

Aggregated Hold-Out

Guillaume Maillard

Université Paris Sud

February 13, 2020

Plan

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

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 $T \subset \{1, \dots, n\}$, $D_n^T = (X_j, Y_j)_{j \in T}$.
- $\mathcal{A} : \cup_{k \in \mathbb{N}} (\mathcal{X} \times \mathcal{Y})^k \rightarrow$ predictors a *learning rule*.

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 $\mathcal{T} \subset \{1, \dots, n\}$. For any learning rule \mathcal{A} ,

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

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

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

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: Agghoo

Idea: aggregate several hold out estimators.

Definition

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

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

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

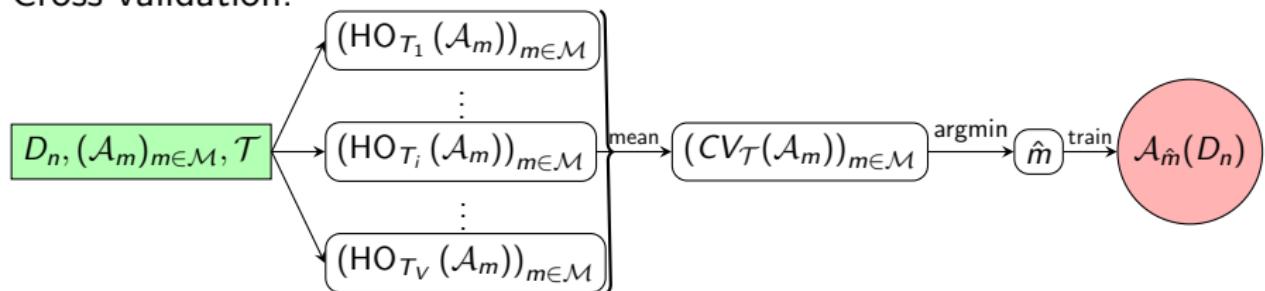
Parameters are V and τ .

V-fold CV \implies V-fold aggregation.

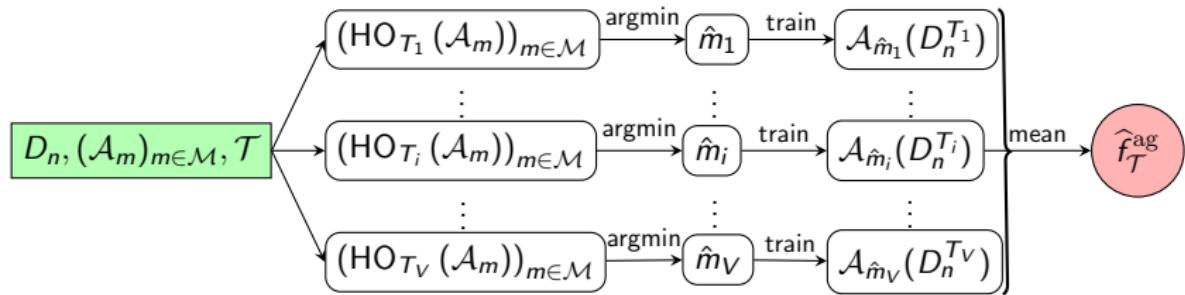
$$CV_T \implies \hat{f}_T^{\text{ag}}.$$

1. Presentation of Agghoo: comparison with cross-validation

Cross-validation:



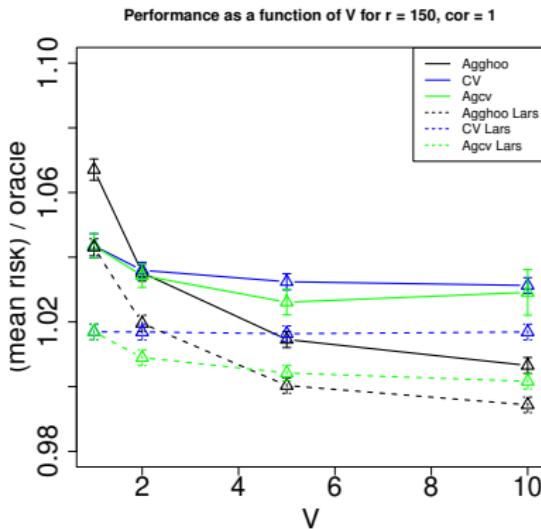
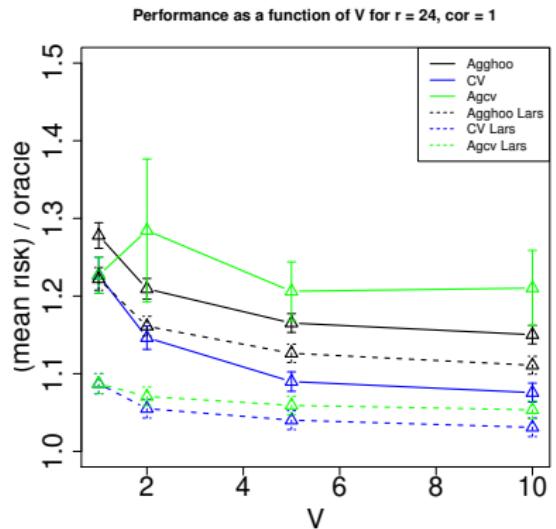
Agghoo:



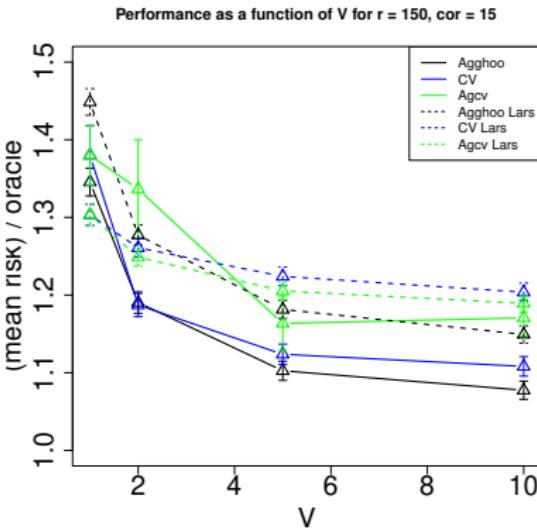
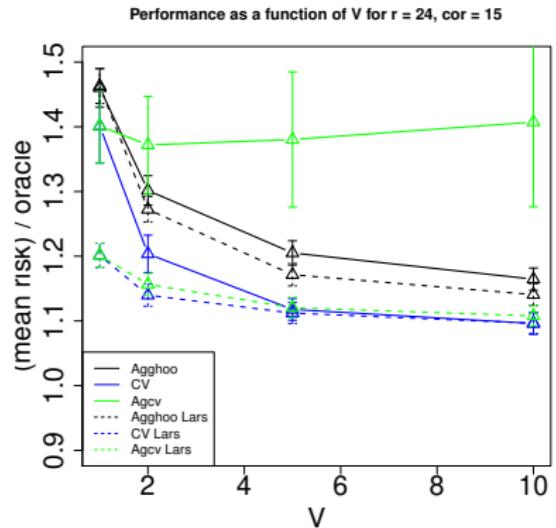
Simulation setting 1

- $Y \sim \langle w_*, X \rangle + \varepsilon$ where $\varepsilon \sim \text{Cauchy}(0, 0.08)$
- $X_j = \sum_{i=1}^{1000} \mathbb{I}_{|i-j| \leq cor} Z_i e^{-2.33^2 \frac{(i-j)^2}{2cor^2}}$ where $Z = \text{standard gaussian vector.}$
- $w_{*,j} = u_{*,g(j)}$ where $g = \text{random permutation}$ and $u_{*,j} = b$ for $1 \leq j \leq r$, $u_{*,j} = \frac{b}{4}$ for $r+1 \leq j \leq 3r$.
- b such that $\|\langle X, w_* \rangle\|_{L^2} = 1$.
- sample size $n = 100$

Agghoo vs CV for two parametrizations of the Lasso



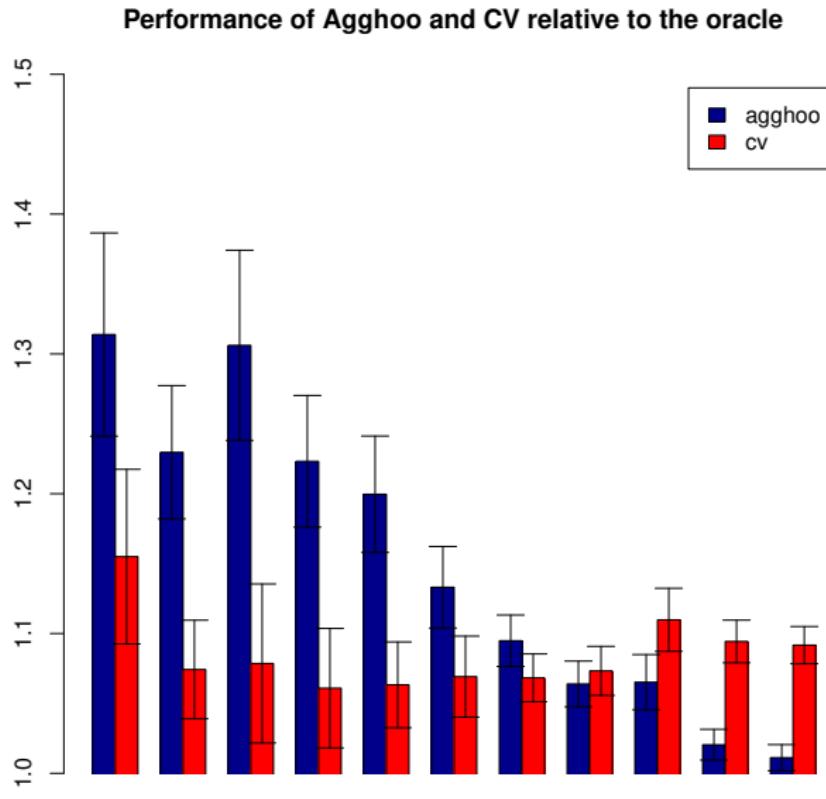
Agghoo vs CV for two parametrizations of the Lasso



Simulation setting 2

- $Y \sim \langle w_*, X \rangle + \varepsilon$ where $\varepsilon \sim \text{Cauchy}(0, 0.3)$
- $X_{js+i} = \sqrt{0.8}Z_{j+1}^0 + \sqrt{0.2}Z_{j,i}$ for $0 \leq j \leq r-1$, $1 \leq i \leq s$,
 $X_i = W_{i-r}$ for $rs < i \leq 1000$, where Z^0, Z, W independent standard gaussian.
- $w_{*,js} = b$ for $0 \leq j \leq r$, 0 otherwise.
- b such that $\|\langle X, w_* \rangle\|_{L^2} = 3$.
- sample size $n = 100$

Agghoo vs CV: impact of confounders



Simulation setting 3

- $Y \sim \langle w_*, X \rangle + \varepsilon$ where $\varepsilon \sim \text{Cauchy}(0, 0.3)$
- $X_j = \sqrt{\rho}Z^0 + \sqrt{1 - \rho}Z_j$ for $1 \leq j \leq r$, $X_i = W_{i-r}$ for $r < i \leq 1000$, where Z^0, Z, W independent standard gaussian.
- $w_{*,j} = b$ for $0 \leq j \leq r$, 0 otherwise.
- b such that $\|\langle X, w_* \rangle\|_{L^2} = 3$.
- sample size $n = 100$

Agghoo vs CV: impact of correlations between predictive covariates

