

# Agrégation d'hold-out

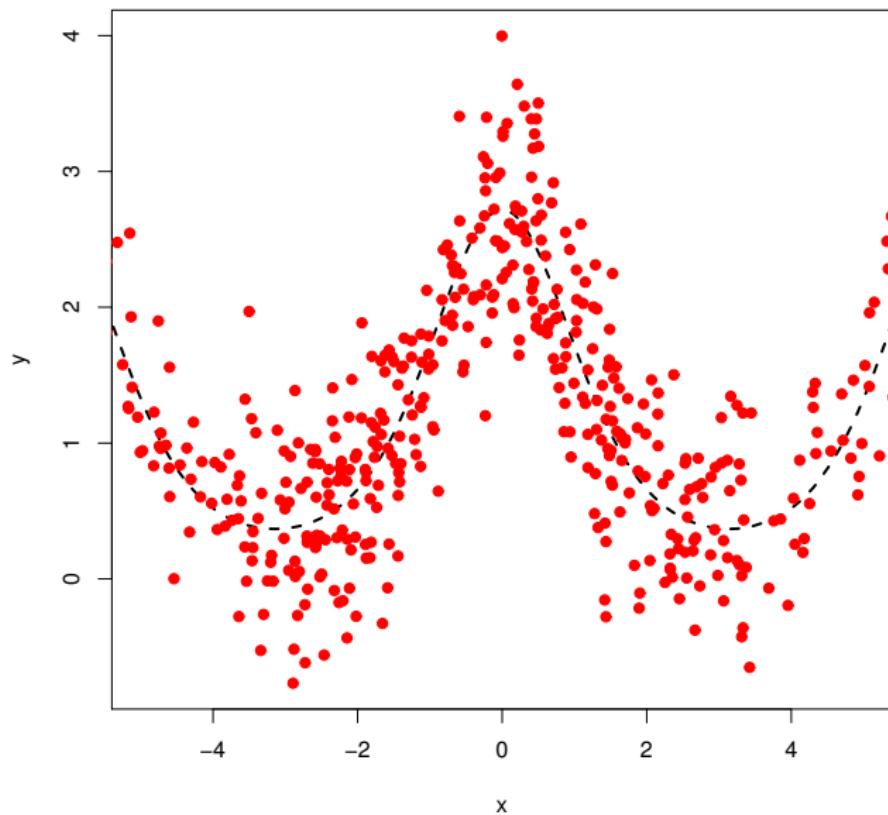
Guillaume Maillard, M.Lerasle, S.Arlot

Université du Luxembourg

May 25, 2021

- Presentation of Agghoo
- Theoretical performance
- Simulations in regression
- Classification

# 1. Presentation of Agghoo: Regression



# 1. Presentation of Agghoo - prediction

$\mathcal{X}, \mathcal{Y}$  are measurable spaces,  $P$  a distribution on  $\mathcal{X} \times \mathcal{Y}$ , and  $\gamma : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  a measurable *contrast function*. For ex:

- $\mathcal{Y} = \{-1; 1\}$  and  $\gamma(u, y) = \mathbb{I}_{u \neq y}$  in classification
- $\mathcal{Y} = \mathbb{R}$  and  $\gamma(u, y) = |y - u|$  in median regression.

## Definition

- A *predictor* is a measurable function  $t : \mathcal{X} \rightarrow \mathcal{Y}$ .
- The *excess risk* of a predictor  $f$  is defined by

$$\ell(s, t) = \mathbb{E} [\gamma(f(X), Y)] - \inf_{g \text{ predictor}} \mathbb{E} [\gamma(g(X), Y)]$$

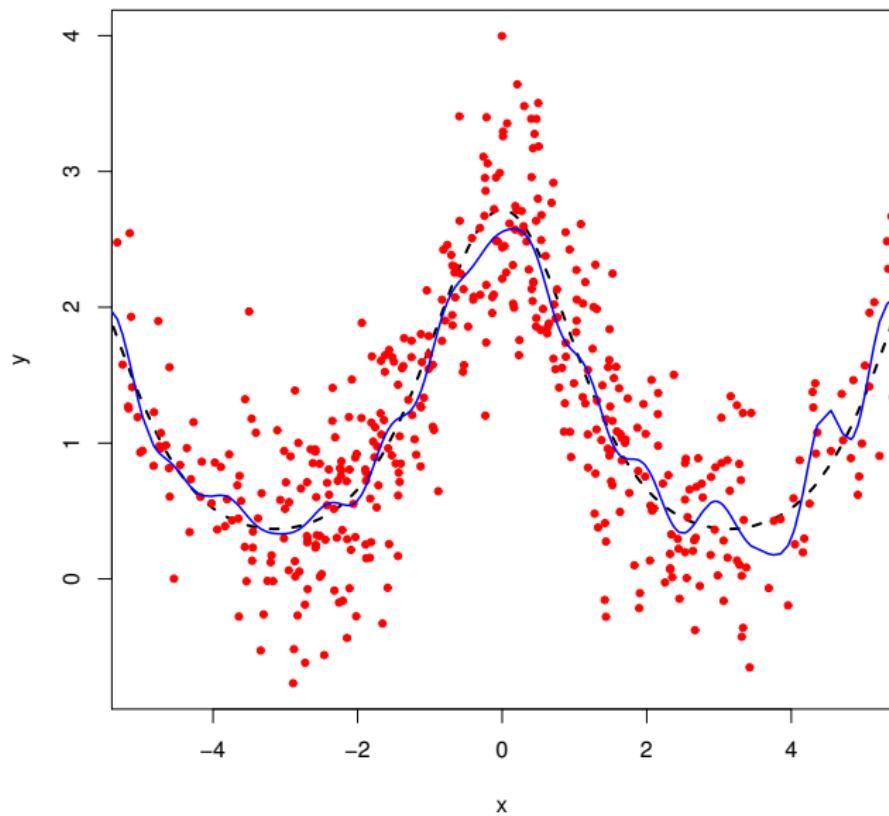
# 1.Presentation of Agghoo: notations

Fix an integer  $n$ . Denote

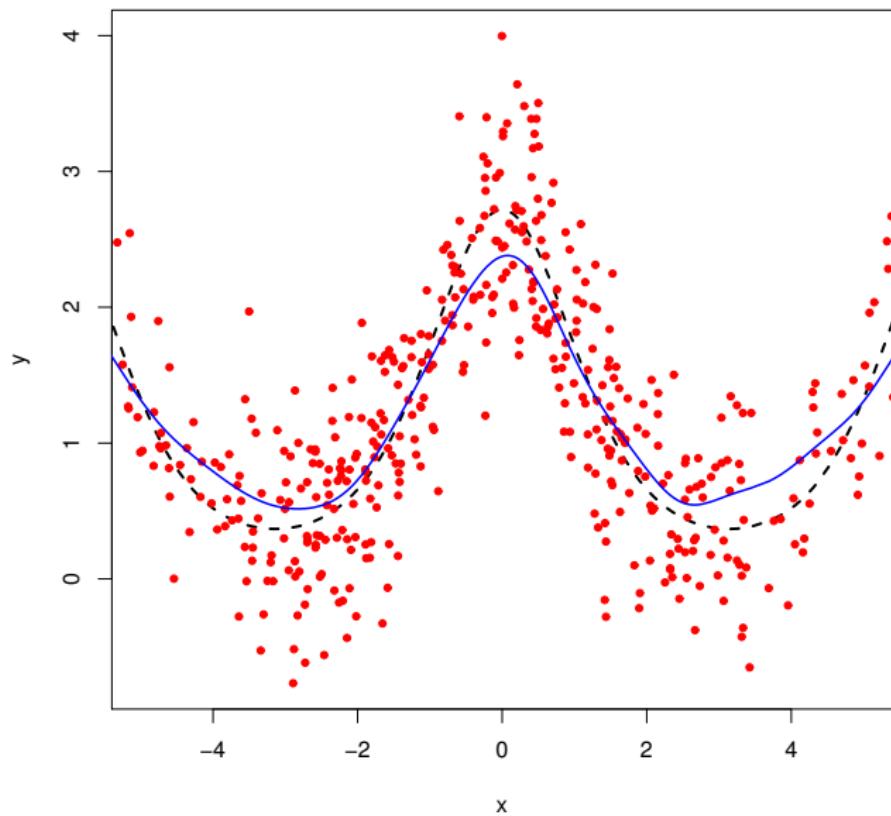
## Definition

- $D_n = (X_i, Y_i)_{1 \leq i \leq n} \sim P^{\otimes n}$  a *sample* of size  $n$ .
- For a sample  $D_n$  and  $T \subset \{1, \dots, n\}$ ,  $D_n^T = (X_j, Y_j)_{j \in T}$ .
- $\mathcal{A} : \cup_{k \in \mathbb{N}} (\mathcal{X} \times \mathcal{Y})^k \rightarrow \text{predictors}$  a *learning rule*.

# 1. Presentation of Agghoo: example estimator



# 1. Presentation of Agghoo: example estimator



# 1. Presentation of Agghoo: hyperparameter selection

- Given learning rules  $(\mathcal{A}_m)_{m \in \mathcal{M}}$  and a sample  $D_n$ , choose parameter  $m$ .
- Optimal choice  $m_*$  realizes  $\inf_{m \in \mathcal{M}} \ell(s, \mathcal{A}_m(D_n))$ .

Examples:

- $m = k$  in  $k$ -NN for classification.
- regularization parameter,  $m = \lambda$  (Lasso, Ridge, SVM...)

# 1. Presentation of Agghoo: CV risk estimation

## Definition

Fix a sample  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$ . Let  $T \subset \{1, \dots, n\}$ . For any learning rule  $\mathcal{A}$ ,

$$\text{HO}_T(\mathcal{A}) = \frac{1}{|T^c|} \sum_{j \notin T} \gamma(\mathcal{A}(D_n^T)(X_j), Y_j).$$

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

$$CV_T(\mathcal{A}) = \frac{1}{|\mathcal{T}|} \sum_{T \in \mathcal{T}} \text{HO}_T(\mathcal{A})$$

# 1. Presentation of Agghoo: CV selection

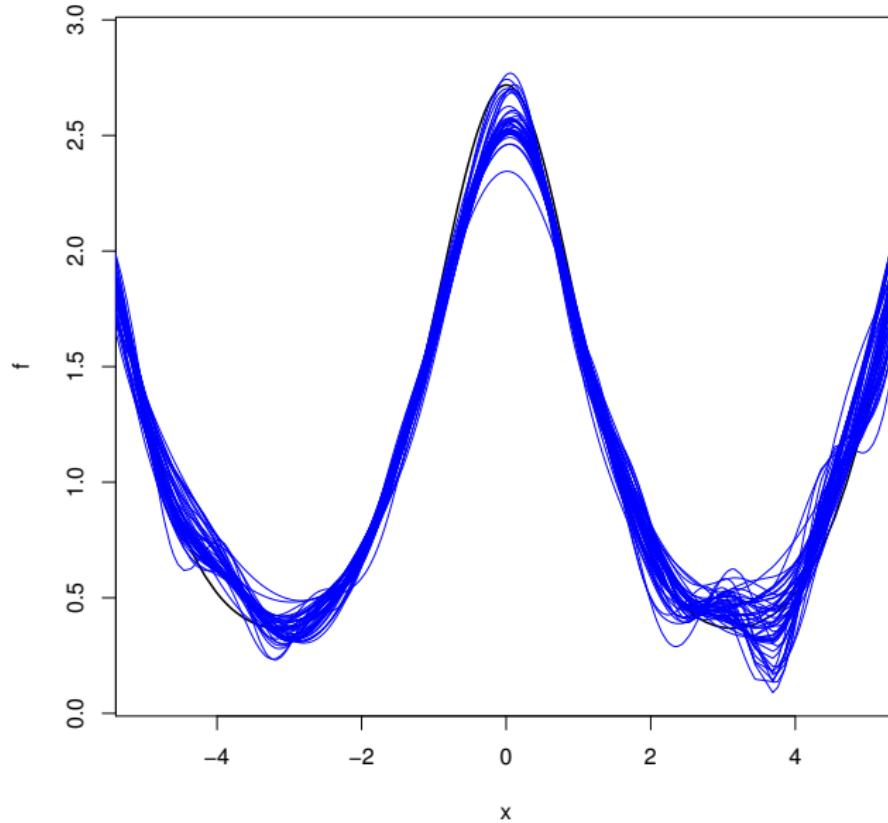
Idea: estimate  $\ell(s, \mathcal{A}_m(D))$  and optimize in  $m$ .

## Definition

Let  $\mathcal{M}$  be a set of hyperparameters.

- A *hold out predictor* is  $\hat{f}_T^{\text{ho}} = \mathcal{A}_{\hat{m}_T^{\text{ho}}}(D_n^T)$  where  $\hat{m}_T^{\text{ho}} = \operatorname{argmin}_{m \in \mathcal{M}} \text{HO}_T(\mathcal{A}_m)$ .
- A *CV predictor* is defined by  $\hat{f}_T^{CV} = \mathcal{A}_{\hat{m}_T^{CV}}(D_n)$  where  $\hat{m}_T^{CV} = \operatorname{argmin}_{m \in \mathcal{M}} \text{CV}_T(\mathcal{A}_m)$

# 1. Presentation of Agghoo: variability of the hold out



# 1. Presentation of Agghoo: Agghoo

Idea: aggregate several hold out estimators.

## Definition

Assume  $\mathcal{Y}$  is convex (regression...).

$T_1, \dots, T_V$  i.i.d  $\sim \mathcal{U}(\{T \subset \{1, \dots, n\} : |T| = \lfloor \tau n \rfloor\})$ .

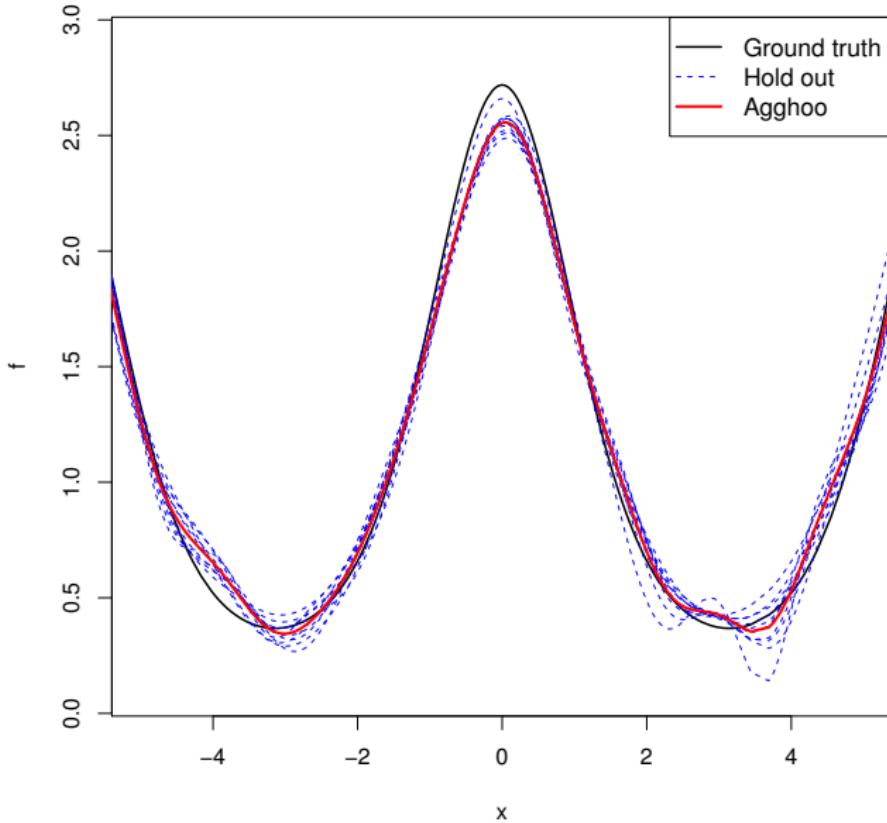
$$\hat{f}_{\tau, V}^{\text{ag}} = \frac{1}{V} \sum_{i=1}^V \hat{f}_{T_i}^{\text{ho}}$$

Parameters are  $V$  and  $\tau$ .

$V$ -fold CV  $\implies$   $V$ -fold aggregation.

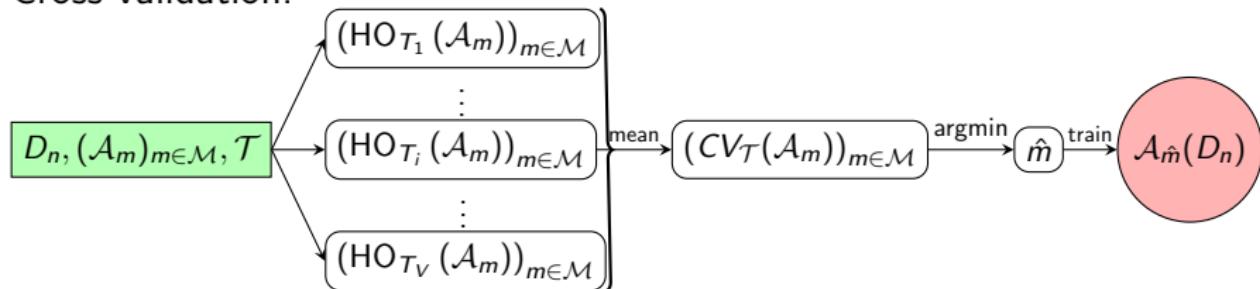
$$CV_T \implies \hat{f}_T^{\text{ag}}$$

# 1. Presentation of Agghoo: an example

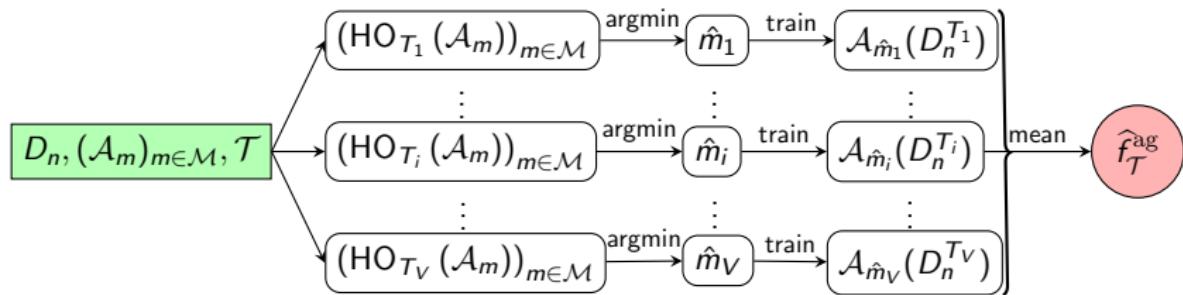


# 1. Presentation of Agghoo: comparison with cross-validation

Cross-validation:



Agghoo:



## 2.Theoretical performance: comparison with hold out

Agghoo is safe:

- $\widehat{f}_{\tau,1}^{\text{ag}} = \widehat{f}_{T_1}^{\text{ho}}$  is just the hold out.
- $\mathbb{E} \left[ \ell(s, \widehat{f}_{\tau,V}^{\text{ag}}) \right] \leq \mathbb{E} \left[ \ell(s, \widehat{f}_{T_1}^{\text{ho}}) \right]$  if  $\gamma$  is convex.
- For any  $k \geq 1$ ,  $\mathbb{E} \left[ \ell(s, \widehat{f}_{\tau,V}^{\text{ag}})^k \right] \leq \mathbb{E} \left[ \ell(s, \widehat{f}_{T_1}^{\text{ho}})^k \right]$  if  $\gamma$  is convex.

Oracle inequalities for the hold-out imply oracle inequalities for Agghoo.

## 2.Theoretical performance: oracle inequalities

Aim is to show that

$$\mathbb{E} \left[ \ell(s, \hat{f}_{\tau, V}^{\text{ag}}) \right] \leq C_n \mathbb{E} \left[ \inf_{m \in \mathcal{M}} \ell(s, \mathcal{A}_m(D_n^T)) \right] + r_n$$

where ideally,

- $C_n \rightarrow 1$ ,  $r_n$  negligible.
- $C_n, r_n$  only depend on  $(\mathcal{A}_m)_{m \in \mathcal{M}}$  through  $|\mathcal{M}|$  (generality).
- Few assumptions are made on  $P$ .

## 2.Theoretical performance: known results on the hold out

- Generality in classification. Margin adaptivity (G.Blanchard, P.Massart, 2006)
- Least squares regression with bounded  $Y$  and linear models (F.Navarro and A.Saumard, 2017)
- Abstract theorems - how do you check the assumptions?
- Specific settings (ex: regressograms, S.Arlot, 2008)

See 2010 survey by S.Arlot and A.Celisse for more references.

A new setting: penalization hyperparameters.

## 2.Theoretical performance: an RKHS setting

For any sample  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$ , let

$$\mathcal{A}_\lambda(D) = \operatorname{argmin}_{t \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n_t} c(t(X_i), Y_i) + \lambda \|t\|_{\mathcal{H}}^2.$$

Penalized empirical risk minimization where

- The model is an RKHS  $\mathcal{H}$
- The penalty is  $\lambda \|\cdot\|_{\mathcal{H}}^2$
- $c$  is convex and Lipschitz continuous in its first argument

Examples:  $c(u, y) = (1 - uy)_+$  (SVM),  $c(u, y) = (|u - y| - \varepsilon)_+$  (SVMR).

## 2.Theoretical performance: an oracle inequality

Now assume that:

- The kernel  $K$  is bounded.
- $c(u, y) = \gamma(u, y) = |u - y|$

### Theorem

Assume that for a.e.  $x \in \mathcal{X}$ , for all  $u \in \mathbb{R}$ ,

$$|\mathbb{P}(Y \leq u | X = x) - \mathbb{P}(Y \leq s(x) | X = x)| \geq \frac{|u - s(x)|}{\delta} \mathbb{I}_{|u - s(x)| \leq \frac{\delta}{4}}.$$

Let  $\Lambda \subset \mathbb{R}_+$  and  $\lambda_m = \min \Lambda$ . Then for all  $\theta \in (0; 1]$ ,

$$\begin{aligned} (1 - \theta)\mathbb{E}[\ell(s, \hat{f}_{\tau, V}^{\text{ag}})] &\leq (1 + \theta)\mathbb{E}\left[\min_{\lambda \in \Lambda} \ell(s, \mathcal{A}_\lambda(D_{n_t}))\right] \\ &\quad + \frac{12\delta}{n_v} + \frac{C_3}{\lambda_m} \left[ \frac{\log^2 n_v}{\theta^3 n_v^2} + \frac{\log^{\frac{3}{2}} |\Lambda|}{\theta n_v \sqrt{n_t}} \right] \end{aligned}$$

where  $n_t = \lfloor \tau n \rfloor$  and  $n_v = n - n_t$ .

# Size of $\lambda_m$ ?

- Universal consistency for  $\lambda_n \rightarrow 0$ ,  $\lambda_n >> \frac{1}{\sqrt{n}}$  (Steinwart and Christmann, Theorem 9.6)
- Gaussian kernel  $K$ ,  $\lambda_n \approx \frac{1}{n}$  minimax-optimal (Eberts and Steinwart, 2013)
- In practice: libsvm, default  $\lambda_n \approx \frac{1}{n}$ .

Conditionnally on  $D_n^{T_1}$ ,  $t_m = \mathcal{A}_m(D_n^{T_1})$  deterministic function. Let

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} P_n \gamma(t_m, \cdot) \text{ (hold-out).}$$

## Definition

Let  $w$  be a function. If  $x \mapsto \frac{w(x)}{x}$  non-increasing, for any  $\lambda > 0$ , let  $\delta(w, \lambda)$  solve  $w(x) = \lambda x^2$ .

# State-of-the-art on the hold-out II

Theorem (Massart, St-Flour lecture notes, 2003)

Assume that

$$\text{Var}_P(\gamma(t_m, \cdot) - \gamma(t_{m'}, \cdot)) \leq \left( w(\sqrt{\ell(s, t_m)}) + w(\sqrt{\ell(s, t_{m'})}) \right)^2$$
$$\|\gamma(t_m, \cdot) - \gamma(t_{m'}, \cdot)\|_\infty \leq B$$

where  $w$  non-decreasing,  $x \mapsto \frac{w(x)}{x}$  non-increasing. With probability  $\geq 1 - e^{-x}$ , for any  $\theta \in [0; 1]$ ,

$$(1 - \theta)\ell(s, t_{\hat{m}}) \leq (1 + \theta) \min_{m \in \mathcal{M}} \ell(s, t_m) + c_1 \frac{\log |\mathcal{M}| + x}{\theta} \delta(w, \sqrt{n})^2$$
$$+ c_2 (\log |\mathcal{M}| + x) \frac{B}{n}.$$

# Weaker assumptions

- $w$  non-decreasing,  $\frac{w(x)}{x}$  non-increasing  $\rightarrow w$  non-decreasing.
- $\delta(w, \lambda)$  solves  $w(\delta) = \lambda\delta^2 \rightarrow$
- $\delta(w, \lambda) = \inf \{\delta \geq 0 : \forall u \geq \delta, w(u) \leq \lambda u^2\}.$
- $\|\gamma(t_m, \cdot) - \gamma(t_{m'}, \cdot)\|_\infty \leq B \rightarrow (H_\infty) :$   
$$\|\gamma(t_m, \cdot) - \gamma(t_{m'}, \cdot)\|_\infty \leq \left( w_1(\sqrt{\ell(s, t_m)}) + w_1(\sqrt{\ell(s, t_{m'}))}) \right)^2$$

# General result on the hold-out

With probability  $\geq 1 - e^{-x}$ , for any  $\theta \in [0; 1]$ ,

$$(1 - \theta)\ell(s, t_{\hat{m}}) \leq (1 + \theta) \min_{m \in \mathcal{M}} \ell(s, t_m) + c_1 r_1(\theta, x + \log |\mathcal{M}|)$$
$$+ c_2 r_2(\theta, x + \log |\mathcal{M}|) \text{ où}$$

$\ \cdot\ _\infty$	$w$	$r_1(\theta, y)$	$r_2(\theta, y)$
borné $B$	$\frac{w(u)}{u} \downarrow$	$\frac{y}{\theta} \delta^2(w, \sqrt{n})$	$\frac{By}{n}$
$(H_\infty)$	$\frac{w(u)}{u} \downarrow$	$\frac{y}{\theta} \delta^2(w, \sqrt{n})$	$\frac{y^2}{\theta} \delta^2(w_1, n)$
$(H_\infty)$	Toute $w \uparrow$	$\theta \delta^2 \left( w, \frac{\theta}{2} \sqrt{\frac{n}{y}} \right)$	$\theta^2 \delta^2 \left( w_1, \frac{\theta^2}{4} \frac{n}{y} \right)$

Bounds for different ws can be combined.

# Application to kernel methods

Let  $\lambda < \mu$ ,  $\gamma(t, (x, y)) = c(t(x), y)$  (to simplify),  $\hat{t}_\lambda = \mathcal{A}_\lambda(D_n^T)$ .

- $\|\gamma(\hat{t}_\mu, \cdot) - \gamma(\hat{t}_\lambda, \cdot)\|_\infty \leq \|\hat{t}_\lambda - \hat{t}_\mu\|_\infty$  ( $\gamma$  1-Lipschitz).
- $\|\hat{t}_\lambda - \hat{t}_\mu\|_\infty \leq \kappa \|\hat{t}_\lambda - \hat{t}_\mu\|_{\mathcal{H}}$  (RKHS, bounded kernel)
- $\|\hat{t}_\lambda - \hat{t}_\mu\|_{\mathcal{H}}^2 \leq \|\hat{t}_\lambda\|_{\mathcal{H}}^2 - \|\hat{t}_\mu\|_{\mathcal{H}}^2$  ( $\mathcal{H}$  Hilbert,  $c$  convex)
- By definition of  $\hat{t}_\lambda$ ,

$$\begin{aligned}\lambda \left( \|\hat{t}_\lambda\|_{\mathcal{H}}^2 - \|\hat{t}_\mu\|_{\mathcal{H}}^2 \right) &\leq P_n^T \gamma(\hat{t}_\mu, \cdot) - P_{n_t} \gamma(\hat{t}_\lambda, \cdot) \\ &\leq \ell(s, \hat{t}_\mu) + (P_n^T - P) [\gamma(\hat{t}_\mu, \cdot) - \gamma(\hat{t}_\lambda, \cdot)].\end{aligned}$$

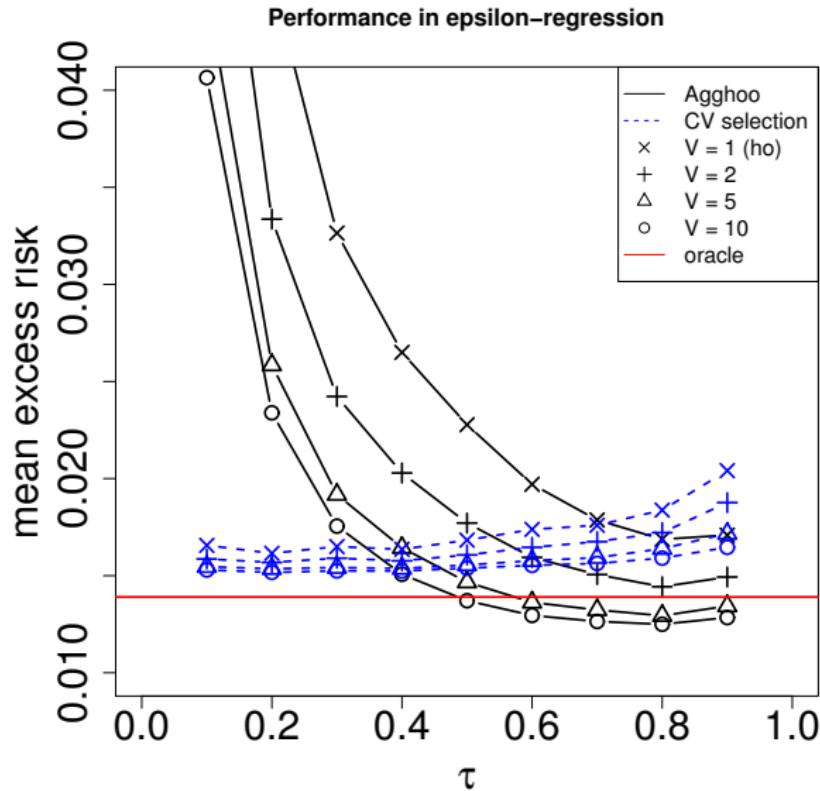
Questions not answered by theory:

- What is the effect of aggregation?
- Is Agghoo competitive with cross-validation?
- How should  $\tau$ ,  $V$  be chosen?

### 3. Simulations in regression: description

- Kernel rules  $(\mathcal{A}_\lambda)_{\lambda \in \Lambda}$  with Gaussian kernel  
 $k(x, y) = e^{-2(x-y)^2}$ , cost  $c(u, y) = (|u - y| - 0.25)_+$  and grid  
 $\Lambda = \left\{ \frac{2^{j-1}}{500\tau n} \mid 1 \leq j \leq 17 \right\}$
- $Y = e^{\cos(X)} + \xi$  where  $X \sim \mathcal{N}(0, \pi)$ ,  $\xi \sim \mathcal{N}\left(0, \frac{1}{2}\right)$ .
- I.i.d sample of size  $n = 500$ .
- Risk evaluation using  $\gamma(u, y) = |u - y|$ .

### 3. Simulations in regression: Results



## 4. Classification: Majhoo

What can be done if  $\mathcal{Y}$  is not convex?  $\mathcal{Y} = \{0; 1\}$ ,  $\gamma(u, y) = \mathbb{I}_{u \neq y}$ .

### Definition

Majority vote: Let  $(f_i)_{i=1..V}$  be predictors.

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

Majhoo:  $\widehat{f}_{T_1}^{\text{ho}}, \dots, \widehat{f}_{T_v}^{\text{ho}} \rightarrow \text{maj}\left(\widehat{f}_{T_1}^{\text{ho}}, \dots, \widehat{f}_{T_v}^{\text{ho}}\right)$ .

## 4. Classification: Theory

Comparison with hold-out: a factor 2 is lost:

$$\forall k \geq 1, \mathbb{E} \left[ \ell(s, \hat{f}_{\mathcal{T}}^{\text{maj}})^k \right] \leq 2^k \mathbb{E} \left[ \ell(s, \hat{f}_{\mathcal{T}_1}^{\text{ho}})^k \right].$$

Using Massart (2007) yields for example:

### Theorem

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

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

(margin hypothesis). Then

$$\mathbb{E}[\ell(s, \hat{f}_{\tau, V}^{\text{ag}})] \leq 3\mathbb{E} \left[ \inf_{m \in \mathcal{M}} \ell(s, \mathcal{A}(D_n^{T_1})) \right] + \frac{29c^{\frac{1}{\beta+2}} \log(e|\mathcal{M}|)}{(1-\tau)^{\frac{\beta+1}{\beta+2}} n^{\frac{\beta+1}{\beta+2}}}$$

## 4. Classification: simulation description

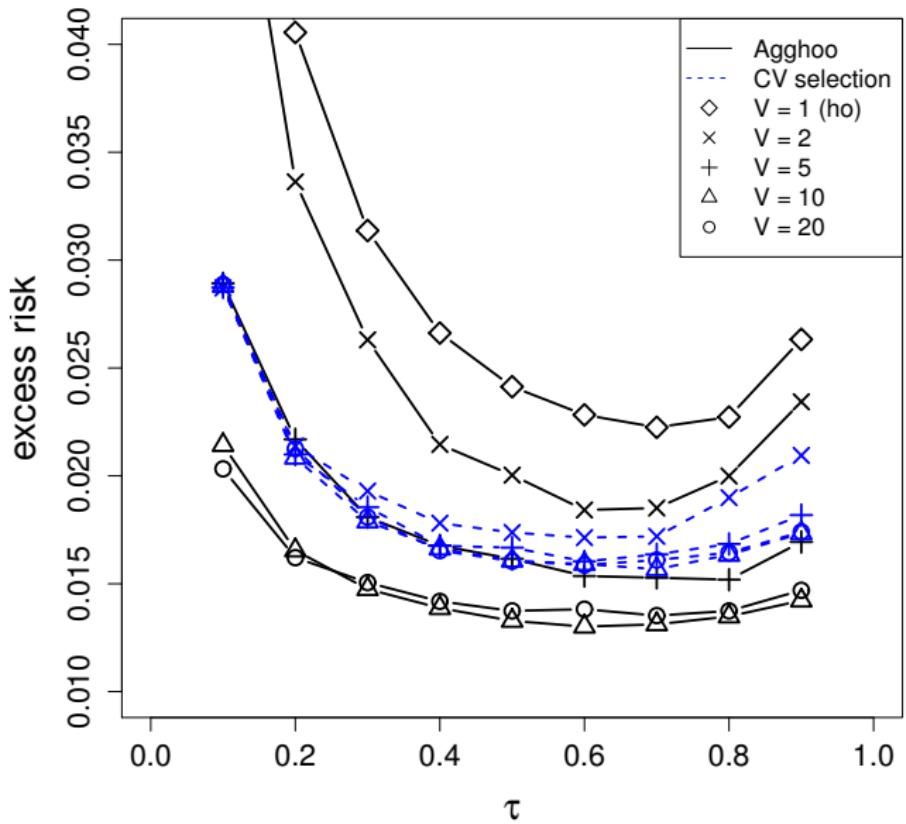
- Binary classification with  $0 - 1$  loss  $\gamma(u, y) = \mathbb{I}_{u \neq y}$ .
- k-NN rules for  $k \in \{k | 1 \leq k \leq n, k \text{ odd.}\}$
- I.i.d sample,  $n = 500$ , with distribution

$$X \sim \mathcal{U}([0; 1]^2)$$

$$\mathbb{P}(Y = 1|X) = \sigma\left(\frac{g(X) - b}{\lambda}\right) \quad \text{where} \quad \sigma(u) = \frac{1}{1 + e^{-u}},$$

$$g(u, v) = e^{-(u^2+v)^3} + u^2 + v^2$$

## 4. Classification: simulation results



- Agghoo is an alternative to cross-validation.
- It verifies an oracle inequality in regularized kernel regression and classification.
- Simulations show that Agghoo can perform better than cross-validation in regularized kernel regression and classification.

Article reference:

 G. Maillard, S. Arlot and M. Lerasle, Aggregated Hold-out, *Journal of Machine Learning Research*, (2021) v.22 , p.1 - 55

Preprint can be found at

<https://hal.archives-ouvertes.fr/hal-02273193>