valeur $u(g, \hat{d})$ nous dit à quel point est-ce que la décision \hat{d} amènera des conséquences bénéfiques si la quantité d'intérêt s'avère valoir g : plus cette valeur est grande, plus les conséquences sont bénéfiques à nos yeux ;
- (ii) En outre, entre deux situations incertaines où l'utilité que nous obtiendrons effectivement est aléatoire, notre préférence va à celle où *l'espérance* de l'utilité est la meilleure (et nous sommes indifférents entre les deux situations lorsque l'espérance est identique).

L'opposé d'une fonction d'utilité est qualifié de *fonction de perte* (traditionnellement notée ' \mathcal{L} ') : c'est exactement le même principe, sauf que cette fois-ci nous avons une préférence pour les valeurs les plus *petites* !^[*] ♡

Commentaire (HL'). Détaillons un peu le point (ii) de la définition ci-dessus. Imaginons que je vous donne le choix entre les deux options suivantes :

Option A Je vous offre un cadeau aléatoire, déterminé par le tirage d'une pièce de monnaie non truquée : ce cadeau sera soit un vélo électrique si la pièce tombe sur pile, soit un jeu de quilles finlandaises si la pièce tombe.

Option B Je vous offre un lave-vaisselle, quoi qu'il arrive. (Dans cette option je lance *quand même* la pièce, mais on n'exploitera pas son résultat^[†]).

Selon le choix que vous aurez fait et le résultat de la pièce, vous gagnerez soit un jeu de quilles, soit un lave-vaisselle, soit un vélo : si nous appelons resp. u_0, u_1, u_2 les bénéfices que le gain de chacun de ces trois cadeaux représente à vos yeux, la fonction d'utilité correspondant à cette situation est alors définie sur $\{\text{pile, face}\} \times \{A, B\}$ par

$$\begin{cases} u(\text{pile}, A) = u_2 ; \\ u(\text{face}, A) = u_0 ; \\ u(\text{pile}, B) = u_1 ; \\ u(\text{face}, B) = u_1 . \end{cases} \quad (\text{HM}')$$

Notez bien ici que ces valeurs u_0, u_1, u_2 reflètent vos préférences *personnelles*, avec la subjectivité que cela comporte : ainsi, s'il y a bien du sens à dire qu'un certain triplet (u_0, u_1, u_2) modélise correctement vos préférences (au sens où il conduit effectivement aux décisions que vous prendriez dans une telle expérience), ce n'est pas pour autant qu'un triplet (u_0, u_1, u_2) qui serait correct *pour vous* serait pour autant correct pour votre voisine, et vice-versa !

Essayons maintenant de comprendre plus précisément ce que les valeurs u_0, u_1, u_2 disent sur vos préférences. Déjà, si (comme on peut raisonnablement le supposer) vous préférez le vélo au lave-vaisselle, et le lave-vaisselle aux quilles, cela devra se refléter par le fait que $u_2 > u_1 > u_0$. Mais à part cela, l'écart entre u_1 et u_2 doit-il être plus petit, plus grand, ou égal à celui entre u_0 et u_1 ? Sur ce point, votre réponse devra refléter celle des deux options A ou B ayant votre préférence, la première ayant pour espérance u_1 tandis que la seconde a pour espérance $(u_0 + u_2)/2$! Ainsi, si vous préférez l'option A (autrement dit si, à vos yeux, avoir une chance sur deux de voir votre lave-vaisselle se transformer en vélo électrique suffit à compenser le fait d'avoir un risque sur deux de le voir se transformer en quilles), cela doit se refléter par le fait que $(u_0 + u_2)/2 > u_1$, autrement dit que $u_2 - u_1 > u_1 - u_0$: ce qui traduit le fait que, à vos yeux, l'écart entre le bénéfice du vélo et celui du lave-vaisselle est supérieur à l'écart entre le bénéfice du lave-vaisselle et celui du jeu de quilles. Si à l'inverse vous préférez l'option B, cela signifie que vous attribuez (dans votre

[*]. En fait, dans ce cours, comme dans la plupart des ouvrages de statistique, nous utiliserons le formalisme de la fonction de perte plutôt que celui de la fonction d'utilité : mais j'ai trouvé que la définition centrale était plus aisée à saisir lorsqu'on l'écrivait avec la convention de l'utilité.

[†]. Ce protocole étrange sert simplement à unifier la présentation mathématique entre les deux options. En pratique bien sûr, si vous choisissez l'option B, je ne lancerai pas *réellement* la pièce : le résultat de la pièce sera alors simplement *virtuel* (ce qui n'a aucune incidence, puisqu'il n'intervient pas dans le cadeau que vous obtenez !), et ne servira qu'à simplifier les calculs de comparaison entre les deux options.

système de préférences personnel) moins d'écart entre le vélo et le lave-vaisselle qu'entre le lave-vaisselle et les quilles, ce qui doit donc se traduire par le fait que $u_2 - u_1 < u_1 - u_0$: et cela correspond bien, au niveau des espérances, à la propriété que $(u_0 + u_2) / 2 < u_1$!

Nous voyons ainsi que votre préférence entre les options A et B reflète lequel des deux écarts $u_2 - u_1$ et $u_1 - u_0$ doit être le plus grand (lorsque u_0, u_1, u_2 correspondent aux utilités pour votre cas personnel, s'entend). Mais on peut même aller plus loin, en interprétant le sens à donner à la valeur *précise* du ratio $(u_2 - u_1) / (u_1 - u_0)$! Pour ce faire, supposons que ma pièce de monnaie soit déséquilibrée et ait en fait une certaine probabilité $p \in]0, 1[$ de tomber sur pile, probabilité que vous comme moi connaissons. (Ce n'est donc pas de la triche, juste que le tirage au sort ne se fait pas de façon uniforme). Il est clair alors que votre préférence entre les options A et B va dépendre de la valeur de p : pour p tendant vers de 0, l'option A devient essentiellement équivalente à gagner un jeu de quilles à coup sûr, donc vous préférez l'option B (rappelons que nous avons supposé, pour les besoins de cet exemple, que vous préféreriez le lave-vaisselle au jeu de quilles), tandis que pour p tendant vers 1, l'option A tend vers le gain du vélo électrique à coup sûr, auquel cas vous la préférez à l'option B. En outre, plus p est élevé, plus l'option A devient intéressante à vos yeux : il y aura donc un certain seuil "critique" p_c en-deçà duquel vous préférerez B et au-delà duquel vous préférerez A, et tel que vous serez indifférent·e si la probabilité de pile vaut exactement p_c . Cette indifférence au seuil p_c soit se traduire par le fait que, pour cette valeur de p , les deux utilités doivent devenir identiques :

$$p_c u_2 + (1 - p_c) u_0 = u_1, \quad (\text{HN}')$$

ce qui est équivalent à

$$\frac{u_2 - u_1}{u_1 - u_0} = \frac{1}{p_c} - 1. \quad (\text{HO}')$$

Ainsi, le ratio entre $u_2 - u_1$ et $u_1 - u_0$ traduit le seuil auquel vous basculez entre l'option B et l'option A : lorsque ce seuil est petit ($p_c \ll 1$), même une faible probabilité de gagner le vélo suffit à vos yeux à compenser le fait de risquer de perdre le lave-vaisselle pour les quilles, et cela se voit au fait que $u_2 - u_1 \gg u_1 - u_0$ (vous mettez beaucoup plus d'écart entre le lave-vaisselle et le vélo qu'entre les quilles et le lave-vaisselle), tandis que lorsque ce seuil est grand $1 - p_c \ll 1$, même une faible probabilité de perdre le lave-vaisselle vous décourage de tenter de gagner le vélo, signifiant alors que vous placez beaucoup plus d'écart entre les quilles et le lave-vaisselle qu'entre le lave-vaisselle et le vélo. ♣

16.2 Considérations sur la notion d'utilité

Remarque (HP'). Comme le commentaire (HL') a dû le rendre clair, notez bien que la notion d'utilité, du point de vue industriel, est un concept fondamentalement *subjectif* : elle formalise *votre* niveau de satisfaction relatif concernant telle situation par rapport à telle autre. Ainsi, si une entreprise décide d'intégrer les notions de développement durable parmi ses objectifs, la façon dont cela se traduira mathématiquement sera au niveau de la fonction d'utilité que l'entreprise cherche à maximiser (ou plus exactement, dont elle cherche à maximiser l'espérance) : cette fonction se mettant alors à refléter une notion de succès *globale* incluant le résultat financier et les aspects environnementaux, sociaux et éthiques !... ♣

Remarque (HQ'). Entre autres notions subjectives, la notion d'utilité permet de modéliser l'*aversion au risque* (ou parfois la propension au risque). Imaginez par exemple que vous dirigiez une entreprise dont la capital est d'environ 10 M€. Une ingénieure vient vous proposer une idée risquée : si vous la mettez en pratique, il y a une chance sur deux que cela fasse gagner 200 k€ à l'entreprise, mais une chance sur deux que cela lui fasse perdre 100 k€... Que décidez-vous ? Normalement, vous allez considérer que, *en moyenne*, l'idée en question est gagnante, et que cela mérite donc qu'on la mette en œuvre ! Mais maintenant, imaginez que cette même ingénieure vienne vous proposer une idée qui a une chance sur deux de faire gagner 18 M€ à l'entreprise, et une chance sur deux de lui faire perdre 9 M€... L'espérance de gain est toujours positive, et même 90 fois plus grande que dans la situation précédente ! Pourtant, cette fois-ci vous allez estimer que le risque ne mérite pas d'être couru : en effet, vu le capital de l'entreprise, si l'idée échoue, on ne va pas "juste" perdre de l'argent, mais se retrouver carrément en situation de faillite ou de quasi-faillite, ce qui aura des conséquences catastrophiques dont il sera très difficile de sortir, même à long terme ! Alors qu'en cas de succès, certes, on triple la valeur de l'entreprise, mais d'un point de vue qualitatif, ce n'est pas un changement si radical que cela...

La notion d'utilité permet de fournir un modèle justifiant cette décision mathématiquement : en effet, ce que vous voulez maximiser, ce n'est pas l'espérance du résultat financier *en tant que tel*, mais l'espérance de *l'utilité* de ce résultat ! Et ici le fonction qui, à un résultat financier donné, associera son utilité sera typiquement *concave*, traduisant une *aversion au risque* : entre deux solutions ayant la même espérance financière "brute", votre préférence va à celle qui implique le moins d'aléa ! Pour l'exemple que j'ai donné ci-dessus, une fonction d'utilité qui modélise bien la décision prise consiste à attribuer à un gain financier g une utilité de

$$\ln\left(1 + \frac{g}{10 \text{ M€}}\right). \quad (\text{HR}')$$

On voit alors que ne rien faire a une utilité nulle, gagner 200 k€ a une utilité de +0,0198, perdre 100 k€ a une utilité de -0,0101, gagner 18 M€ a une utilité de +1,03, et perdre 9 M€ a une utilité de -2,30 : on vérifie bien que, en termes d'espérance, appliquer l'idée de l'ingénieure dans la première situation (gain d'utilité moyen : +0,0048) par rapport à ne rien faire, mais que dans la seconde situation, le gain d'utilité moyen (-0,64) est en fait une lourde perte !! ♣

Remarque (HS'). À l'inverse, notez qu'il peut exister des situations de propension au risque : imaginez ainsi un élève de Mines Nancy en grande difficulté sur ses résultats, qui risque très fortement de se faire exclure de l'École, à moins qu'il n'obtienne au moins 15/20 à l'examen de statistique. Cet élève va évidemment faire son maximum pour réussir, mais se pose la question de savoir s'il doit réviser "normalement" en travaillant tous les chapitres du cours, ou faire des "impasses" et privilégier certains sujets très précis en espérant qu'ils tombent à l'examen... Dans le premier cas, il considère qu'il devrait obtenir typiquement environ 12/20 (à plus ou moins un point près) [‡] ; alors que dans le second cas, il a une chance sur deux d'obtenir environ 16/20 si le sujet qu'il a travaillé en priorité tombe à l'examen, et une chance sur deux d'obtenir environ 6/20 si ce sujet ne tombe pas ! Ici, faire l'impasse est une situation *à la fois* plus défavorable en moyenne "brute" et plus risquée ; pourtant, il est clair que c'est dans l'intérêt de l'élève de suivre cette stratégie : car, foutu pour foutu, qu'est-ce que ça changera pour lui d'avoir eu 12 ou 6 s'il doit être exclu dans les deux cas... ? ! Donc, dans une telle situation (rare en pratique, néanmoins ☹), la fonction d'utilité reflèterait une propension au risque. ♣

[‡]. Ah ben oui, hein, on a dit que c'était un élève en difficulté : donc, même en travaillant dur, il va avoir du mal à obtenir un résultat excellent !...

Remarque (HT'). Lorsqu'on introduit une fonction d'utilité, la valeur zéro pour l'utilité n'a aucune signification particulière, et il n'y a pas non plus d'unité de mesure pour quantifier les utilités : ainsi, du point de vue "physique", transformer les valeurs d'une fonction d'utilité selon une bijection affine croissante conduit exactement à la *même* modélisation — et, en particulier, les conclusions qu'on tirera quant aux décisions qu'il convient de prendre seront *exactement* les mêmes ! En termes d'utilité, tout ce qui a du sens, c'est :

- Le sens de comparaison entre deux situations : telle situation est-elle préférable à telle autre ?
- Entre trois situations A, B, C pour lesquelles on préfère A à B et B à C, le *ratio* entre l'ampleur de notre préférence pour A sur B et l'ampleur de notre préférence pour B sur C [§].

♣

Remarque (HU'). Dans le contexte statistique, on peut aussi observer que cela ne change rien, du point de vue des décisions à prendre, de remplacer une fonction d'utilité $u(g, \hat{d})$ par une variante $u(g, \hat{d}) + c(g)$ qu'on a translatée selon un terme ne dépendant que de g . Cela expliquera pourquoi, lorsqu'on s'intéressera aux fonctions de perte concernant l'estimation, on pourra toujours se permettre de supposer qu'une reconstitution exacte correspond à une perte de zéro. ♣

Remarque (HV'). Cette modélisation de la prise de décision en contexte incertain par le fait qu'on souhaite maximiser l'espérance d'une utilité est clairement intéressante ; mais on peut se demander néanmoins s'il n'existerait pas des situations de prise de décision qu'on ne pourrait pas modéliser ainsi... Le *théorème de von Neumann-Morgenstern* énonce que non : sous certaines hypothèses très raisonnables exprimant des contraintes de rationalité minimales sur un processus de prise de décision, on montre que nos décisions correspondront nécessairement à optimiser l'espérance d'une certaine fonction d'utilité. C'est pourquoi, en mathématiques et en économie, on utilise quasiment toujours ce formalisme pour modéliser les questions de décision optimale. ♣

[§]. Du point de vue philosophique, cela entraîne des considérations particulièrement intéressantes. À cause de la propriété d'invariance que nous venons d'expliquer, si, demain, vous vous mettiez à tout ressentir de façon deux fois plus intense, aussi bien la joie que la tristesse, en toute chose, alors votre façon de vous comporter resterait *exactement* la même. Pourtant, du point de vue d'une personne souhaitant suivre une morale « utilitariste » (c-à-d. une morale dont le but est de maximiser le bonheur total de l'humanité), l'exacerbation de vos ressentis devrait justifier que vos besoins soient satisfaits en priorité par rapport à ceux des personnes dont la fonction d'utilité varie sur des plages moins larges... Mais comment une telle personne pourra-t-elle savoir que vos ressentis sont exacerbés, puisque cette exacerbation est *intrinsèquement* impossible à mesurer ? Et d'ailleurs, est-il juste de fonder une philosophie morale sur un concept (à savoir, le bonheur total de l'humanité) qu'il est ainsi fondamentalement impossible de mesurer ? C'est là une objection philosophique très importante à l'utilitarisme, que Robert NOZICK a popularisée avec son expérience de pensée du « monstre utilitariste », une personne dont les utilités seraient si énormément exacerbées que cela justifierait que tout le reste de l'humanité se sacrifie pour le moindre de ses désirs... Il est tentant, de ce fait, de vouloir refuser de comparer l'utilité entre deux personnes différentes. Mais pourtant, lorsqu'on doit prendre une décision favorisant une personne sur une autre (exemple : « si je fais des crêpes, ma fille sera contente et mon garçon déçu ; mais si je fais des gaufres, ce sera l'inverse ; or je n'ai pas le temps de faire les deux à la fois : que dois-je donc choisir ? »), notre souhait de justice va bel et bien nous amener à comparer les utilités entre des personnes distinctes ! (Ainsi, si mon garçon n'a qu'une très légère préférence pour les gaufres sur les crêpes, alors que ma fille adore les crêpes et déteste les gaufres, il semble logique que je prépare des crêpes, dans l'intérêt collectif !). Est-ce que cela est possible ? Est-ce que cela a seulement du sens ? Est-ce que cela ne risque pas de nous faire préférer injustement ceux qui ressentent les émotions plus violemment ?... La notion mathématique d'utilité, sous son apparence anodine, amène de grandes questions philosophiques ! ☺

Remarque (HW'). Dans le cadre de ce cours, il ne vous sera jamais demandé de proposer vous-mêmes des fonctions d'utilité ; tout au plus de justifier le choix d'une fonction d'utilité qui vous aura été donnée par l'énoncé. ♣

16.3 Un exemple en contexte industriel

Pour finir ce chapitre, présentons un exemple de contexte industriel dans lequel il y a une fonction de perte qui s'impose assez clairement (ce qui est souvent le cas en pratique) :

Exemple (HX'). Considérons un gisement contenant une quantité φ de minerai. Après avoir entrepris une campagne de prospection qui lui a donné une idée (entachée néanmoins d'incertitude) sur la valeur de φ , une entreprise doit prendre la décision, notée δ , d'exploiter ou non le gisement (la valeur VRAI signifiant qu'on exploite le gisement, et la valeur FAUX qu'on ne l'exploite pas). On suppose ici qu'exploiter le gisement engendre des frais f fixes, et rapporte a par unité de minerai extrait, les valeurs de f et a étant connues ; on suppose également que, une fois l'exploitation du gisement entamée, même si le gisement se révèle décevant, il sera toujours préférable poursuivre l'exploitation jusqu'au bout, dans la mesure où c'est amont de l'exploitation qu'on engage l'essentiel des frais. La question est : dans ce contexte, quelle fonction d'utilité (ou de perte) prendre par rapport à la quantité d'intérêt φ et à la décision δ ?

Ici, il semble assez raisonnable d'associer la notion d'utilité au bénéfice financier réalisé par l'entreprise. Si on n'exploite pas le gisement, il n'y a ni frais ni perte ; et si on l'exploite, le bénéfice final vaura $a\varphi_{\vee} - f$ (où φ_{\vee} désigne la quantité véritable de minerai contenue dans le gisement) : on prendra donc

$$u(\varphi, \delta) = \mathbf{1}_{\delta}(a\varphi - f). \quad (\text{HY}')$$

Voyons maintenant comment reformuler cette fonction d'utilité dans le formalisme de la fonction de perte normalisée. De manière générale, la perte est l'opposée de l'utilité, soit $\mathbf{1}_{\delta}(f - a\varphi) =: \ell_{\text{brut}}(\varphi, \delta)$ en l'occurrence. Maintenant, normaliser la fonction de perte signifie qu'on va retrancher de celle-ci la valeur minimale qu'elle prend pour une quantité d'intérêt donnée (ce qui ne change pas la signification "physique" de la perte), de sorte que, à quantité d'intérêt donnée, le minimum de la fonction de perte renormalisée soit toujours 0.

Il est évident que, à valeur fixée de φ , le minimum de la perte est atteint pour $\delta = \text{FAUX}$ lorsque $\varphi < f/a$ (il n'est pas intéressant d'exploiter le gisement lorsque la quantité de minerai s'avère trop faible pour couvrir les frais), resp. pour $\delta = \text{VRAI}$ lorsque $\varphi \geq f/a$: le minimum valant alors respectivement 0 et $f - a\varphi$. On obtient donc la fonction de perte normalisée suivante :

$$\ell_{\text{norm}}(\varphi, \delta) = \ell_{\text{brut}}(\varphi, \delta) - \mathbf{1}_{\varphi \geq f/a}(f - a\varphi) = (\mathbf{1}_{\varphi \geq f/a} - \mathbf{1}_{\delta}) \cdot (a\varphi - f). \quad (\text{HZ}')$$

Ce résultat peut se comprendre ainsi : lorsque $\delta = \{\varphi \geq f/a\}$, cela signifie que l'entreprise a pris la "bonne" décision ^[¶], et dans ce cas sa perte est nulle, ce qui

[¶]. Ici je veux dire que la décision apparaît *rétrospectivement* comme la bonne, au sens où cette décision coïncide avec celle qu'on aurait dû prendre si on avait connu à l'avance la quantité de minerai ! Bien entendu, en pratique, il peut arriver qu'on prenne une décision *correctes* (au sens où il est logique de la prendre au vu des informations dont on dispose) qui s'avèrera rétrospectivement

signifie qu'on a fait aussi bien que possible. Lorsque $\varphi < f/a$ mais que $\delta = \text{VRAI}$ (autrement dit, l'entreprise a pris la décision d'exploiter un gisement qui ne contenait pas assez de minerai pour couvrir les frais), la perte vaut $f - a\varphi$ et représente l'argent globalement perdu par l'entreprise (en prenant en compte ce que le minerai a tout de même rapporté) dans cette exploitation qu'il aurait mieux valu ne pas faire. Et lorsque $\varphi \geq f/a$ mais que $\delta = \text{VRAI}$ (autrement dit, l'entreprise a renoncé à exploiter un gisement qui était en fait rentable), la perte vaut $a\varphi - f$: autrement dit, elle représente le *manque à gagner* par rapport à la situation où on aurait exploité le gisement !

Pour finir, on peut aussi définir une fonction de perte liée à l'estimation de la quantité de minerai contenue dans le jugement. Ci-dessus la fonction de perte que nous considérons s'entendait par rapport à une quantité d'intérêt et à une décision. Maintenant, imaginons qu'on cherche plutôt à estimer à quel point c'est grave de croire qu'il y a une quantité de minerai $\hat{\varphi}$ lorsque la quantité réelle est φ . (On cherche donc une fonction de perte où on compare deux quantités φ et $\hat{\varphi}$ qui vivent dans le même espace). En fait, cela se ramène à la situation précédente : on va juste dire que, si $\varphi < f/a$, cela signifie qu'on croit qu'il y a trop peu de minerai dans le gisement pour qu'il vaille la peine d'exploiter celui-ci : dans ce cas, l'entreprise prendra la décision de ne pas exploiter et la perte (normalisée) sera $\mathbf{1}_{\varphi \geq f/a}(a\varphi - f)$; et que si $\varphi \geq f/a$, à l'inverse, notre estimation conduira l'entreprise à exploiter le gisement, d'où une perte de $\mathbf{1}_{\varphi < f/a}(f - a\varphi)$. En résumé, si on note $\delta(\hat{\varphi}) := \{\hat{\varphi} \geq f/a\}$ la décision d'exploitation qui découle logiquement de l'estimation $\hat{\varphi}$, notre perte vaut

$$\ell(\varphi, \hat{\varphi}) = (\mathbf{1}_{\delta(\hat{\varphi})} - \mathbf{1}_{\varphi \geq f/a}) \cdot (f - a\varphi) = \begin{cases} \mathbf{1}_{\varphi \geq f/a}(a\varphi - f) & \text{si } \hat{\varphi} < f/a ; \\ \mathbf{1}_{\varphi < f/a}(f - a\varphi) & \text{si } \hat{\varphi} \geq f/a. \end{cases} \quad (\text{IA}')$$

En l'occurrence, cela peut se ré-écrire plus élégamment à l'aide des fonctions « partie positive » \bullet^+ et « partie négative » \bullet^- :

$$\ell(\varphi, \hat{\varphi}) = \begin{cases} (a\varphi - f)^+ & \text{si } \hat{\varphi} < f/a ; \\ (a\varphi - f)^- & \text{si } \hat{\varphi} > f/a. \end{cases} \quad (\text{IB}')$$

♣

Remarque (IC'). Ici j'ai choisi d'assimiler la fonction de perte à la perte financière de l'entreprise en tant que telle : autrement dit, je suppose que passer d'un bénéfice de 10 M€ à un bénéfice de 12 M€ est exactement aussi intéressant que passer d'un déficit de 1 M€ à un bénéfice de 1 M€... En pratique, dans bien des cas l'intérêt véritable qu'on a à gagner de l'argent diminue à mesure qu'on en a gagné déjà beaucoup, alors qu'à l'inverse être ruiné est une catastrophe qu'il convient d'éviter absolument... Autrement dit, la véritable fonction d'utilité est plutôt une fonction concave du gain financier : en termes d'économie, on parle dans une telle situation d'*aversion au risque* [1]. Cependant, lorsqu'on ne parle que de petits montants

“mauvaises” (au sens où, après coup, on se rend compte qu'une autre décision aurait conduit à un meilleur résultat, à cause d'éléments dont on ne disposait pas encore au moment où on a pris notre décision) : dans ce cas, ce n'est pas parce qu'on n'a pas pris la “bonne” décision qu'on en serait blâmable pour autant ! ☹

[1]. À noter qu'il existe aussi, quoique ce soit plus rare, des situations de *propension au risque* ; conf. p.ex. la remarque ?? *supra*.

(par exemple, si la mine que vous hésitez à exploiter n'est qu'une des centaines de concessions que possède votre entreprise multinationale), le phénomène de concavité est si peu marqué pour les montants mis en jeu par une seule mine qu'on peut alors le négliger et assimiler la fonction de perte à (l'opposé de) le bilan financier de l'exploitation, comme je l'ai fait ici. ♣

Remarque (ID'). Dans l'exemple (HX'), nous avons pu déterminer la fonction de perte pertinente pour notre problème sans vraiment d'ambiguïté. Cependant, autant la *définition* de la fonction de perte est bien claire dans un contexte industriel (c'est ce qu'on « gagne » si on sait que φ vaut φ et qu'on agit en conséquence, moins ce qu'on « gagne » si on croit que φ vaut $\hat{\varphi}$), autant sa *quantification* sous forme numérique n'est pas évidente : car évaluer numériquement « ce qu'on gagne » est assez subjectif... Déjà, il y a l'aversion au risque dont j'ai parlé dans la remarque précédente. Mais il y a aussi, notamment, tous les enjeux autour de la *responsabilité sociétale de l'entreprise*, qui sont difficiles à quantifier lorsqu'on doit les comparer aux enjeux financiers ! Imaginons ainsi que vous ayez reçu une offre d'un fournisseur en Amérique du Sud, offre intéressante financièrement et dont la légalité ne pose pas question, mais que vous ayez de bonnes raisons de soupçonner que ce fournisseur soit en fait lié au narcotrafic local, et que lui donner des débouchés aiderait indirectement les gangs à prospérer. Sous réserve que vous ayez un minimum d'éthique, vous devrez alors, dans la construction de votre fonction de perte, tenir compte, non seulement du bilan financier (clairement positif) qu'un tel changement de fournisseur impliquerait pour votre entreprise, mais aussi de l'impact sociétal (clairement négatif) de celui-ci en Amérique du Sud ! Mais comment pondérer les deux ?... Une réponse naïve consisterait à introduire un poids *infini* à l'aspect éthique, en disant « j'ai une éthique, donc je ne veux faire que des choses parfaitement propres ». Mais on se rend vite compte que le choix de l'infini ne marche pas en pratique, parce que les choses sont plus compliquées que cela : par exemple, peut-être que le lien de votre fournisseur avec le narcotrafic est en fait très ténu, et que passer un accord avec lui ne dégradera pas grand-chose à la situation des gens sur place, alors que d'un autre côté, les bénéfices de votre entreprise seront si considérables que cela lui permettra, grâce à ses produits bon marché, d'aider la vie de millions de clients, d'offrir de meilleurs salaires à des milliers d'employés, de faire des dons humanitaires... Donc, dans tous les cas, il y a bien une question de *pondération*. Mais la déterminer comportera nécessairement une part de subjectivité... ♣

Chapitre 17

Pièges et paradoxes des tests

17.1 Que signifie « accepter une hypothèse nulle » ?

Déjà, soulignons que, dans la mesure du possible, il est méthodologiquement préférable de conclure l'analyse d'un test booléen en parlant du rejet (ou pas) de l'hypothèse nulle, plutôt que de la validation (ou pas) de l'hypothèse alternative :

Point (IE'). Lorsqu'on écrit la phrase de conclusion d'un test positif, il est généralement recommandé, par souci de clarté méthodologique, de préférer la formulation « on est enclin à rejeter \mathcal{H}_0 » à la formulation « on valide l'hypothèse \mathcal{H}_1 » ; et de même, pour la phrase de conclusion d'un test négatif, il est recommandé de préférer la formulation « l'observation est compatible avec \mathcal{H}_0 » à la formulation « on n'est pas en mesure de prouver l'hypothèse \mathcal{H}_1 ».

En effet, c'est là la philosophie même du principe de test statistique : la conclusion vaut par rapport à l'hypothèse nulle, *et uniquement par rapport à celle-ci* ; l'hypothèse « alternative », comme son nom l'indique, est simplement la réunion de toutes les hypothèses autres que l'hypothèse nulle... Certes, c'est généralement le fait de valider l'hypothèse *alternative* qui nous intéresse *réellement* ; néanmoins, écrire la conclusion en fonction de l'hypothèse nulle limitera les risques de mécompréhension du sens de votre analyse par ceux qui la liront ! ♣

Remarque (IF'). Le point méthodologique ci-dessus est même encore plus flagrant lorsqu'il s'applique à un test par p -valeur, dans la mesure où la probabilité à laquelle cette p -valeur se réfère concerne la loi de la statistique de test *sous l'hypothèse nulle* !. ♣

Par ailleurs, profitons-en pour remettre une couche sur la façon dont il convient d'interpréter resp. la positivité et (surtout !) la négativité d'un test. Nous avons certes déjà évoqué cela dans le chapitre 12 ; mais il s'agit d'une erreur d'interprétation si classique que remettre une couche dessus ne peut pas faire de mal :

Point (IG'). Lorsqu'on dit « on rejette l'hypothèse nulle au risque α », cela signifie « on a observé un phénomène qui a une chance moindre que α de se produire sous \mathcal{H}_0 ; et donc il est peu vraisemblable que l'hypothèse \mathcal{H}_0 soit effectivement vérifiée ». Cela est une conclusion d'autant plus forte que α sera petit. ♣

Point (IH'). À l'inverse, « accepter l'hypothèse nulle au risque α » ne signifie rien d'autre que « ne pas être en mesure de rejeter l'hypothèse nulle au risque α ». En particulier, ce n'est absolument pas synonyme de « rejeter l'hypothèse alternative au risque α » ! ♣

!

!!

Remarque (II'). Notons au passage, que, pas contraposition du point (IG'), la conclusion « accepter l'hypothèse nulle au risque α » est celle fois-ci une conclusion d'autant plus *faible* que α est petit ! Néanmoins, gardez bien à l'esprit qu'une telle conclusion d'acceptation n'a qu'une portée très faible dans tous les cas, peu importe la valeur de α ...

Remarque (IJ'). On notera ainsi qu'il est tout-à-fait possible et pas choquant du tout d'accepter (disons, au risque 10 %) plusieurs hypothèses contradictoires entre elles : cela signifie juste qu'aucune de ces hypothèses ne peut être exclue par le résultat du test !

Exemple (IK'). Pour illustrer la remarque ci-dessus, imaginons qu'un détective enquête sur le meurtre qui vient d'avoir lieu au cours d'une soirée dans un manoir. L'unique témoin n'a pu voir qu'une vague silhouette de l'assassin : celui-ci avait, selon son témoignage, la stature d'un homme. (On suppose ici que tous les convives masculins ont grosso-modo la même stature, plus grande que celle de toutes les convives féminines). Bien sûr, la témoin a pu se tromper ; mais on admettra qu'il y a moins d'une chance sur dix pour que, voyant la silhouette d'une femme, elle la prenne pour celle d'un homme.

Du point de vue statistique, la quantité d'intérêt (explicative) à laquelle s'intéresse notre détective est alors l'identité de l'assassin ; et l'observation sur laquelle il s'appuie pour inférer cette quantité d'intérêt est le récit de la témoin. Dans ces conditions, lorsque le détective se demande si Madame Pervenche est l'assassin, il peut rejeter l'hypothèse nulle {M^{me} Pervenche est l'assassin} au risque 10 % : puisque, fût-elle réellement coupable, il y avait moins de 10 % de risque que sa silhouette fût jugée masculine. Ainsi, le test innocenté cette suspecte ! (au risque 10 %).

En revanche, si on applique le même test aux hypothèses nulles respectives {le Colonel Moutarde est l'assassin}, {le Révérend Olive est l'assassin} et {Le Professeur Violet est l'assassin}, cette fois-ci ces trois hypothèses seront toutes acceptées ! Est-ce que cela signifie qu'ils sont tous coupables ? Évidemment non, puisque le témoignage assure ici que l'assassin est unique ! Cela veut juste dire qu'on n'est, à ce stade, en mesure d'innocenter aucun des trois ☹

17.2 Comment choisir l'hypothèse nulle

Caractère fermé de l'hypothèse nulle

Lorsqu'on procède à un test statistique pour trancher entre deux hypothèses complémentaires, une question qui se pose naturellement est de savoir qui choisir comme hypothèse nulle et qui choisir comme hypothèse alternative. Bien souvent, il s'avère qu'on n'a pas vraiment le choix — ce qui peut d'ailleurs être fâcheux, car il y a souvent une hypothèse qu'on aimerait pouvoir démontrer, et pour cela on aurait besoin d'en faire notre hypothèse alternative... On a en effet la règle générale suivante (plus d'explications dans la suite) :

! **Principe (IL')**. On choisira toujours^[*] une hypothèse nulle correspondant à un sous-ensemble topologiquement fermé de l'espace du paramètre caché.

En particulier, lorsque le test est défini par une ou plusieurs inégalités, il convient que ces inégalités soient larges du côté de l'hypothèse nulle (et donc strictes du côté de l'hypothèse alternative).

D'autre part, lorsque le test consiste à se demander si une certaine égalité (parfaite) est vérifiée ou pas, le cas « les deux termes sont égaux » devra toujours correspondre à l'hypothèse nulle, et le cas « les deux membres sont distincts » à l'hypothèse alternative... !

De manière plus générale, on peut retenir que « l'hypothèse nulle doit contenir le cas de référence », où par « cas de référence », je fais référence au cas d'égalité d'indépendance, etc. : soit que cette hypothèse corresponde au cas de référence lui-même, soit qu'elle corresponde à une inégalité large consistant à être au moins aussi [quelque chose] que ce cas de référence. \diamond

Remarque (IM'). En fait, quand on dit que l'hypothèse nulle est « fermée », il faudrait préciser par rapport à quelle topologie on entend cela... L'énoncé techniquement précis du principe ci-dessus serait que l'hypothèse nulle doit être fermée par rapport à la topologie de la distance $(\theta, \theta') \mapsto \|\text{Loi}_{\theta'}(X) - \text{Loi}_{\theta}(X)\|_{\text{VT}}$, où $\|P' - P\|_{\text{TV}}$ représente la « distance en variation totale » (dont nous ne donnerons pas la définition dans le cadre de ce cours) entre deux distributions de probabilité P et P' sur le même espace \mathcal{X} .

En fait, dans tout ce cours, nous faisons implicitement l'hypothèse que l'application $\theta \mapsto \text{Loi}_{\theta}(X)$ est continue depuis la topologie sur Θ (dont nous rappelons que, dans le cadre de ce cours, il s'agit toujours d'une partie d'un certain espace \mathbb{R}^d : cela définit donc la topologie sur Θ sans ambiguïté) vers la topologie de la distance en variation totale sur $\mathcal{M}_1(\mathcal{X})$. (Et, de fait, tous les exemples présentés satisfont bien cette hypothèse). Dans ces conditions, « être une partie fermée de Θ au sens de la distance en variation totale entre les $\text{Loi}_{\theta}(X)$ » implique « être une partie fermée de Θ vu comme un sous-espace de \mathbb{R}^d » : donc, dire qu'il faut choisir l'hypothèse nulle fermée par rapport à la topologie de la variation totale implique bien qu'il faut choisir cette hypothèse nulle fermée pour la topologie « du modèle » sur Θ . (Même s'il est possible que certaines parties de Θ fermées au sens du modèle ne soient en réalité pas fermées pour la norme en variation totale : auquel cas, nonobstant leur caractère fermé, elles ne pourront pas convenir pour définir une hypothèse nulle intelligente). \clubsuit

Remarque (IN'). Attention : ici, quand je dis que l'hypothèse nulle doit être fermée, c'est par rapport à la topologie sur Θ , pas sur l'espace \mathbb{R}^d (dans lequel Θ est plongé) tout entier ! Par exemple, considérons le modèle du chasseur où on réduit Θ à l'ensemble $]0, 0,5[\sqcup]0,8, 0,9]$: cela correspondrait à la situation où le candidat qui se présente au club prétend être un certain tireur d'élite (dont on sait que, le cas échéant, sa fiabilité est comprise, au sens large, entre 80 % et 90 %), mais qu'on se demande s'il s'agit bien de lui où si ce ne serait pas plutôt son frère jumeau (dont on sait qu'il n'est pas capable d'atteindre 50 % de ses cibles, même si, là encore, on ignore sa fiabilité exacte). Eh bien, dans ce modèle-ci, l'hypothèse $]0, 0,5[$ est, nonobstant les apparences, bel et bien *fermée*^[†], dans la mesure où, *vu comme un sous-espace de Θ* , cet ensemble $]0, 0,5[$ est en fait identique à l'ensemble $[0, 0,5]$ (ou encore, par exemple, à l'ensemble défini par l'inégalité large $\{\theta \leq 0,6\}$). \clubsuit

Remarque (IO'). Attention, l'expression « cas de référence » peut être un peu trompeuse. Mettons par exemple qu'on cherche à rejeter l'hypothèse nulle selon laquelle les garçons seraient, en moyenne strictement plus forts en mathématiques que les filles (pour une certaine définition du concept de « force en mathématique », évaluée sous forme de nombre réel). Comme cette hypothèse n'est pas fermée, il procède du principe (IL') que cela ne constitue pas un bon choix d'hypothèse

[*]. « Toujours » doit s'entendre ici comme « toujours, en pratique » : du point de vue technique, absolument rien n'empêche de choisir une hypothèse nulle non fermée ; néanmoins, comme nous l'expliquerons plus loin, cela ne présenterait aucun intérêt le cas échéant...

[†]. Elle est néanmoins ouverte *aussi*, ainsi qu'on pouvait le voir immédiatement d'après sa forme

nulle. Nous pourrions alors être tentés de dire qu'on n'a qu'à tester l'hypothèse nulle que « les garçons soient au moins aussi forts que les filles » ; sauf que cette hypothèse nulle ne nous intéresse pas non plus : car ce qu'on aimerait prouver, ce n'est pas que les filles soient strictement plus fortes que les garçons (on ne s'attend d'ailleurs pas à ce que cela soit vrai, pas plus qu'on ne s'attend à ce qu'il soit vrai que les garçons soient strictement plus forts que les filles) ; c'est que les niveaux moyens des filles et des garçons soient quasiment identiques...!

Dans ce cas, il s'avère que l'hypothèse nulle qu'il sera pertinent de tester est en fait « les garçons ont un avantage strict en mathématiques sur les filles, et pas complètement infime qui plus est ^[‡] » : autrement dit, si dans notre modèle, le niveau moyen des garçons est notée μ_G , et que le niveau moyen des filles est noté μ_F , notre hypothèse nulle correspondra à $\{\mu_G - \mu_F \geq \varepsilon\}$, où ε est une certaine valeur strictement positive, mais néanmoins extrêmement petite, de sorte qu'on peut considérer que les cas où $\{0 < \mu_G - \mu_F < \varepsilon\}$ correspondent *en pratique* à une égalité essentiellement parfaite entre garçons et filles !

Dans ce cas, vous pourriez m'objecter « hé ; cette nouvelle hypothèse nulle correspond certes à une inégalité large ; mais elle ne contient pas le cas de référence où garçons et filles ont le même niveau moyen ! ». Certes. Sauf qu'ici, ce n'est pas cela qu'il faut entendre par « cas de référence » : le « cas de référence », en fait, désigne la situation, pouvant être décrite par une égalité, *par rapport à laquelle notre hypothèse nulle se place* : et en l'occurrence, il s'agit donc du cas d'égalité $\{\mu_G = \mu_F + \varepsilon\}$ (qui signifie « les garçons ont le plus petit avantage imaginable sur les filles qui mériterait, le cas échéant, d'être signalé »), cas qui est bel et bien inclus dans notre hypothèse nulle ! ☺

Expliquons maintenant les « preuves » du fait qu'il convient de toujours prendre l'hypothèse nulle fermée :

Premier argument pour le principe (IL '). La raison pour laquelle il convient de placer le cas de référence au sein de l'hypothèse nulle provient de la philosophie même du test statistique (booléen), qui nous demande de tirer une conclusion qui soit valable au risque α *quelle que soit l'hypothèse ponctuelle* θ au sein de l'hypothèse nulle \mathcal{H}_0 . (Et ici, notez que la dissymétrie entre les hypothèses nulle et alternative se manifeste clairement, puisqu'on ne pose, en revanche, aucune exigence qui soit uniforme sur Θ_1). Par conséquent, si l'hypothèse \mathcal{H}_0 est trop « grande », et qu'elle peut conduire à un ensemble de conséquences trop diverses, on ne parviendra pas construire de test statistique pertinent... Idéalement, il faut donc que l'hypothèse nulle soit réduite à un seul point, ou qu'on puisse plus ou moins se ramener à cette situation : ce point auquel on peut se réduire correspondant au fameux « cas de référence », il est naturel de l'inclure au sein de l'hypothèse nulle, puisque, de toutes façons, elle devra s'y ramener... ! ☺

Deuxième argument pour prendre l'hypothèse nulle fermée. Une seconde raison pour toujours prendre l'hypothèse nulle fermée est qu'on peut toujours fermer l'hypothèse nulle sans qu'on y perde rien ! Plus précisément, soit Θ_0 un sous-ensemble non nécessairement fermé de Θ et *Test* un test ^[§] de l'hypothèse $\{\theta \in \Theta_0\}$. Alors,

[‡]. Avec le vocabulaire qui sera introduit dans la § 17.5, cette idée que l'avantage soit « pas complètement infime » correspond à dire qu'on exige qu'il y ait une *taille d'effet* pas trop ridicule pour le phénomène de différence éventuelle entre filles et garçons.

[§]. J'écris ce point avec le cas d'un test booléen en tête ; mais en fait cela fonctionne exactement de la même façon dans le cas d'un test par *p*-valeur.

si nous notons $\bar{\Theta}_0$ l'adhérence de Θ_0 , il se trouve que *Test* est aussi un test de l'hypothèse Θ_0 sans qu'on ait besoin d'en changer le niveau ou la p -valeur ! Et comme $\bar{\Theta}_0$ est un sur-ensemble de Θ_0 , cela signifie qu'en considérant qu'on a en fait testé l'hypothèse $\{\theta \in \bar{\Theta}_0\}$, dans les cas où on rejette l'hypothèse nulle, la conclusion obtenue est plus forte, puisqu'on rejette plus de possibilités... \diamond

L'argument qui précède peut même être mis en valeur en tant que théorème :

Théorème (IP'). *Si Test est un test binaire au niveau α (resp., si P est un test par p -valeur) de l'hypothèse nulle $\{\theta \in \Theta_0\}$, où Θ_0 est une partie non nécessairement fermée de l'espace Θ du paramètre caché, alors, notant $\bar{\Theta}_0$ la fermeture topologique de Θ_0 au sein de Θ , Test constitue automatiquement un test de niveau α (resp. P constitue un test par p -valeur) de l'hypothèse nulle $\{\theta \in \bar{\Theta}_0\}$ aussi.* \diamond

Troisième argument pour prendre l'hypothèse nulle fermée. Un dernier argument pour dire qu'une hypothèse nulle non fermée n'est pas pertinente et liée à l'idée de *consistance* du test, ou plus généralement à son risque de seconde espèce. En effet, il découle du théorème (IP') que, si θ_1 est un point situé dans $\bar{\Theta}_0$ mais pas dans Θ_0 , alors on aura $\mathbb{P}(\text{Test} \mid \theta = \theta_1) \leq \alpha$, de sorte que la fonction de risque de seconde espèce, en θ_1 , vaudra au moins $1 - \alpha$: ce qui est catastrophique... Certes, le théorème ?? nous dit que, de manière générale, le *supremum* de cette fonction de risque de seconde espèce vaudra au moins $1 - \alpha$; néanmoins, autant faire en sorte que, dans la mesure du possible, ce supremum ne soit pas *atteint* pour autant ! C'est d'autant plus que, à supposer qu'on s'intéresse à un régime asymptotique « $n \rightarrow \infty$ » pour notre modèle, on voudrait que, dans ce régime, notre test soit consistant, autrement dit, que le risque de deuxième espèce $\beta^{(n)}(\theta)$ tende, lorsque $n \rightarrow \infty$, vers 0, et ce pour tout $\theta \in \Theta_1$. Néanmoins, s'il existe un θ_1 dans $\bar{\Theta}_0 \setminus \Theta_0$, un tel θ_1 sera bien dans Θ_1 , mais vérifiera cependant $\beta^{(n)}(\theta) \geq 1 - \alpha$ pour tout n : et dès lors, ce point de Θ_1 empêchera le test d'être consistant : même en prenant un paramètre n arbitrairement grand, lorsque θ est égal à θ_1 , il n'est pas possible de s'apercevoir qu'on n'est pas dans l'hypothèse nulle... !

Ainsi, non seulement fermer l'hypothèse nulle ne nuit pas au niveau de notre test (confer l'argument précédent), mais *ne pas* la fermer, par contre, nuirait à sa puissance ! Il est donc toujours préférable de fermer l'hypothèse nulle. \diamond

Remarque (IQ'). L'argument ci-dessus a été écrit dans le formalisme des tests booléens ; mais il se transpose également (au prix d'un peu de complexité technique) au cas des tests par p -valeur \smile \clubsuit

Remarque (IR'). L'argument précédent montre que le fait de prendre l'hypothèse nulle fermée n'est pas tant un *choix* du statisticien, qui ne serait affaire que de convention, mais bien une *contrainte*, due au principe même du test statistique, s'il souhaite pouvoir effectivement démontrer l'hypothèse alternative (quitte à devoir se placer en régime asymptotique) chaque fois qu'elle se présente... !

Par exemple, dans le modèle du pédagogue, imaginons que notre enseignant se dise « d'accord ; je n'arrive peut-être pas à prouver que ma méthode fait strictement progresser les élèves, mais je pourrais au moins essayer de prouver qu'elle ne les fait pas régresser ! », en souhaite ainsi tester l'hypothèse nulle $\{\mu < \mu_{\text{réf}}\}$ en vue de la réfuter. La problématique s'il choisit une telle hypothèse nulle, c'est que la méthodologie même du test impose qu'il y ait une probabilité au moins $(1 - \alpha)$ d'accepter l'hypothèse nulle lorsqu'on a $\{\mu = \mu_{\text{réf}} - 10\}$, mais aussi lorsqu'on a $\{\mu = \mu_{\text{réf}} - 1\}$,

« Tous les champignons sont mortels, sans exception.
Mais certains mettent plus de cent ans à agir. »



FIGURE 17.1 – Si nous appelons $\tau_v \in [0, \infty]$ le temps au bout duquel l'ingestion d'un champignon finit par tuer, l'affirmation « Tel champignon est mortel » correspond à l'hypothèse $\{\tau < \infty\}$, tandis que l'affirmation « Tel champignon est sans danger » correspond à l'hypothèse $\{\tau = \infty\}$. On voit donc que l'hypothèse de mortalité n'est pas fermée dans $[0, \infty]$; dès lors, elle ne saurait constituer un choix d'hypothèse nulle judicieux pour un test : en effet, si on admet que l'effet d'un champignon peut être arbitrairement long, il devient impossible de réfuter sa létalité... Bien entendu, Anicet est de mauvaise foi, car en pratique on peut mettre une borne au temps d'effet d'un champignon, de sorte que l'affirmation de sa mortalité devient alors une hypothèse parfaitement testable...! (source : *Anicet le Pingouin*, « À en mourir de rire »)

lorsqu'on a $\{\mu = \mu_{\text{réf}} - 0,1\}$, lorsqu'on a $\{\mu = \mu_{\text{réf}} - 0,01\}$, ... Et on pourrait alors montrer, par passage à la limite, quand dans ce cas, on aura nécessairement une probabilité au moins $(1 - \alpha)$ d'accepter l'hypothèse nulle lorsqu'on sera dans le cas $\{\mu = \mu_{\text{réf}}\}$. Or, *justement*, ce qui intéresserait notre enseignant, c'est d'être capable, même lorsque l'impact de sa méthode sur les élèves est nul, de *prouver* qu'elle n'a pas un impact négatif. Mais ça, si grand que soit le nombre d'élèves qu'il teste, ce n'est juste pas possible : puisque dans tous les cas, il aura une très grande probabilité de conclure que le résultat observé est compatible avec l'hypothèse d'une régression !

On peut le comprendre de façon plus intuitive en disant qu'un test de détection de la chute du niveau des élèves qui devrait réagir à la moindre chute, si petite soit-elle, sera tellement sensible qu'il sera obligé de réagir même quand il n'y a plus de chute du tout...

Dans le gag de la figure 17.1, on retrouve à nouveau cette idée qu'un test est incapable de détecter les points de l'hypothèse alternative qui seraient au bord de l'hypothèse nulle : même si, cette fois-ci, l'exploitation qui en est faite est de nature purement humoristique! ☺

Intérêt pour l'hypothèse alternative

La sous-section ci-dessus concernait une caractéristique *technique* de l'hypothèse nulle : il faut la choisir fermée, pour des raisons purement *mathématiques*. Mais mettons que nous ne soyons pas soumis à cette contrainte, par exemple parce que nous devons trancher entre deux hypothèses simples (voir aussi l'importante remarque ?? *infra* sur ce point). Dans ce contexte, qui faut-il choisir comme hypothèse nulle, et qui faut-il choisir comme hypothèse alternative ?

!! Principe (IS'). Dans le cadre de la statistique fréquentiste, si on a le choix concernant les rôles à jouer par les deux hypothèses complémentaires qu'on considère, deux considérations doivent entrer en jeu :

- (i) *S'il y a une des deux hypothèses dont la preuve apporte une information plus intéressante que l'autre, par exemple, si la conclusion « on a prouvé \mathcal{H}_a » apporte, le cas échéant (car, bien sûr, \mathcal{H}_a peut aussi être fausse, ou trop difficile à prouver), une information plus intéressante que la conclusion « on a prouvé $\neg\mathcal{H}_a$ », alors cette hypothèse doit jouer le rôle d'hypothèse alternative : en effet, si notre est négatif, la conclusion qu'il apportera ne sera de toutes façons que très peu intéressante : donc autant se concentrer sur ce qui serait le plus intéressant pour nous dans le cas d'un test positif !*
- (ii) *S'il y a une des deux hypothèses en faveur de laquelle pencher à tort serait plus fâcheux que l'autre (par exemple, si trancher en faveur de \mathcal{H}_a alors qu'on a en réalité $\neg\mathcal{H}_a$ a des conséquences plus fâcheuses que trancher en faveur de $\neg\mathcal{H}_a$ alors qu'on a en réalité \mathcal{H}_a , alors cette hypothèse doit jouer le rôle d'hypothèse nulle : en effet, de par la nature du test, le seul risque qu'on soit effectivement en mesure de contrôler est le risque de première espèce, donc il convient d'appliquer ce contrôle au risque qu'on cherche le plus fortement à éviter !*

◇

Exemple (IT'). Le cas (i) est typique de la recherche pharmaceutique : le laboratoire qui développe une molécule donnée souhaite *prouver* que sa molécule apporte un bénéfice thérapeutique, afin d'obtenir son autorisation de mise sur le marché ! Dès lors, il procèdera à des tests où l'hypothèse nulle sera que la molécule est inefficace, voire nuisible. En effet, si l'hypothèse d'une efficacité trop faible était l'hypothèse alternative, on serait dans une situation stupide : soit le test serait négatif et alors on ne pourrait pas conclure de façon ferme, soit il serait positif et alors ça n'aurait de toutes façons quasiment aucun intérêt d'apprendre qu'on a travaillé pour rien... !

♣

Exemple (IU'). Le cas (ii) est quant à lui typique de l'expertise judiciaire : notre philosophie contemporaine de la justice, en effet, est fondée sur le principe qu'il est plus grave de condamner un innocent que d'acquitter un criminel ! (D'où le principe « le doute doit profiter à l'accusé »^[¶] — ou, en latin : « *In dubio pro reo* »). Par conséquent, lorsqu'une expertise judiciaire est appliquée (par exemple, pour analyser une trace de sang sur les lieux du crime), c'est l'hypothèse que l'accusé n'a rien à voir avec ce qui est reproché (donc, en l'occurrence, ce sang n'est pas le sien) qui doit être prise comme hypothèse nulle !

♣

Remarque (IV'). Si on se trouve à la fois dans la situation du cas (i) et dans la situation du cas (ii), “normalement” ils doivent jouer dans le même sens. En effet, s'il est plus grave de rejeter \mathcal{H}_a à tort que de la croire à tort, la prudence consiste à agir en partant du principe que \mathcal{H}_a est vraie : et dès lors, avoir la confirmation que \mathcal{H}_a est effectivement vraie ne serait pas si intéressant que ça, puisque de toutes façons nous aurions agi comme si elle était vraie... Ainsi, l'hypothèse qu'il est le plus grave de rejeter à tort est aussi celle qu'il est le moins intéressant de prouver : dans les deux cas, cela nous indique de prendre cette hypothèse comme hypothèse nulle !

♣

[¶]. Confer p.ex. l'article 304 du code de procédure pénale, ou l'arrêt « Boutaffala v^s Belgique » rendu par la Cour Européenne des Droits de l'Homme le 28 juin 2022.

Remarque (IW'). On peut néanmoins s'amuser à chercher des cas où ces deux principes jouent en sens inverse. Prenons par exemple le problème auquel s'est retrouvée confrontée l'armée américaine au cours de la seconde guerre mondiale : y a-t-il, ou pas, la possibilité d'exploiter l'énergie nucléaire pour créer une arme révolutionnaire (à savoir, la bombe atomique)?^[I] On a clairement l'impression qu'il est plus intéressant de savoir que la bombe atomique est bel et bien possible (il faut alors mettre le paquet dessus!) que d'apprendre qu'elle est en fait impossible (le cas échéant, *business as usual*, quoi!). Pourtant, croire à tort en l'impossibilité de la bombe atomique aurait des conséquences beaucoup plus grave (si l'Allemagne nazie, de son côté, se mettait à la développer, ce serait la défaite des Alliés) que de croire à tort à sa possibilité (on y dilapiderait certes une partie de l'effort de guerre, mais rien de critique...!). ♣

Remarque (IX'). Cela dit, lorsqu'il est possible de faire jouer *soit* le rôle d'hypothèse nulle, *soit* le rôle d'hypothèse alternative, à chacun de nos deux hypothèses complémentaires, autant faire les deux tests correspondants : cela ne coûte en effet qu'une augmentation de l'effort d'analyse statistique^[**], pour un résultat qui a deux fois plus de chances d'être exploitable...! ♣

! *Remarque (IY')*. Il arrive fréquemment que le principe (IS') nous donne envie de tester une hypothèse nulle de la forme $\{\varphi > \varphi_{\text{réf}}\}$ ou $\{\varphi \neq \varphi_{\text{réf}}\}$: par exemple, si le pédagogue souhaite s'assurer que sa méthode n'est pas nuisible, ou si une entreprise veut s'assurer qu'elle ne discrimine pas ses employés en fonction de l'âge... Or, de telle hypothèses ne sont pas fermées, et semblent donc inaptes à jouer le rôle d'hypothèse nulle!

Néanmoins, dans un tel cas, nos hypothèses nulles sont équivalentes *en pratique* à des hypothèses nulles fermées, où l'on prend juste un tout petit peu de marche par rapport à l'égalité parfaite du cas de référence. Par exemple, au lieu de tester l'hypothèse nulle « la nouvelle méthode fait strictement régresser les élèves », l'enseignant peut choisir de plutôt l'hypothèse nulle « la nouvelle méthode fait régresser les élèves *d'au moins* 0,1 pt » : car, bien franchement, si elle fait régresser les élèves de 0,05 pt, *en pratique*, c'est comme si elle n'avait aucun effet! De même, notre entreprise pourra tester l'hypothèse nulle « il existe un niveau discrimination avec l'âge suffisamment important pour avoir un impact réel sur la vie des employés » : car, à supposer qu'on arrive à démontrer que, toutes choses égales par ailleurs, un employé en fin de carrière touche 1 € par an en moins qu'une employée en début de carrière, on n'a pas envie de dire qu'un tel niveau de "discrimination" pose quelque souci que ce soit...! Ce très léger déplacement du seuil de référence que pouvoir se ramener à des hypothèses nulles fermées est très lié à l'idée de *taille d'effet* que nous présenterons dans la § 17.5.

De ce fait, *en pratique*, il est toujours possible de définir n'importe quelle hy-

[I]. De notre côté de l'histoire, la question paraît stupide! Mais en réalité, lorsque le projet *Manhattan* a été lancé, il n'était pas clair du tout qu'il fût effectivement possible de réaliser effectivement une bombe atomique, et encore moins, le cas échéant, qu'une telle arme apportât un avantage militaire réel le cas échéant : ainsi, la fameuse lettre Einstein-Szilárd, enjoignant les États-Unis d'entreprendre des recherches dans cette direction, sous-estime le potentiel réel des bombes atomiques à un niveau frappant : ces physiciens, nonobstant leur exceptionnel génie à tous les deux, y écrivent en effet simplement qu'« il est concevable, quoique bien moins certain, que des bombes d'un nouveau type et extrêmement puissantes [puissent] être assemblées. Une seule bombe de ce type, transportée par bateau et explosant dans un port, pourrait très bien détruire l'ensemble du port ainsi qu'une partie de la zone aux alentours. Toutefois, de telles bombes pourraient très bien s'avérer trop lourdes pour un transport aérien »...

[**]. Sachant que, bien souvent, c'est plutôt la collecte des données qui est le facteur limitant...

pothèse comme nulle ou alternative tout en restant mathématiquement cohérent, quitte à utiliser la petite perturbation expliquée ci-dessus ☺

Voir aussi le théorème (IP') à ce sujet.

17.3 Trop de tests tuent le test !

Si vous procédez à deux tests différents pour une même hypothèse \mathcal{H}_0 , disons un test T_1 au risque α_1 et un test T_2 au risque α_2 , vous pouvez alors définir un test de synthèse T à partir de ces deux tests : ce test T sera défini comme étant positif si *au moins un des deux* tests T_1 et T_2 l'est — attendu qu'une seule conclusion positive suffit à rejeter définitivement \mathcal{H}_0 , alors qu'à l'inverse une conclusion négative ne permet de l'accepter que jusqu'à preuve du contraire. Mais attention : quel est le niveau du test global T ?... C'est $\alpha_1 + \alpha_2$ [††], puisque sous l'hypothèse nulle, il y a un risque α_1 d'échouer au premier test et un risque α_2 d'échouer au second...

Remarque (IZ'). Cela est évident, me direz-vous. Pourtant, on peut facilement se faire piéger en oubliant cet argument... Imaginons ainsi que je procède à 20 tests T_0, \dots, T_{19} de la même hypothèse nulle, chaque test ayant un risque de 5 %. Le test global aura un risque de 100 %, c.-à-d. qu'il ne signifiera strictement rien : de la sorte, même si l'hypothèse nulle est effectivement vérifiée, il n'y aura rien d'étonnant à ce qu'au moins un test T_i s'avère positif — disons T_{13} pour fixer les idées. Mais si vous vous contentez de clamer « Le test T_{13} est positif, donc je peux rejeter l'hypothèse nulle au risque 5 % », votre conclusion sera malhonnête, car vous n'aurez pas mentionné que pour arriver à un tel résultat vous avez *aussi* procédé à 19 autres tests négatifs, et donc que la valeur de 5 % n'est pas le vrai niveau du test réellement opéré... C'est exactement ce qui se passe dans la bande dessinée de la figure 17.2. ☺

Principe (JA'). *Moralité : quand on fait plusieurs tests de la même hypothèse, il faut mentionner tous les tests opérés, y compris ceux qui se sont révélés négatifs, pour avoir le vrai niveau de confiance, lequel est la somme des niveaux de chaque test. Notez qu'il s'agit là d'un vrai souci en recherche scientifique, car on tend à omettre de publier les résultats non concluants (ceux pour lesquels les tests se sont relevés négatifs), ce qui conduit à considérer comme significatifs (c.-à-d. ayant un petit niveau de risque) des résultats qui ne le sont en fait pas...* ◇ !!

Principe (JB'). *Dans le même esprit, il est essentiel de définir votre protocole de test avant d'observer les données : sinon, en bricolant une région d'acceptation ad hoc conçue précisément pour accepter ou rejeter les résultats que vous avez obtenus, vous pourriez accepter ou rejeter n'importe quoi à votre guise...* ◇ !!

Remarque (JC'). Lorsqu'on est face à une situation comme celle de la figure 17.2 où l'hypothèse alternative de notre test apparaît comme une union de plusieurs sous-hypothèses alternatives individuelles “élémentaires”, il existe heureusement des méthodes pour pouvoir quand même regrouper les sous-tests individuels en un test global ayant le bon niveau de significativité ! La plus simple de ces méthodes est

[††]. $(\alpha_1 + \alpha_2)$ est la meilleure borne sur le niveau du test global dont on puisse être assuré dans le cadre présenté ici. On notera toutefois qu'en pratique, les tests T_1 et T_2 seront souvent indépendants voire positivement corrélés, auquel cas on pourra un peu (et parfois même beaucoup) améliorer cette borne...

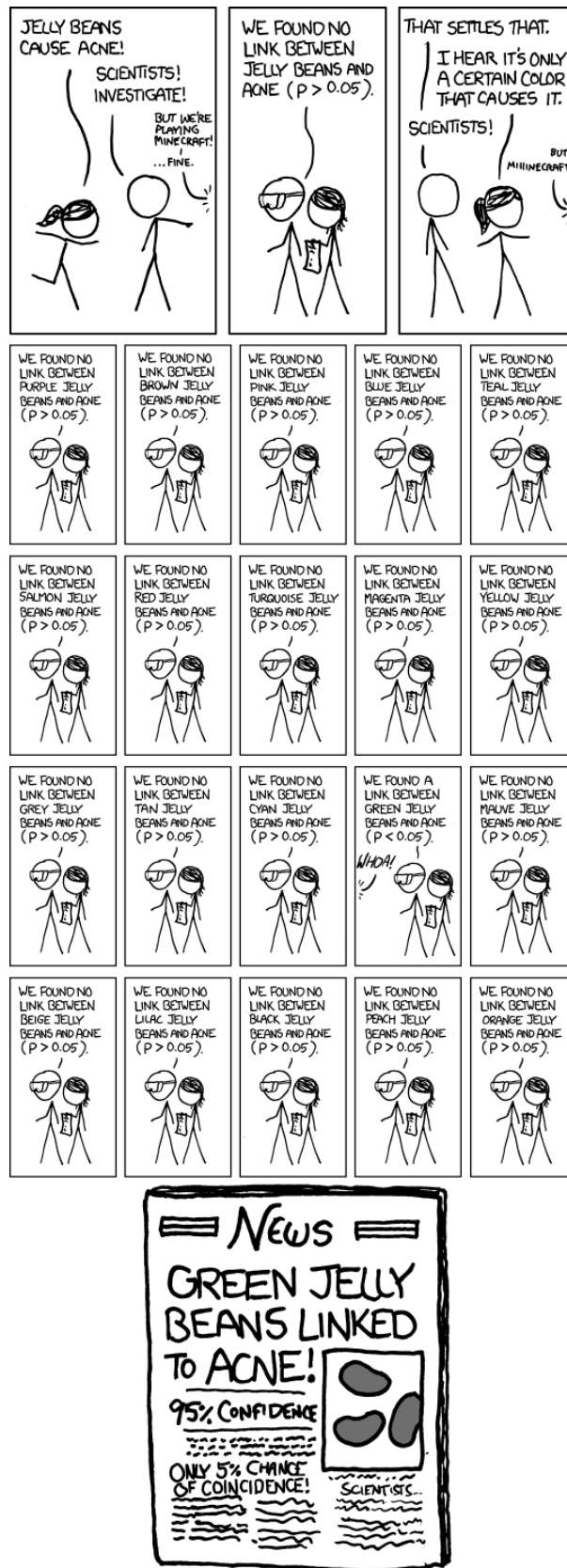


FIGURE 17.2 – Trop de tests tuent le test... (source : <http://xkcd.com/882/>)

la *correction de Bonferroni*, qui consiste, s'il y a m hypothèses alternatives individuelles ($m = 20$ dans l'exemple de la figure 17.2), à rejeter l'hypothèse nulle globale au niveau α lorsque l'une au moins des hypothèses nulles individuelles est rejetée au niveau α/m . (Ce qui revient à dire que, si les p -valeurs trouvées pour chaque sous-test individuel sont p_1, \dots, p_m , on renverra une p -valeur globale de $(\min_{1 \leq i \leq m} p_i) \times m$). Cela dit, on n'a rien sans rien : et lorsqu'on veut tester une hypothèse nulle conjonction de m sous-hypothèses élémentaires en utilisant cette correction de Bonferroni, la puissance du test devient souvent médiocre comparé à ce qu'elle était pour chaque sous-test individuel... ♣

Une autre façon de se faire piéger par l'adage « trop de tests tuent le test », plus insidieuse dans la mesure où fréquemment on ne s'en aperçoit pas, est d'oublier qu'IL FAUT DÉFINIR LE TEST AVANT DE REGARDER LES DONNÉES ! Ne pas le faire est une des erreurs les plus faciles à commettre. Vous observez un jeu de données, par exemple les décimales de π . Là, vous observez un phénomène apparemment remarquable : entre la position 762 et la position 767, il y a 6 décimales '9' successives^[‡‡] ! Dès lors vous vous demandez : est-ce vraiment une coïncidence extraordinaire, ou est-ce que cela peut se produire facilement sous l'effet du seul hasard ? Pour cela, vous décidez de regarder la statistique de test « plus longue série de chiffres consécutifs identiques parmi les 767 premiers chiffres », sous l'hypothèse nulle que « les chiffres sont statistiquement indistinguables d'une série de chiffres i.u.d. » : on trouve une p -valeur de 6,84 ‰^[*], ce qui est tout de même très faible. Faut-il en déduire qu'on a levé un grand mystère dans les décimales de π sur lequel la science doit se pencher dare-dare ?...

... En fait, sans doute pas. Certes, la probabilité d'observer 6 chiffres identiques à la suite dans une série de 767 chiffres aléatoires *est* très faible, rien à redire là-dessus. Mais aurions-nous choisi *spontanément* de regarder la longueur de la plus longue séquence de chiffres identiques parmi les 767 premiers chiffres si nous n'avions pas d'abord *vu* que les décimales de π présentaient cette série de six '9' ? Certainement pas avec la valeur 767 en tout cas : à la rigueur, nous aurions peut-être regardé la plus longue séquence parmi les 1 000 premiers chiffres (ce qui fait déjà légèrement augmenter la p -valeur calculée, à 8,92 ‰). Mais pourquoi 1 000, après tout ?... Nous aurions pu plutôt choisir de nous pencher sur les 100 premiers chiffres, ou les 10 000 premiers chiffres, ou... ! Et puis, nous aurions pu trouver tout aussi remarquable, par exemple, qu'un chiffre apparaisse 130 fois ou plus parmi les 1 000 premiers, ou qu'une séquence de 5 chiffres y apparaisse deux fois consécutivement (du genre, « ... 4124541245 ... »), ou qu'il y ait une série de 17 chiffres consécutifs tous compris entre 0 et 4, ou que la même séquence de 8 chiffres apparaisse à deux endroits différents (du genre, « ... 93779974 ... 99377974 ... »), etc., autant d'évènements dont la probabilité de survenue est de l'ordre de 1 ‰ — et qui, eux, n'arrivent pas pour les vraies décimales de π ^[†]... Autrement dit, puisque nous avons fait notre test *en fonction* des données, c'est comme si nous avions fait tous les tests possibles selon ce qu'on aurait pu remarquer et eu envie de tester ! Et on se retrouve dans la situation précédente : trop de tests tuent le test...

Il faut donc vraiment insister sur ce point :

[‡‡]. Authentique : Vous pouvez vérifier ! ☺

[*]. Et non 7,62 ‰ comme on serait enclin à le penser, car le calcul de cette p -valeur est assez piégeux...

[†]. Oui, j'ai vérifié ☺

Principe (JD'). Quand on procède à une étude statistique, il faut décider de ce qu'on va tester avant de connaître les données, sinon on voit des phénomènes significatifs là où il n'y a rien! \diamond

Remarque (JE'). En pratique, comme il est difficile de définir complètement à l'avance le protocole d'analyse, ce qu'on fait souvent est qu'on regarde un premier jeu de données, en fonction de ce qu'il semble montrer, on se dit qu'on aurait aimé tester telle chose, puis on fait le test sur un second jeu de données indépendant du premier : là, il n'y a pas de biais. (Cette pratique est appelée *validation croisée*; elle est très utilisée dans le domaines de l'apprentissage automatique que vous verrez en 2^e année). Pour reprendre l'exemple de π , on pourrait par exemple essayer de confirmer l'hypothèse qu'il y a tendance à avoir de longues séries de chiffres identiques parmi les $\sim 1\,000$ premiers chiffres en regardant le développement de π en base 11 (lequel n'a, à priori, rien à voir avec celui en base 10) : et là, patatras, on ne trouve rien de mieux que 3 chiffres identiques consécutifs, soit une p -valeur supérieure à 0,5... \clubsuit

Remarque (JF'). En réalité, lorsqu'on ne peut vraiment pas se permettre de procéder à de la validation croisée (par exemple, parce qu'on considère les données liées à un évènement passé qu'on ne peut pas reproduire à volonté), on peut quand même faire de l'analyse statistique, mais il faut alors être très prudent pour que les méthodes proposées ne tombent pas dans les pièges évoqués ci-dessus : on sera alors amené à utiliser des techniques spécifiques dans le style de la correction de Bonferroni... On qualifie de *méthodes d'analyse post hoc* les techniques destinées à faire de la statistique correcte sur des données collectées avant d'avoir défini les éléments que nous souhaitons regarder. \clubsuit

17.4 Le paradoxe du test à droite

Plaçons-nous dans la situation suivante. On considère des données $x_{0\vee}, \dots, x_{15\vee}$ issues de la réalisation de variables aléatoires X_0, \dots, X_{15} de même loi $\mathcal{N}(\mu_{\vee}, \sigma_{\vee}^2)$, données pour lesquelles on a mesuré une moyenne empirique $m_{\vee} = 0,5$ et un écart-type empirique $s_{\vee} = 1$. On voudrait regarder ce qui se passe quand on teste l'hypothèse nulle $\{\mu = 0\}$ et quand on teste l'hypothèse nulle $\{\mu \leq 0\}$ (sachant que σ_{\vee} est inconnu). Dans les deux cas, on va utiliser que

$$\frac{M - \mu_{\vee}}{S / \sqrt{15}} \stackrel{\mathbb{P}_{\vee}}{\sim} T_{\text{St}}(15), \quad (\text{JG}')$$

où $T_{\text{St}}(d)$ désigne la loi de Student de paramètre d .

Pour tester l'hypothèse nulle $\{\mu = 0\}$, on écrit que si $\mu_{\vee} = 0$, on doit avoir

$$\frac{M}{S / \sqrt{15}} \stackrel{\mathbb{P}_{\vee}}{\sim} T_{\text{St}}(15), \quad (\text{JH}')$$

d'où, avec un risque de 5 %,

$$\frac{M}{S / \sqrt{15}} \in [-2,14, 2,14]. \quad (\text{JI}')$$

Ici la réalisation du membre de gauche vaut 1,94, de sorte qu'on ne peut pas rejeter l'hypothèse $\{\mu = 0\}$.

Si maintenant on souhaite tester l'hypothèse nulle $\{\mu \leq 0\}$, on doit alors écrire que sous cette hypothèse, on a

$$\frac{M}{S / \sqrt{15}} \leq \frac{M - \mu_{\vee}}{S / \sqrt{15}} \stackrel{\mathbb{P}_{\vee}}{\sim} T_{\text{St}}(15). \quad (\text{JJ}')$$

Dans ce cas-là, on n'a pas d'autre choix que de faire un test à droite^[‡] : observant que pour $X \sim T_{St}(15)$, on a $X \leq 1,76$ au risque 5 %, on en déduit à fortiori que, sous l'hypothèse que $\{\mu \leq 0\}$, on doit avoir

$$\frac{M}{S/\sqrt{15}} \leq 1,76; \quad (\text{JK}')$$

et puisque la réalisation du membre de gauche toujours 1,94, cette fois-ci l'hypothèse est rejetée.

Voilà qui extrêmement choquant : on n'est pas en mesure de rejeter l'hypothèse $\{\mu = 0\}$; mais à partir de la même statistique de test, et pour le même risque, on peut rejeter l'hypothèse pourtant *plus large* $\{\mu \leq 0\}$...!! Comment expliquer ce paradoxe ?

La première remarque, c'est que le test consistant à regarder si $M/(S/\sqrt{15}) > 1,76$ est *aussi* un test de l'hypothèse nulle $\{\mu = 0\}$: la logique est donc sauve puisque, eussions-nous utilisé *exactement* le même test pour les deux hypothèses nulles, le rejet de l'hypothèse la plus large aurait bien entraîné celui de l'hypothèse la plus étroite...

Se pose alors naturellement la question : pourquoi ne pas avoir aussi utilisé un test à droite pour l'hypothèse $\{\mu = 0\}$? La réponse est à chercher dans le risque de seconde espèce. En effet, le défaut du test à droite est qu'il sera très peu efficace pour détecter des valeurs *strictement négatives* de μ : dans un tel cas en effet, aussi grand soit $|\mu_{\nu}|$, le test sera réussi avec au moins 95 % de probabilité, ce qui est très mauvais... Ce souci rend le test à droite, quoique *correct* formellement, *inadapté* en pratique au test de l'hypothèse $\{\mu = 0\}$, alors qu'il ne se pose pas pour tester l'hypothèse $\{\mu \leq 0\}$.^[§]

« Oui, mais quand même... », me direz-vous. Le paradoxe subsiste : si vous me demandez sans plus de précision s'il est plausible que μ soit nul, la procédure adéquate m'amènera à vous répondre « oui » ; alors que si vous me demandez s'il est plausible que μ soit négatif ou nul, je répondrai « non »... Y a-t-il une explication « morale » à ce phénomène... ?

Eh bien, oui ! Car *poser une question, c'est déjà y répondre un peu...*^[¶] Quand je vous demande « μ est-il nul ? », vous en déduisez intuitivement que cela sous-entend que μ a de bonnes chances d'être nul aussi bien que d'être non nul (disons avec une probabilité à priori de 1/2 pour chaque), sans que cela ne vous dise rien sur le fait que μ soit positive ou négative dans le cas non nul (disons qu'il y a une probabilité à priori de 1/4 pour chacune de ces deux possibilités). À l'inverse, quand je vous demande « μ est-il négatif ou nul ? », vous avez tendance à en déduire qu'il y a des chances équilibrées que μ soit négatif-ou-nul d'une part, ou strictement positif d'autre part (disons avec une probabilité à priori de 1/2 pour chaque), et dans le cas négatif ou nul, vous n'avez pas de raison de savoir si c'est strictement négatif ou pas, ce qui fait que vous pouvez attribuer, disons, une probabilité à priori de 1/4 à chacune de ces deux possibilités. Cela correspond aux probabilités à priori résumées par le tableau suivant :

Question	$\mathbb{P}_{pr}(\mu < 0)$	$\mathbb{P}_{pr}(\mu = 0)$	$\mathbb{P}_{pr}(\mu > 0)$
$\{\mu = 0\}$?	25 %	50 %	25 %
$\{\mu \leq 0\}$?	25 %	25 %	50 %

[‡]. En effet, si on avait obtenu une conclusion plausible sur $T_{St}(15)$ qui aurait pris une autre forme qu'une majoration, on n'aurait été en mesure de rien en déduire sur $M/(S/\sqrt{15})$, vu que tout ce qu'on a est une majoration de cette variable aléatoire par une quantité suivant la loi $T_{St}(15)$.

[§]. On peut même être un peu plus précis en regardant la puissance du test en termes *asymptotiques*, c.-à-d. en faisant tendre n vers l'infini : asymptotiquement, le test à droite sera rejeté à coup sûr dès que $\mu_{\nu} > 0$, mais pas si $\mu_{\nu} < 0$. Ainsi, le risque asymptotique de seconde espèce est uniformément nul (i.e., le test est consistant) sur l'hypothèse alternative $\{\mu > 0\}$, mais pas sur l'hypothèse alternative $\{\mu \neq 0\}$, qui est trop large...

[¶]. Que la question posée soit révélatrice du résultat auquel on s'attend est l'évidence même dans la vie de tous les jours : ainsi, pour mesurer les connaissances d'un enfant, on lui posera à priori des questions plus faciles que face à un adulte, parce qu'on s'attend à des performances plus faibles...

Maintenant, faisons une modélisation grossière où nous disons que la probabilité d'obtenir une valeur de l'ordre de 1,94 vaut selon le cas :

Situation	Proba conditionnelle
$\mu_{\mathcal{J}} < 0$	1 %
$\mu_{\mathcal{J}} = 0$	4 %
$\mu_{\mathcal{J}} > 0$	75 %

En appliquant le théorème de Bayes, on trouve alors des probabilités à postériori de :

Question	$\mathbb{P}_{\text{post}}(\mu < 0)$	$\mathbb{P}_{\text{post}}(\mu = 0)$	$\mathbb{P}_{\text{post}}(\mu > 0)$
« $\mu = 0 ?$ »	1,2 %	9,5 %	89,3 %
« $\mu \leq 0 ?$ »	0,6 %	2,6 %	96,8 %

Conclusion : les probabilités à postériori sont tellement différentes selon la façon dont est posée la question (à cause de la différence dans les probabilités à priori), que la probabilité à postériori de $\{\mu \leq 0\}$ n'est que de 3,2 % lorsqu'on a posé la seconde question, alors que lorsqu'on a posé la première question, la probabilité à postériori de la seule hypothèse $\{\mu = 0\}$ vaut 9,5 %... Il n'est donc pas surprenant que le second test échoue, mais pas le premier ζ

Cette explication n'est pas du tout aussi foireuse qu'elle n'y paraît : en réalité, au prix de quelques modifications mineures, on peut réinterpréter une bonne partie de la théorie des tests en termes bayésiens...

17.5 p -valeur vs taille d'effet

Dans beaucoup de logiciels appliquant les statistiques, pour un test d'indépendance par exemple, le résultat de l'analyse est donné uniquement sous la forme de p -valeur. On est alors enclin à penser que lorsque la p -valeur est minuscule (disons, inférieure à 1×10^{-6}), l'hypothèse nulle est "carrément fausse" ; et à l'inverse, si la p -valeur reste raisonnable, que l'hypothèse nulle est juste, ou tout au plus "légèrement fausse". Or ce serait là une grave erreur d'interprétation :

!! *Point (JL')*. Le niveau de certitude qu'on a sur la fausseté de l'hypothèse nulle n'est pas synonyme de "à quel point" l'hypothèse nulle est fausse ! \clubsuit

Commençons par le second sophisme (une p -valeur raisonnable impliquerait que l'hypothèse nulle est "presque" vraie) : nous allons réfuter celui-ci à l'aide de l'exemple suivant.

Exemple (JM'). Imaginons que je sois chargé d'établir les risque d'inondations dans différentes zones : la législation prévoyant des mesures de sécurité particulières lorsqu'il y a en moyenne au moins une inondation tous les 70 ans, je me plonge dans les archives (qui remontent, dans notre exemple, jusqu'à il y a 231 ans) pour savoir combien il y a eu d'inondations historiquement. Mon modèle est donc le suivant : $\theta_{\mathcal{J}}$ est la véritable fréquence à laquelle interviennent les inondations, qui vit dans $\Theta := \mathbb{R}_+$; X est le nombre d'inondations au cours des t dernières années avec $t := 231$ (je prends l'année pour unité de temps), qui suit la loi Poisson($\theta_{\mathcal{J}}t$) sous $\mathbb{P}_{\mathcal{J}}$. (Il semble en effet raisonnable de supposer que, à l'échelle de plusieurs années, les inondations sont des événements ponctuels et sans mémoire). Mon hypothèse nulle est $\{\theta \geq \theta_{\text{réf}}\}$ avec $\theta_{\text{réf}} = 1/70$; et dans ce cas il est clair qu'il y a un test uniformément le plus puissant, consistant à prendre pour statistique de test la valeur X elle-même, et à considérer comme suspectes (c'est-à-dire inclinant à supposer qu'il y a *peu* d'inondations) les plus petites valeurs de X . Pour finir, parmi les différentes sous-hypothèses simples composant l'hypothèse nulle, c'est évidemment

toujours dans le cas-limite $\{\theta = \theta_{\text{réf}}\}$ que la probabilité d'observer X en-dessous d'une valeur donnée sera maximale.

Cela étant établi, on calcule que la p -valeur pour observer un nombre d'inondations n_{\checkmark} sur la durée des archives est égale à $\mathbb{P}(\text{Poisson}(t\theta_{\text{réf}}) \leq n_{\checkmark}) \stackrel{\text{déf}}{=} \text{répartPoisson}(t\theta_{\text{réf}}; n_{\checkmark}+)$. Nous supposons ici qu'on a observé $n_{\checkmark} = 1$ observation sur la période des archives, ce qui donne le calcul numérique suivant pour la p -valeur :

```
> ppois(1, 231 / 70)
[1] 0.1585976.
```

Avec une p -valeur de 16 %, l'affaire semble entendue : si on a observé 1 inondation depuis le début des archives, c'est certes plutôt peu, mais absolument pas assez pour avoir une quelconque certitude sur le fait qu'il y ait moins d'une inondation tous les 70 ans... Mais est-ce à dire que, même si la période moyenne entre deux inondations est peut-être supérieure à 70 ans, elle n'est pas "trop" supérieure pour autant ? Ce n'est pas sûr du tout non plus ! Par exemple, si la fréquence réelle des inondations est d'une tous les 1 500 ans, cela donne déjà 14 % de chances d'avoir observé une inondation dans les 231 dernières années : autrement dit, quand le risque *réel* d'inondation est d'une tous les 1 500 ans, j'ai 14 % de chances d'avoir une p -valeur supérieure à 0,16 concernant l'hypothèse nulle qu'il y a (en moyenne) une inondation tous les 70 ans ou moins ! Autrement dit, il se peut très bien, sans avoir affaire à quelque coïncidence extraordinaire que ce soit, que l'hypothèse nulle soit *complètement* fautive et que pourtant on ne puisse pas la rejeter, ou seulement "mollement". ♣

Voyons maintenant le phénomène inverse. Je me demande si la saison de naissance impacte significativement l'intelligence. Pour ce faire, j'utilise les résultats des tests de l'armée (je suppose que j'ai eu le droit d'y accéder pour mes recherches), ce qui me permet d'avoir accès à un volume considérable de données : 4 millions de tests sur les 5 dernières années ! Les tests de l'armée divisent les résultats des tests en trois catégories : « intelligence faible », « intelligence moyenne » et « intelligence supérieure ». Je fais un test du χ^2 sur toutes mes données (on peut par exemple utiliser R pour ce faire, le test en que non plusstion y étant pré-implémenté), et je trouve une p -valeur de

Pearson's Chi-squared test

```
data: M
X-squared = 4452.796, df = 6, p-value < 2.2e-16.
```

Damned! La p -valeur est absolument minuscule : il semble donc que l'effet de la saison de naissance sur l'intelligence soit absolument dramatique !... Mais regardons maintenant les données de base :

Saison	Faible	Moyenne	Supérieure
Printemps	26 %	50 %	24 %
Été	24 %	52 %	24 %
Automne	24 %	50 %	26 %
Hiver	26 %	48 %	26 %

La différence entre la distribution de l'intelligence selon la saison de naissance est en fait très faible ! Simplement, comme nous avons un très grand nombre d'observations, nous pouvons voir avec une *très grande certitude* que cette différence

n'est pas *rigoureusement* nulle ! Mais cela nous fait une belle jambe, car en réalité nous n'avons pas envie de nous exciter sur une variation entre un taux de 24 % ou de 26 %... Ce dernier s'explique vraisemblablement par le fait que, les enfants faisant tous leur rentrée en septembre, ceux nés en automne commencent l'école un chouïa plus jeunes que ceux nés au printemps, et voient donc leur cerveau stimulé par l'école un peu plus tôt, ce qui leur confère un infime avantage... mais extrêmement inférieur aux fluctuations naturelles d'un individu à l'autre : si je prends un individu né en automne et un né au printemps, il y a 32,8 % de chances que celui d'automne soit dans une meilleure catégorie que celui né au printemps, contre 29,8 % de chances que ce soit l'inverse : autant dire que c'est kif-kif !

! **Définition (JN')**. Dans le vocabulaire de la statistique, on parle de *taille d'effet* pour dire à quel point la paramètre caché s'écarte de la situation de référence définie par l'hypothèse nulle. ♥

Le choix d'un indicateur de taille d'effet dépend du modèle et est par ailleurs affaire de convention. Dans le cas des inondations, la taille d'effet (qui quantifie le niveau de sécurité de notre zone) pourrait être définie comme le nombre de périodes de 70 ans survenant, en moyenne, entre deux inondations, i.e. comme $\theta_{\text{réf}} / \theta$: avec cette convention, l'indicateur de taille d'effet varie entre 0 et l'infini, notre critère étant que nous voulions que cet indicateur soit supérieur à 1. En l'espèce, notre conclusion est que, avec 1 inondation sur la période des archives, nous n'avons pas de preuve *significative* que la taille d'effet soit plus grande que 1, mais qu'il reste potentiellement parfaitement *plausible* que la taille d'effet soit *beaucoup* plus grande que 1, disons de l'ordre de 20 ! En d'autres termes, nous avons encore énormément d'incertitude sur la valeur de θ , ce qui est assez logique puisque nous avons très peu d'observations dans les archives (1 seule inondation, sachant que le nombre d'inondations est forcément entier...).

Pour l'exemple des tests d'intelligence, un indicateur de taille d'effet pertinent est de quantifier à quel point la connaissance de la valeur θ_{\checkmark} permet de deviner qui, de deux conscrits pris au hasard, et dont on ne connaît que la saison de naissance, a la meilleure catégorie d'intelligence. J'explique un peu mieux l'idée : nous avons un « devin » qui doit désigner le plus intelligent de deux conscrits pris au hasard : il marquera 1 point s'il a vu juste, -1 point s'il a eu faux, et 0 point si les deux conscrits ont le même niveau. Notre devin, ayant connaissance de la vraie valeur de θ_{\checkmark} , donnera sa réponse à chaque fois en fonction de ce qui lui donnera le meilleur résultat en espérance. Dans ces conditions, un bon indicateur de taille d'effet serait le score moyen (exprimé en pourcentage) obtenu par le devin : cet indicateur varie naturellement entre 0 % et 100 % [||]. Comme les données collectées par l'armée sont très volumineuses, nous avons que la taille d'effet est non nulle avec une significativité extrêmement forte ; en revanche, si nous essayons d'estimer cette taille d'effet, nous trouvons qu'elle n'est que de 1,1 %, ce qui est sans intérêt...

[||]. En fait, pour ce modèle, le score maximal de 100 % n'est pas mathématiquement possible — même si la corrélation entre saison et intelligence est parfaite — : car il y aura forcément des cas où les deux conscrits ont le même niveau d'intelligence, ou qu'ils sont nés à la même saison, et que le devin n'aura alors aucun indice pour le guider dans sa réponse... Plus précisément, vu les effectifs dans chaque catégorie d'intelligence, on peut montrer que la valeur maximale possible pour l'indicateur de taille d'effet vaut, en l'occurrence, 62,5 %. Mais en pratique, l'ordre de grandeur de ce maximum étant comparable à 100 %, cela ne changera essentiellement rien quand il s'agira de décider si une taille d'effet donnée doit être considérée comme digne d'intérêt ou pas ☺

Chapitre 18

Visualisation des données

18.1 Représentation des distributions de probabilité

Dans un contexte d'ingénierie et de communication, il est évidemment très important d'être capable de *représenter* graphiquement les distributions de probabilité qui nous intéressent, lorsque cela est possible, et de le faire de façon *pertinente* ! Il y a deux cas où les représentations sont tellement classiques que j'irais presque jusqu'à dire qu'elles *doivent* être utilisées lorsque cela est possible :

Définition (JO') (Diagramme en bâtons). On représente généralement les distributions de probabilités discrètes sur (un sous-ensemble de) \mathbb{R} par des *diagrammes en bâtons* : en tout point ω doté d'une masse non nulle, on trace un segment vertical basé en ω sur l'axe des abscisses et de hauteur (proportionnelle à) la masse $P(\{\omega\})$. Notez que le légendage de l'axe des ordonnées n'est pas indispensable (quoique conseillé), dans la mesure où on sait que la somme des longueurs des segments, qui représente la masse de probabilité totale, doit nécessairement valoir 1. ♡

Exemple (JP'). Soit $\Omega = \mathbb{R}_+$. La *distribution binomiale négative* de paramètres 3 et $\frac{1}{2}$, notée $\text{NégBin}(3, \frac{1}{2})$, est la distribution de probabilité P portée par $\Omega' = \mathbb{N}$ telle que, pour tout $n \in \mathbb{N}$,

$$P(\{n\}) = (n+1)(n+2)2^{-(n+4)} : \quad (\text{JQ}')$$

a ainsi $P(\{0\}) = \frac{1}{8}$, $P(\{1\}) = \frac{3}{16}$, $P(\{2\}) = \frac{3}{16}$, $P(\{3\}) = \frac{5}{32}$, ... [*].

Cette distribution de probabilité peut être représentée par le diagramme en bâtons de la distribution de la figure 18.1. ♡

Remarque (JR'). Sur la figure 18.1, j'ai à la fois gradué l'axe des ordonnées et le légendage "hors-sol" d'un fragment de longueur : c'est évidemment redondant... En pratique, choisissez l'option qui a votre préférence. ♡

Définition (JS') (Graphe de densité). On peut représenter une distribution de probabilité à densité sur \mathbb{R} simplement en traçant le graphe de la densité f . Alors la masse attribuée par cette probabilité à une partie $A \subseteq \mathbb{R}$ correspond à l'aire comprise entre la courbe et l'axe des abscisses située au-dessous de A . De même que dans le cadre discret, l'échelle n'est pas absolument nécessaire, puisqu'elle est automatiquement impliquée par le fait que l'aire sous la courbe vaille 1. ♡

[*]. On peut montrer que la somme des $P(\{n\})$ vaut bien 1 : une méthode pour ce faire consiste à évaluer en $\frac{1}{2}$ le développement en série entière de la fonction $x \mapsto \frac{1}{8}(1-x)^{-3}$.

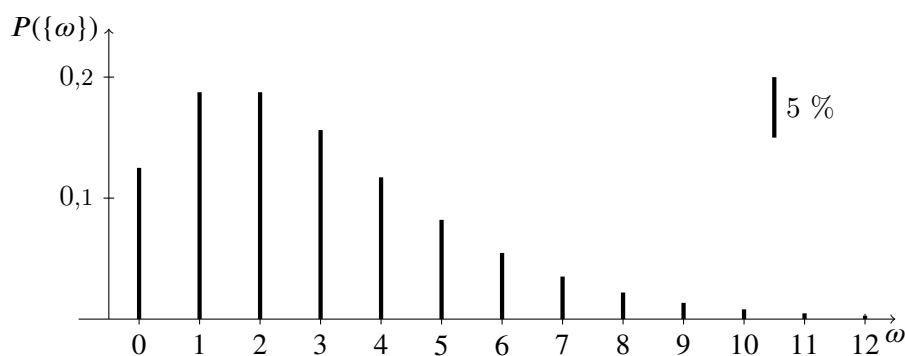


FIGURE 18.1 – Représentation par diagramme en bâtons de la distribution de probabilité $\text{NégBin}(3, \frac{1}{2})$.

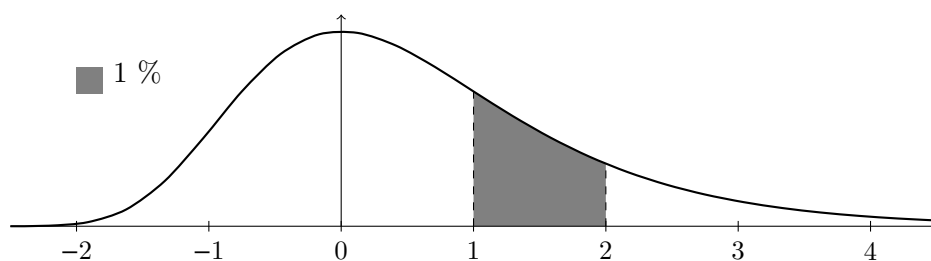


FIGURE 18.2 – Représentation graphique de la distribution de probabilité $\text{Gumbel}(0, 1)$. L'aire grisée sous la courbe correspond à la masse donnée par cette distribution au segment $[1, 2]$.

Exemple (JT'). La distribution de Gumbel standard, notée $\text{Gumbel}(0, 1)$, est la distribution de probabilité sur \mathbb{R} qui donne au voisinage infinitésimal dx de x la masse $\exp(-x - e^{-x}) \text{vol}(dx)$ [†]. On calcule alors que, par exemple, la masse donnée par cette distribution au segment $[1, 2]$ est

$$\begin{aligned} \int_{x \in [1, 2]} \exp(-x - e^{-x}) \text{vol}(dx) &= \int_{x=1}^2 \exp(-x - e^{-x}) dx \text{ [‡]} \\ &= \left[\exp(-e^{-x}) \right]_1^2 = e^{-1/e^2} - e^{-1/e} \approx 18,122 \% ; \end{aligned}$$

À l'inverse, la masse attribuée par cette distribution à un simple point comme $\{0\}$, ou à tout ensemble fini (ou même dénombrable) de points, serait nulle. La figure 18.2 donne la représentation par diagramme de densité d'une loi de Gumbel standard.

♣

Remarque (JU'). Vous noterez que je n'ai pas gradué l'axe des ordonnées en figure 18.2, mais que j'ai plutôt légendé "hors-sol" un petit bout de surface. En fait, je vous *déconseille* de graduer l'axe des abscisses dans un tel cas, car l'information sur la valeur de densité en elle-même est essentiellement sans intérêt, d'autant

[†]. On vérifie que cette densité s'intègre bien à 1 : $\int_{-\infty}^{\infty} \exp(-x - e^{-x}) dx = \left[\exp(-e^{-x}) \right]_{-\infty}^{\infty} = 1 - 0 = 1$.

[‡]. Noter que la notation ' dx ', qui se réfère à une *zone* infinitésimale dans le premier membre, se réfère à un *accroissement* infinitésimal à partir du second membre.

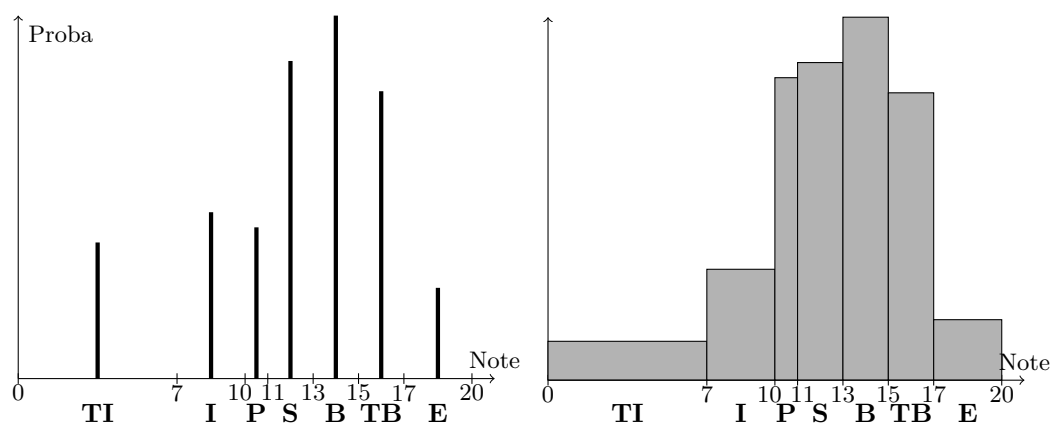


FIGURE 18.3 – Représentation de la distribution de notes de l'exemple (JW') selon qu'on opte pour une représentation en bâtons (à gauche) ou en histogramme (à droite).

qu'elle s'exprimera en général dans une unité bizarre et pas parlante du tout...! (Néanmoins en pratique, on ne rencontre que trop rarement le légendage par un petit bout de surface "hors-sol").

Remarque (JV'). Attention ; bien comprendre que les représentations des distributions de probabilité discrètes et à densité sont fondamentalement *différentes* : dans le premier cas en effet, la probabilité correspond à une *longueur*, alors que dans le second cas elle correspond à une *aire* ! C'est d'ailleurs une subtilité à prendre en compte lorsqu'on doit tracer une distribution de probabilité qu'on ne connaît que sur des classes assez grossières : soit on utilise l'approche en bâtons et alors la *hauteur* des bâtons doit correspondre à la proportion de chaque classe ; soit on utilise un histogramme et alors c'est l'*aire* qui doit être proportionnelle à l'effectif — ce qui, évidemment, donne des résultats très différents dès que les classes n'ont pas toutes la même largeur ! Voir la figure ?? à titre d'illustration.

Exemple (JW'). Supposons qu'à un examen, les mentions 'Excellent', 'Très Bien', 'Bien', 'Satisfaisant', 'Passable', 'Insuffisant' et 'Très Insuffisant' soient associées resp. aux intervalles de notes $[17, 20]$, $[15, 17[$, $[13, 15[$, $[11, 13[$, $[10, 11[$, $[7, 10[$ et $[0, 7[$, et que les proportions des notes dans les différentes mentions aient été de resp. 6 %, 19 %, 24 %, 21 %, 10 %, 11 % et 9 %. Observez comment la représentation de ces notes diffère selon qu'on les représente par un diagramme en bâtons ou par un histogramme...!

Remarque (JX'). Lorsqu'on représente une distribution à densité sur \mathbb{R} , c'est toujours l'*aire* qui doit être proportionnelle à la probabilité. Il arrive dans certains cas qu'on souhaite utiliser une échelle *logarithmique* sur l'axe des abscisses : dans ce cas, conformément au principe ci-devant, la fonction tracée ne sera pas la densité de notre distribution de probabilité elle-même, mais celle de sa mesure-image par l'application logarithme... (voir plus loin).

Cinquième partie

Modèles statistiques classiques

Annexe A

La régression linéaire

A.1 Régression linéaire simple

Définition (JY'). Le modèle de la régression linéaire simple s'appuie sur un paramètre caché traditionnellement noté (α, β, σ) , où α , β et σ sont resp. à valeurs dans \mathbb{R} , \mathbb{R} et \mathbb{R}_+^* . Il possède pour paramètres du modèle, dans sa version explicative, un entier n (que nous supposons ≥ 3) et des réels x_0, \dots, x_{n-1} . L'observation passée est notée (Y_0, \dots, Y_{n-1}) ; et suit la loi suivante pour la valeur (α, β, σ) du paramètre caché :

$$\vec{Y}_{\llbracket 0, n \llbracket} \sim \bigotimes_{i=0}^{n-1} \text{Normale}(\alpha x_i + \beta, \sigma^2) : \quad (\text{JZ}')$$

autrement dit, les Y_i sont indépendants, chaque Y_i suivant une loi normale centrée sur $\alpha x_i + \beta$ et de variance σ^2 (avec la même variance pour tous les Y_i). \heartsuit

Définition (KA'). Pour traiter le modèle de la régression linéaire simple, on a coutume de s'appuyer sur les paramètres et statistiques suivants :

- \bar{x} est la moyenne empirique des x_i . (On peut donc le voir comme un paramètre du modèle, s'exprimant à partir des paramètres de base);
- s_x est l'écart-type empirique des x_i ^[*]. (Là aussi, on peut le voir comme un paramètre du modèle dérivé des paramètres de base);
- \bar{Y} est la moyenne empirique des observations passées Y_i .
- $\hat{\alpha}$, qu'on peut appeler « pente de régression estimée », est la statistique définie par

$$\hat{\alpha} := \frac{\text{cov}_{\text{emp}}((x_i, Y_i)_{i \in \llbracket 0, n \llbracket})}{s_x^2}; \quad (\text{KB}')$$

- R est le coefficient de corrélation empirique entre les x_i et les Y_i :

$$R := \frac{\text{cov}_{\text{emp}}((x_i, Y_i)_{i \in \llbracket 0, n \llbracket})}{s_x \text{var}_{\text{emp}}^{1/2}(Y_i)_{i \in \llbracket 0, n \llbracket}}. \quad (\text{KC}')$$

- On appelle « droite de régression estimée » la droite (aléatoire, fonction de l'observation passée) de pente $\hat{\alpha}$ et passant par le point de coordonnées (\bar{x}, \bar{Y}) ;
- On note $\hat{\beta}$ l'ordonnée à l'origine de la droite de régression : $\hat{\beta} := \bar{Y} - \hat{\alpha}\bar{x}$.

[*]. Dans de nombreux ouvrages de référence, aucune notation spécifique n'est introduite pour s_x : on écrit alors simplement « $\sum_i (x_i - \bar{x})^2$ » pour ce que nous noterons ns_x^2 .

- Pour tout $i \in \llbracket 0, n \rrbracket$, on définit l'« erreur estimée » comme l'écart entre Y_i et l'ordonnée de la droite de régression à l'abscisse correspondante :

$$\hat{E}_i := Y_i - \hat{\alpha}x_i - \hat{\beta}. \quad (\text{KD}')$$

- On note

$$\hat{\sigma}_{\text{emp}} := \left(\sum_{i=0}^{n-1} \hat{E}_i^2 / n \right)^{1/2};$$

$$\hat{\sigma}_{\text{B}} := \left(\sum_{i=0}^{n-1} \hat{E}_i^2 / (n-2) \right)^{1/2}.$$

♡

Estimation de la pente de la droite de régression

Procédure (KE'). Dans ce contexte, l'estimateur standard pour α est $\hat{\alpha}$. Il correspond à la fois à l'estimateur des moindres carrés et à l'estimateur du maximum de vraisemblance. ♡

Remarque (KF'). Plus précisément, l'estimateur des moindres carrés pour le paramètre caché dans son ensemble est donné par le triplet $(\hat{\alpha}, \hat{\beta}, \bullet)$, où l'estimateur pour σ n'est pas précisé dans la mesure où tous les estimateurs possibles pour ce paramètre conduisent aux mêmes résultats en ce qui concerne les moindres carrés. En l'occurrence, appliquer l'estimation par moindres carrés consiste à trouver $\hat{\alpha}_{\checkmark}$ et $\hat{\beta}_{\checkmark}$ comme les valeurs $\hat{\alpha}$ et $\hat{\beta}$ minimisant

$$\sum_{i=0}^{n-1} (y_{i\checkmark} - (\hat{\alpha}x_i + \hat{\beta}))^2. \quad (\text{KG}')$$

♣

Remarque (KH'). De même, l'estimateur du maximum de vraisemblance pour α provient (comme toujours) de l'estimateur du maximum pour le paramètre caché dans son ensemble : celui-ci est donné plus précisément par le triplet $(\hat{\alpha}, \hat{\beta}, \hat{\sigma}_{\text{emp}})$. ♣

Théorème (KI'). Dans le contexte de la régression linéaire, pour tester une hypothèse de la forme $\{\alpha = \alpha_0\}$, $\{\alpha \leq \alpha_0\}$ ou $\{\alpha \geq \alpha_0\}$, ou pour déterminer un intervalle de confiance pour α , on utilise la propriété suivante :

$$\text{Loi}_{\checkmark} \left(\frac{\hat{\alpha} - \alpha_{\checkmark}}{\hat{\sigma}_{\text{emp}} / \sqrt{n-2} s_x} \right) = T_{\text{St}}(n-2). \quad (\text{KJ}')$$

◇

Remarque (KK'). La statistique qui apparaît au dénominateur de la variable aléatoire ci-dessus est appelée *erreur standard* du paramètre α . Ce nom provient de ce que cette quantité mesure l'ordre de grandeur typique de la différence entre $\hat{\alpha}$ et α , d'une façon que le quotient suive une certaine loi de Student. En fait, le carré de l'erreur standard correspond à un estimateur sans biais de la variance (sous la véritable loi) de l'estimateur $\hat{\alpha}$.

On notera que l'erreur standard peut aussi s'écrire sous la forme

$$\hat{\sigma}_{\text{B}} / \sqrt{n} s_x. \quad (\text{KL}')$$

♣

Remarque (KM'). Le test de l'hypothèse nulle $\{\alpha = 0\}$ est particulièrement important, car le fait que α_{\checkmark} soit non nul, resp. nul, signifie que la valeur de l'abscisse x_i a un impact, resp. pas d'impact, sur la valeur de l'observation Y_i . En fait, α_{\checkmark} quantifie à quel point est-ce qu'une augmentation de x a tendance (sous la vraie loi) à entraîner une augmentation de Y . \clubsuit

Estimation de l'intensité du bruit

Théorème (KN'). On a la relation suivant, qui permet de simplifier le calcul de $\hat{\sigma}$:

$$n\hat{\sigma}^2 = n \operatorname{var}_{\text{emp}}(y_i) - \hat{\alpha}^2 s_x^2. \quad (\text{KO}')$$

◇

Procédure (KP'). Dans ce contexte, l'estimation standard pour σ est $\hat{\sigma}_B$; parfois aussi on rencontre l'estimateur $\hat{\sigma}_{\text{emp}}$. L'estimateur $\hat{\sigma}_{\text{emp}}$ correspond au maximum de vraisemblance, tandis que $\hat{\sigma}_B$ en est une variante qui présente la propriété que $\hat{\sigma}_B^2$ est un estimateur sans biais de σ^2 (mais en tant qu'estimateur de σ lui-même, $\hat{\sigma}_B$ est en revanche biaisé vers le haut). \heartsuit

Théorème (KQ'). Dans le contexte de la régression linéaire, pour tester une hypothèse de la forme $\{\sigma = \sigma_0\}$, $\{\sigma \leq \sigma_0\}$ ou $\{\sigma \geq \sigma_0\}$, ou pour déterminer un intervalle de confiance pour σ , on utilise la propriété suivante :

$$\text{Loi}_{\checkmark} \left(\frac{n\hat{\sigma}_{\text{emp}}^2}{\sigma_{\checkmark}^2} \right) = \chi^2(n-2). \quad (\text{KR}')$$

◇

Remarque (KS'). Le numérateur de la variable aléatoire ci-dessus peut aussi se ré-écrire comme $(n-2)\hat{\sigma}_B^2$. \clubsuit

Estimation de l'ordonnée de la droite de régression en une abscisse donnée La valeur β_{\checkmark} correspond à l'ordonnée à l'origine de la véritable droite « $y = \alpha_{\checkmark}x + \beta$ » : pour cette raison, les logiciels anglophones l'appellent souvent « intercept » (car c'est la hauteur à laquelle la droite de régression *intercepte* l'axe des ordonnées). Cependant, en général dans les contextes industriels considérés, la valeur zéro de l'axe des abscisses ne pas de pertinence particulière : de manière générale, ce qu'on cherchera en général est valeur de la quantité d'intérêt $\alpha x_{\star} + \beta$ pour une certaine abscisse (connue) x_{\star} qui pourra très bien ne pas être nulle. Il se trouve en fait que le problème n'est pas plus compliqué à résoudre pour une valeur générale de x_{\star} que pour une valeur quelconque : c'est donc au cas général que nous nous intéresserons ici.

Procédure (KT'). Dans le contexte de la régression linéaire, l'estimateur standard pour $\alpha x_{\star} + \beta$ est $\hat{\alpha}x_{\star} + \hat{\beta}$: autrement dit, c'est l'estimateur qu'on obtient par substitution à partir des moindres carrés ou du maximum de vraisemblance. \heartsuit

Théorème (KU'). Dans le contexte de la régression linéaire, pour tester une hypothèse de la forme $\{\alpha x_{\star} + \beta = y_{\text{réf}}\}$, $\{\alpha x_{\star} + \beta \leq y_{\text{réf}}\}$ ou $\{\alpha x_{\star} + \beta \geq y_{\text{réf}}\}$, ou pour

déterminer un intervalle de confiance pour $\alpha x_\star + \beta = y_{\text{réf}}$, on utilise la propriété suivante :

$$\text{Loi}_{\checkmark} \left(\frac{(\hat{\alpha}x_\star + \hat{\beta}) - (\alpha_{\checkmark}x_\star + \beta_{\checkmark})}{\hat{\sigma}_{\text{emp}} / \sqrt{n-2} \times (1 + ((x_\star - \bar{x}) / s_x)^2)^{1/2}} \right) = T_{\text{St}}(n-2). \quad (\text{KV}')$$

◇

Remarque (KW'). On observe donc que l'erreur standard sur $\alpha x_\star + \beta$ dépend assez sensiblement de la valeur de x_\star : elle est minimal lorsque x_\star est "au milieu" des valeurs des x_i , et commence à croître sensiblement dès lors que x_\star s'écarte de la plage des valeurs "typiques" de x_i . Informellement, on peut donc dire que la régression linéaire arrive à déterminer plus précisément l'emplacement de la droite de régression pour les abscisses correspondant au milieu du nuage de points que pour les abscisses correspondant à la périphérie du nuage de points. ♣

Remarque (KX'). Il est très important, lorsqu'on doit estimer un intervalle de confiance pour $\alpha x_\star + \beta$, de ne pas commencer par déterminer un intervalle de confiance sur β avant de « reporter » celui-ci dans la formule « $\alpha x_\star + \beta$ » : outre que cette stratégie de construction d'intervalle de confiance en deux temps est presque toujours une stratégie donnant des résultats imprécis de manière générale, ici cela peut être particulièrement catastrophique si les abscisses du nuage de point sont toutes éloignées de zéro, car dans ce cas on aura une très forte incertitude sur β , bien plus forte que l'incertitude qu'on peut obtenir sur $\alpha x_\star + \beta$ en procédant convenablement... ♣

Intervalle de prédiction

Définition (KY'). On peut introduire une version prédictive du modèle de la régression linéaire : dans ce cas, on introduit un paramètre supplémentaire x_n dans le modèle ; et l'observation future, notée Y_n , suit la loi Normale($\alpha x_n + \beta, \sigma^2$), indépendamment des autres Y_i : il s'agit donc en substance du même modèle, mais avec une observation de plus qui n'est observée que dans le futur. ♥

Remarque (KZ'). Dans ce cas, les valeurs de \bar{x} , s_x et des statistiques \bar{Y} , &c. sont toujours déterminées à partir des indices i relatifs aux seules observations *passées*. ♣

Procédure (LA'). Le prédicteur standard pour Y_n est alors $\hat{\alpha}x_n + \hat{\beta}$. ♥

Théorème (LB'). Dans ce contexte, pour déterminer un intervalle de confiance pour Y_n , on utilise la propriété suivante :

$$\text{Loi}_{\checkmark} \left(\frac{(\hat{\alpha}x_n + \hat{\beta}) - Y_n}{\hat{\sigma}_{\text{emp}} / \sqrt{n-2} \times (n+1 + ((x_n - \bar{x}) / s_x)^2)^{1/2}} \right) = T_{\text{St}}(n-2). \quad (\text{LC}')$$

◇

Remarque (LD'). Pour peu que n ne soit pas trop petit, cette fois-ci l'erreur standard sur Y_n ne dépendra pas beaucoup de la valeur x_n , car la somme $n+1 + ((x_n - \bar{x}) / s_x)^2$ sera toujours proche de n en pratique. ♣

Notes bibliographiques

Jusqu'en 2016, le cours de statistique de Mines Nancy était dispensé par Thierry VERDEL : son polycopié [6] (qui s'appuyait lui-même sur un document plus ancien de Claude CHAMBON) est toujours disponible en ligne, et constitue une approche tout à fait recommandable à la statistique pour l'ingénieur, en particulier pour découvrir les techniques statistiques particulières les plus utilisées (tests et intervalles de confiance de Student, de Fisher, du chi-deux ; régression linéaire ; analyse de la variance ; etc.), qui ne sont pas exposées dans le présent document. Cependant, l'approche est très différente entre mon polycopié (qui s'intéresse surtout aux *concepts* statistiques) et celui de M. Verdel (qui s'intéresse surtout aux *procédures*) ; c'est pourquoi je ne pense pas qu'il soit avisé de lire ces deux documents en parallèle : mais plutôt, de garder en tête l'existence du polycopié de M. Verdel pour le jour où vous aurez besoin d'en savoir davantage sur une méthode particulière !

Pour tout ce qui concerne la partie fréquentiste de ce cours, je me suis appuyé d'assez près sur deux excellents ouvrages que je vous recommande (et que vous trouverez à la bibliothèque de l'École) : celui de CADRE & VIAL [1], particulièrement léger, idéal pour une introduction aux concepts de la statistique fréquentiste ; et celui de STOLTZ & RIVOIRARD [4], plus fouillé, qui présente un grand nombre d'exemples d'applications détaillés, avec des implémentations informatiques dont les codes sont téléchargeables librement, parfait pour voir comment ces concepts se mettent en action ! Ces deux ouvrages doivent eux-mêmes beaucoup à l'école de statistique mathématique française, et notamment M^{me} Dominique PICARD, dont l'auteur de ces lignes a lui-même eu la chance de suivre le cours [3] !

Pour la partie bayésienne, malheureusement, il ne semble pas exister d'ouvrage accessible (en tout cas pas en français) recoupant la deuxième partie de ce polycopié... Peut-être parce que ce qui y est dit est "trop facile" pour constituer l'objet d'un livre entier ? Peut-être à cause de la prééminence de la statistique fréquentiste au cours du XX^e siècle?... Un ouvrage très célèbre est celui de Christian ROBERT [5] ; mais il part tout de suite d'un niveau avancé, et sera utile seulement pour celui souhaitant se spécialiser en statistique bayésienne. En fait, la statistique bayésienne étant souvent décriée comme une approche trop peu rigoureuse (mais ces critiques tendent à s'estomper, notamment à cause de l'augmentation de la puissance de calcul qui rend l'approche bayésienne particulièrement performante dans des cas où l'analyse purement mathématique rendrait les armes, en particulier en intelligence artificielle), on trouve plus facilement des ouvrages expliquant *pourquoi* la statistique bayésienne est une bonne idée, plutôt que pour détailler *comment* s'en servir ! Dans cette catégorie, outre le livre [5] déjà cité de Robert, je recommande à tous ceux qui sont intéressés par la découverte de l'extraordinaire puissance de la méthode bayésienne de lire le tout récent livre du célèbre vulgarisateur Lê Nguyễn

HOANG [2], qui mêle mathématique, informatique et philosophie avec un enthousiasme passionnant !

Bibliographie

- [1] Benoît CADRE et Céline VIAL : *Statistique mathématique — Cours et exercices corrigés*. Ellipses, 2012. ISBN 978-2-7298-7323-3.
- [2] Lê Nguyễn HOANG : *La formule du savoir — Une philosophie unifiée du savoir fondée sur le théorème de Bayes*. Edp Sciences, 2018. ISBN 978-2-7598-2260-7.
- [3] Dominique PICARD : Statistique et modèles aléatoires.
<https://www.lpsm.paris/pageperso/picard/pubps/Poly0304.pdf>,
2004.
- [4] Vincent RIVOIRARD et Gilles STOLTZ : *Statistique mathématique en action*. Vuibert, 2^e édition, 2012. ISBN 978-2-311-00720-6.
- [5] Christian P. ROBERT : *Le choix bayésien — Principes et pratique*. Springer, 2006. ISBN 978-2-287-25173-3.
- [6] Thierry VERDEL : Décision et prévision statistiques.
<https://drive.google.com/open?id=1SAhHtaC8tvXhW03v824kuR0ajhoP8HdT>,
2016.