Aggregated cross-validation

Guillaume Maillard

Université Paris Sud

October 23, 2017

Let \mathcal{X} be a measure space.

- A sample is a finite collection $(x_i, y_i)_{1 \le i \le k}$, where $y_i \in \mathcal{Y} = \{0, 1\}$ and $x_i \in \mathcal{X}$.
- A *classifier* is a measurable function $f : \mathcal{X} \to \mathcal{Y}$
- A *learning rule* takes a sample D as input and produces a classifier $G(D) : \mathcal{X} \to \mathcal{Y}$.
- Let (X, Y) ∈ X × Y be a r.v. The excess risk of a classifier f, denoted ℓ(f*, f), equals

$$\mathbb{P}(f(X) \neq Y) - \inf_{g ext{ classifieur}} \mathbb{P}(g(X) \neq Y)$$

Cross-Validation uses data to estimate the risk of a learning rule.

Definition

$$D_n$$
 a sample s.t $|D_n| = n$, $T \subset \{1, ..., n\}$.

$$HO_{T}(G) = \frac{1}{|cT|} |\{j \in C T : G((x_{i}, y_{i})_{i \in T})(x_{j}) \neq y_{j}\}|$$

$$1 \leq p \leq n, \ \mathcal{T} \subset \{\mathcal{T} \subset \{1, ..., n\} : |\mathcal{T}| = p\}$$

$$CV_{\mathcal{T}}(G) = \frac{1}{|\mathcal{T}|} \sum_{\mathcal{T} \in \mathcal{T}} HO_{\mathcal{T}}(G)$$

Let \mathcal{G} be a set of learning rules.

- The hold out classifier is $\hat{f}_T^{\text{ho}} = \hat{G}_T^{ho}(D_n^T)$ where $\hat{G}_T^{ho} = \operatorname{argmin}_{G \in \mathcal{G}} \operatorname{HO}_T(G)$.
- A cross-validated classifier is defined as $\hat{f}_{\mathcal{T}}^{CV} = \hat{G}_{\mathcal{T}}^{CV}(D_n)$ where $\hat{G}_{\mathcal{T}}^{CV} = \operatorname{argmin}_{G \in \mathcal{G}} CV_{\mathcal{T}}(G)$

Aggregation and ensemble methods

Many classifiers \Rightarrow one prediction

Definition

Let $(f_i)_{i=1..V}$ be classifiers.

$$\operatorname{maj}(f_1,\ldots,f_V) = x \to \operatorname{argmax}_{y \in \{0;1\}} |\{i: f_i(x) = y\}|$$

Bagging: $D_n, G, \mathcal{T} \xrightarrow{\text{train}} (G(D_n^{\mathcal{T}}))_{\mathcal{T} \in \mathcal{T}} \xrightarrow{\text{maj}} f_{\mathcal{T}}^{\text{bag}}$

- ₹ 🖬 🕨

The idea is to aggregate several hold-out classifiers.

- Base method: hold out (not retrained).
- The training set varies
- Compared with CV: selection and averaging switched.
- Compared with bagging: hold out computed from the whole sample.

$$T_1, ..., T_V$$
 i.i.d ~ $\mathcal{U}(\{T \subset \{1, .., n\} : |T| = n - p\}).$

$$\widehat{f}_{V}^{\mathrm{ag}} = \mathsf{maj}\left((\widehat{f}_{\mathcal{T}_{i}}^{\mathrm{ho}})_{1 \leq i \leq V}\right)$$

Parameters are V and
$$au = \frac{n-p}{n}$$

V-fold CV \implies V-fold aggregation.

 $CV_{\mathcal{T}} \Longrightarrow \widehat{f}_{\mathcal{T}}^{\mathrm{ag}}.$

□ ▶ ▲ 臣 ▶ ▲

Comparison with cross-validation



Agghoo:

$$\begin{array}{c} (\operatorname{HO}_{T_{1}}(G_{m}))_{m \in \mathcal{M}} \xrightarrow{\operatorname{argmin}} (\widehat{m}_{1}) \xrightarrow{\operatorname{train}} (G_{\widehat{m}_{1}}(D_{n}^{T_{1}})) \\ \vdots & \vdots & \vdots \\ (\operatorname{HO}_{T_{i}}(G_{m}))_{m \in \mathcal{M}} \xrightarrow{\operatorname{argmin}} (\widehat{m}_{i}) \xrightarrow{\operatorname{train}} (G_{\widehat{m}_{i}}(D_{n}^{T_{i}})) \xrightarrow{\operatorname{maj}} (\widehat{f}_{\mathcal{T}}^{\operatorname{ag}}) \\ \vdots & \vdots & \vdots \\ (\operatorname{HO}_{T_{V}}(G_{m}))_{m \in \mathcal{M}} \xrightarrow{\operatorname{argmin}} (\widehat{m}_{V}) \xrightarrow{\operatorname{train}} (G_{\widehat{m}_{V}}(D_{n}^{T_{V}})) \end{array}$$



Benchmark:

 $\inf_{G\in\mathcal{G}}\ell(f_*,G(D_n))$

called *oracle*. Hold out: oracle on D_n^T :

 $\inf_{G\in\mathcal{G}}\ell(f_*,G(D_n^T))$

oracle inequality = bounding the risk by an affine function of the oracle.

Theorem

Let
$$\eta(x) = \mathbb{P}[Y = 1 | X = x]$$
 Suppose that

$$orall h \in [0;1], \mathbb{P}[|2\eta(X)-1| \leq h] \leq ch^eta$$

(margin hypothesis). Then for D_n i.i.d ~ (X, Y),

$$\mathbb{E}[\ell(f_*, \widehat{f}_T^{\text{ho}})] \le 1.5 \mathbb{E}[\inf_{G \in \mathcal{G}} \ell(f_*, G(D_n^T))] + \frac{14.5c^{\frac{1}{\beta+2}} \log(e|\mathcal{G}|)}{p^{\frac{\beta+1}{\beta+2}}}$$

where n - p = |T|

伺 ト く ヨ ト く ヨ ト

3

Proposition

Let $(f_i)_{i=1...V}$ be classifiers, and $f^{maj} = maj(f_1, ..., f_V)$. We obtain

$$\ell(f_*, f^{maj}) \leq \frac{2}{V} \sum_{i=1}^V \ell(f_*, f_i)$$

In particular if the classifiers f_i are equal in distribution,

 $\mathbb{E}[\ell(f_*, f^{maj})] \le 2\mathbb{E}[\ell(f_*, f_1)]$

An oracle inequality for \hat{f}_V^{ag} :

$$\mathbb{E}[\ell(f_*, \widehat{f}_V^{\mathrm{ag}})] \leq 3\mathbb{E}[\inf_{G \in \mathcal{G}} \ell(f_*, G(D_n^T))] + \frac{29c^{\frac{1}{\beta+2}}\log(e|\mathcal{G}|)}{p^{\frac{\beta+1}{\beta+2}}}$$

Relevance:

- The oracle is *adaptative*.
- It is *simultaneously minimax* for well chosen *G*.

• So is the hold-out if
$$p^{\frac{\beta+1}{\beta+2}} = o(\text{oracle})$$
.

Let
$$\mathcal{P}_{\gamma,L,\beta,c}$$
 be the set of r.v (X, Y) such that:
• $\eta(x) = \mathbb{P}[Y = 1 | X = x]$ is γ -Hölder with constant L .
• $\forall h \in [0; 1], \mathbb{P}(|2\eta(X) - 1| \le h) \le ch^{\beta}$
• $\operatorname{supp}(X) \subset [0; 1]^{d}$.

The following result is due to Audibert et Tsybakov:

Theorem

The minimax risk over
$$\mathcal{P}_{\gamma,L,\beta,c}$$
 is of order $n^{-\frac{\gamma(1+\beta)}{\gamma(2+\beta)+d}}$.

- No such guarantee for cross-validation.
- Optimal up to a constant.
- Importance of aggregation not shown
- Why not hold-out?
- Choice of parameters?

- monte carlo aggregated cross validation is compared to monte carlo cross-validation and the oracle.
- $\mathcal{G} = k$ -NN for $1 \le k \le n$, k odd.
- i.i.d sample, n = 500, with distribution

$$X \sim \mathcal{U}([0;1]^2)$$
$$\mathbb{P}(Y = 1|X) = \sigma\left(\frac{g(X) - b}{\lambda}\right) \text{ where } \sigma(u) = \frac{1}{1 + e^{-u}},$$
$$g(u, v) = e^{-(u^2 + v)^3} + u^2 + v^2$$

Picture of the distribution



The Bayes risk is approximately 0.2418

Excess risks as a function of τ



Guillaume Maillard Aggregated cross-validation

An example in regression

Least squares, \mathcal{G} contains regressograms with step $h \in \{\frac{1}{k} : 1 \le k \le n\}$

- $X \sim U([0; 1])$
- $Y \sim t(X) + \varepsilon$ where $\varepsilon \sim \mathcal{U}([-2\sqrt{3}, 2\sqrt{3}])$ indt of X

•
$$t(x) = 10(1 - 2|x - 0.5|)$$



Performance in least-squares regression

Performance in least squares regression



Guillaume Maillard Agg

Aggregated cross-validation

э

- Aggregation of hold-out estimators.
- An alternative to cross-validation in classification and regression.
- Satisfies an oracle inequality in classification
- Aggregation useful in practice, not in theory
- Can beat the oracle!
- More theory to understand aggregation