

# Comparaison de procédures de validation croisée (« V-fold »)

Sylvain Arlot (collaboration avec Matthieu Lerasle)

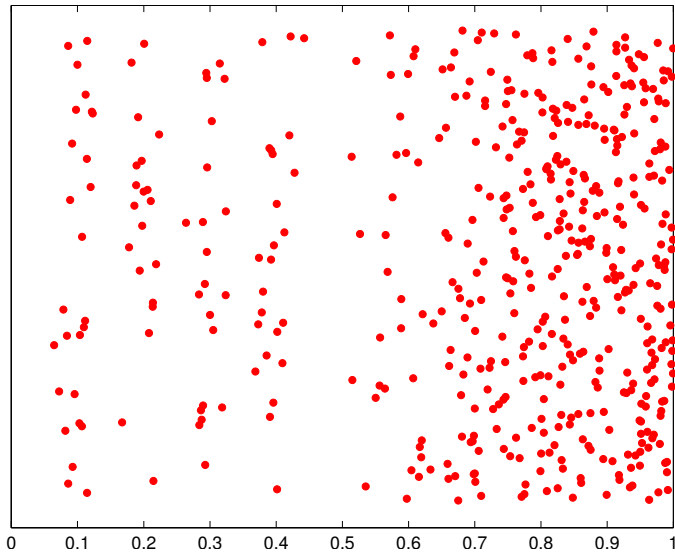
<sup>1</sup>CNRS

<sup>2</sup>École Normale Supérieure (Paris), DI/ENS, Équipe SIERRA

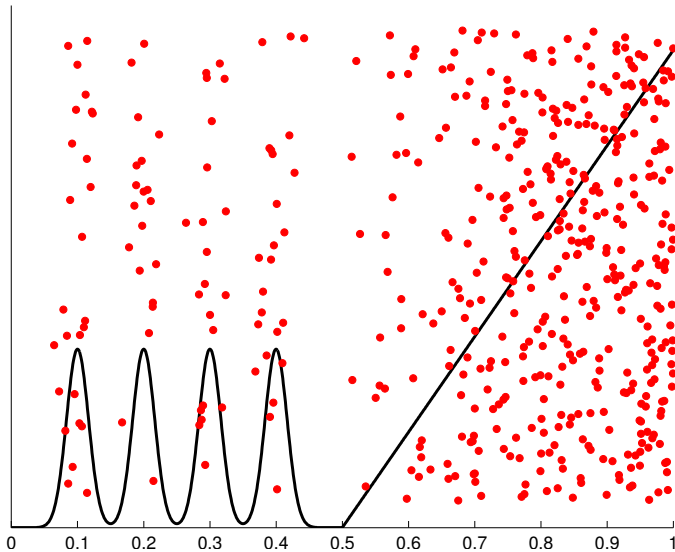
Nice, 20 février 2015

# Plan

- 1 Sélection de modèles par validation croisée
- 2 Analyse au premier ordre : biais
- 3 Analyse au deuxième ordre : variance
- 4 Conclusion

Estimation de densité : données  $\xi_1, \dots, \xi_n$ 

But : estimer la densité  $s^*$  des observations  $\xi_i$



# Problème : estimation de densité

- **Données**  $D_n$  :  $\xi_1, \dots, \xi_n \in \Xi$  (i.i.d.  $\sim P$ , densité  $s^*$  par rapport à  $\mu$ )
- **Contraste des moindres carrés**  $\gamma(t, \xi) = \|t\|_{L^2(\mu)}^2 - 2t(\xi)$
- Objectif : apprendre  $t \in \mathbb{S} = \{\text{fonctions mesurables } \Xi \rightarrow \mathbb{R}\}$   
t.q.  $\mathbb{E}_{\xi \sim P} [\gamma(t; \xi)] =: P\gamma(t)$  est minimale.

# Problème : estimation de densité

- **Données**  $D_n$  :  $\xi_1, \dots, \xi_n \in \Xi$  (i.i.d.  $\sim P$ , densité  $s^*$  par rapport à  $\mu$ )
- **Contraste des moindres carrés**  $\gamma(t, \xi) = \|t\|_{L^2(\mu)}^2 - 2t(\xi)$
- Objectif : apprendre  $t \in \mathbb{S} = \{\text{fonctions mesurables } \Xi \rightarrow \mathbb{R}\}$  t.q.  $\mathbb{E}_{\xi \sim P} [\gamma(t; \xi)] =: P\gamma(t)$  est minimale.

$$P\gamma(t) = \int t^2 d\mu - 2 \int ts^* d\mu = \int (t - s^*)^2 d\mu - \|s^*\|_{L^2(\mu)}^2$$

$\Rightarrow$  densité  $s^* \in \operatorname{argmin}_{t \in \mathbb{S}} P\gamma(t)$  et la **perte relative** vaut

$$\ell(s^*, t) := P\gamma(t) - P\gamma(s^*) = \|t - s^*\|_{L^2(\mu)}^2 .$$

# Problème : estimation de densité

- **Données**  $D_n$  :  $\xi_1, \dots, \xi_n \in \Xi$  (i.i.d.  $\sim P$ , densité  $s^*$  par rapport à  $\mu$ )
- **Contraste des moindres carrés**  $\gamma(t, \xi) = \|t\|_{L^2(\mu)}^2 - 2t(\xi)$
- Objectif : apprendre  $t \in \mathbb{S} = \{\text{fonctions mesurables } \Xi \rightarrow \mathbb{R}\}$  t.q.  $\mathbb{E}_{\xi \sim P} [\gamma(t; \xi)] =: P\gamma(t)$  est minimale.

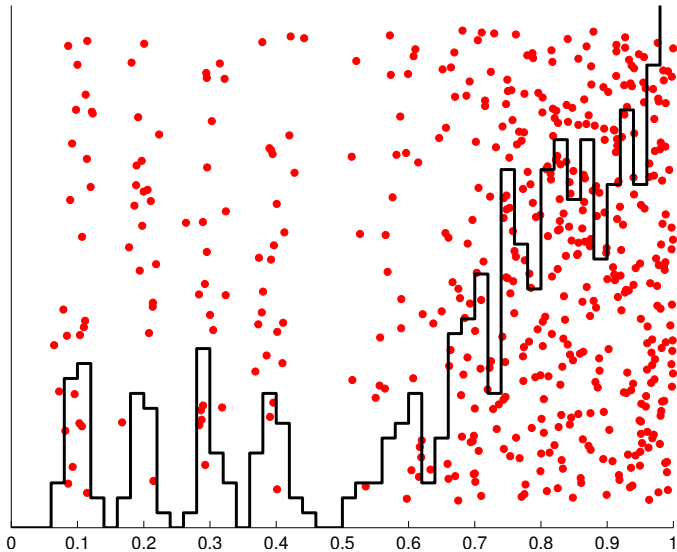
$$P\gamma(t) = \int t^2 d\mu - 2 \int ts^* d\mu = \int (t - s^*)^2 d\mu - \|s^*\|_{L^2(\mu)}^2$$

$\Rightarrow$  densité  $s^* \in \operatorname{argmin}_{t \in \mathbb{S}} P\gamma(t)$  et la **perte relative** vaut

$$\ell(s^*, t) := P\gamma(t) - P\gamma(s^*) = \|t - s^*\|_{L^2(\mu)}^2 .$$

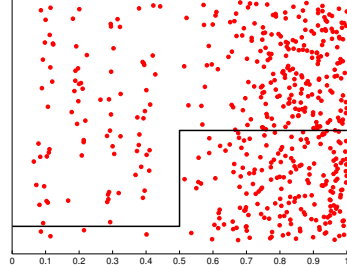
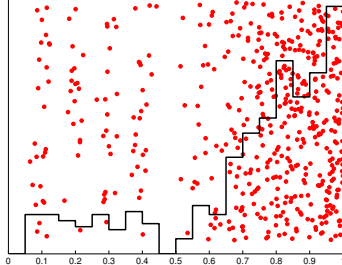
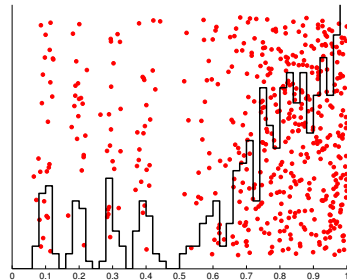
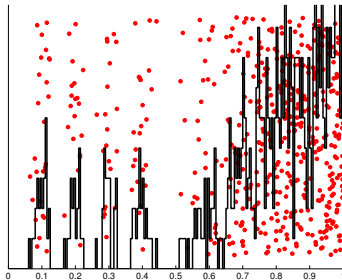
- Cas particulier d'un **cadre général** incluant aussi :
  - prédiction (régression/classification)
  - contraste log-vraisemblance en densité

# Un estimateur par histogramme





# Sélection de modèles : histogrammes réguliers



# Sélection de modèles

- **Estimateur des moindres carrés** sur un **modèle**  $\mathcal{S}_m \subset \mathbb{S}$

$$\hat{s}_m \in \operatorname{argmin}_{t \in \mathcal{S}_m} \{P_n \gamma(t)\} \quad \text{où} \quad P_n \gamma(t) := \frac{1}{n} \sum_{\xi \in D_n} \gamma(t; \xi)$$

Exemples de modèles : histogrammes, base tronquée (Fourier, ondelettes, etc.).

# Sélection de modèles

- Estimateur des moindres carrés sur un modèle  $S_m \subset \mathbb{S}$

$$\hat{s}_m \in \operatorname{argmin}_{t \in S_m} \{P_n \gamma(t)\} \quad \text{où} \quad P_n \gamma(t) := \frac{1}{n} \sum_{\xi \in D_n} \gamma(t; \xi)$$

Exemples de modèles : histogrammes, base tronquée (Fourier, ondelettes, etc.).

- Collection de modèles  $(\hat{s}_m)_{m \in \mathcal{M}} \Rightarrow$  choisir  $\hat{m} = \hat{m}(D_n)$  ?

# Sélection de modèles

- Estimateur des moindres carrés sur un modèle  $S_m \subset \mathbb{S}$

$$\hat{s}_m \in \operatorname{argmin}_{t \in S_m} \{P_n \gamma(t)\} \quad \text{où} \quad P_n \gamma(t) := \frac{1}{n} \sum_{\xi \in D_n} \gamma(t; \xi)$$

Exemples de modèles : histogrammes, base tronquée (Fourier, ondelettes, etc.).

- Collection de modèles  $(\hat{s}_m)_{m \in \mathcal{M}} \Rightarrow$  choisir  $\hat{m} = \hat{m}(D_n)$  ?
- Objectif : minimiser le risque de l'estimateur final, *i.e.*,  
**Inégalité oracle** (en espérance ou avec grande probabilité) :

$$\ell(s^*, \hat{s}_{\hat{m}}) \leq C \inf_{m \in \mathcal{M}} \{\ell(s^*, \hat{s}_m)\} + R_n$$

# Compromis biais-variance

$$\mathbb{E}[\ell(s^*, \hat{s}_m)] = \text{Biais} + \text{Variance}$$

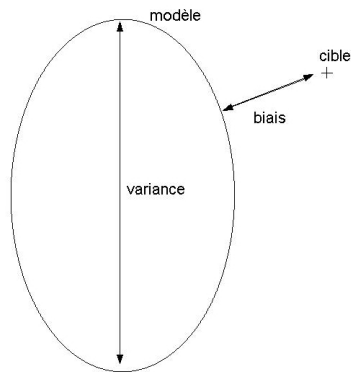
**Biais** ou Erreur d'approximation

$$\ell(s^*, S_m) = \inf_{t \in S_m} \{\ell(s^*, t)\}$$

**Variance** ou Erreur d'estimation

histogrammes réguliers sur  $\mathbb{R}$  de pas  $d_m^{-1}$  :

$$\frac{d_m - \|s_m^*\|_{L^2(\mu)}^2}{n} \approx \frac{d_m}{n}$$



# Compromis biais-variance

$$\mathbb{E}[\ell(s^*, \hat{s}_m)] = \text{Biais} + \text{Variance}$$

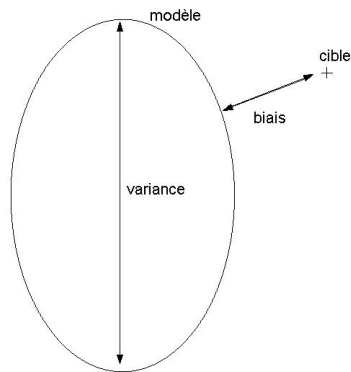
Biais ou Erreur d'approximation

$$\ell(s^*, S_m) = \inf_{t \in S_m} \{\ell(s^*, t)\}$$

Variance ou Erreur d'estimation

histogrammes réguliers sur  $\mathbb{R}$  de pas  $d_m^{-1}$  :

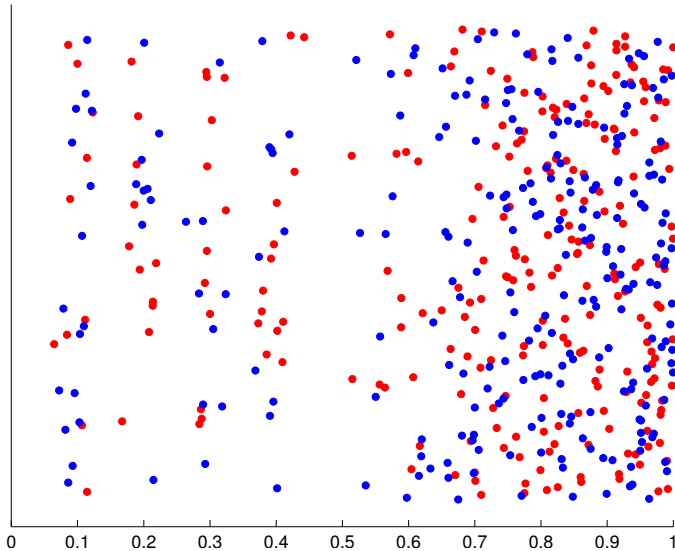
$$\frac{d_m - \|s_m^*\|_{L^2(\mu)}^2}{n} \approx \frac{d_m}{n}$$



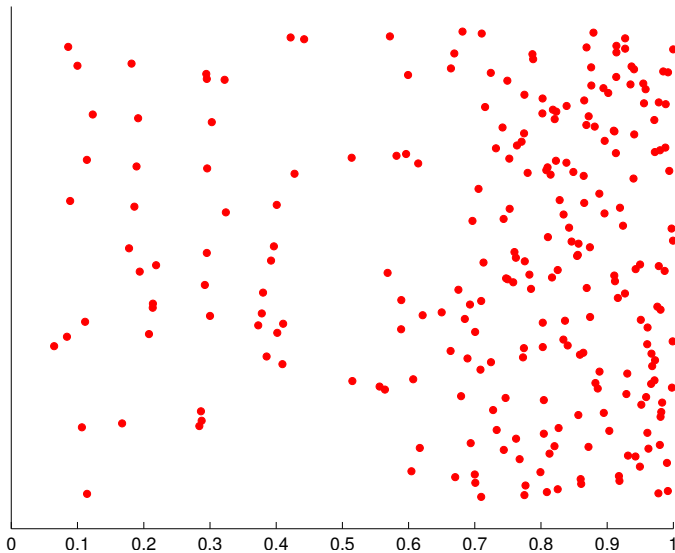
**Compromis biais-variance**

⇔ éviter le **sur-apprentissage** et le **sous-apprentissage**

# Principe de la validation simple

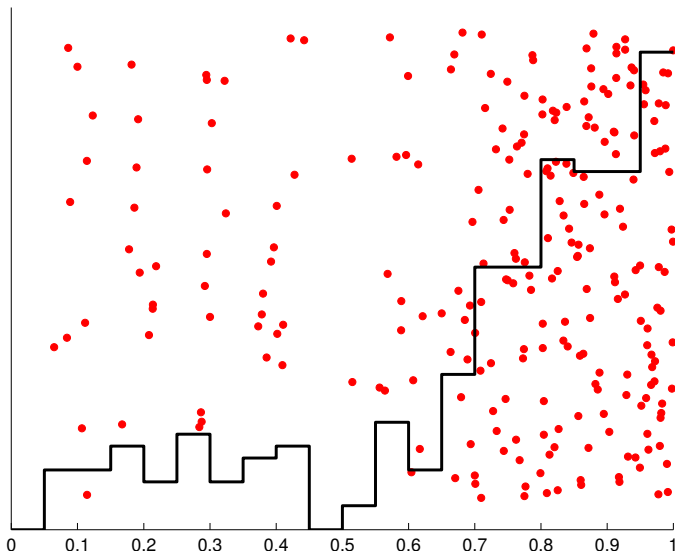


# Principe de la validation : échantillon d'entraînement

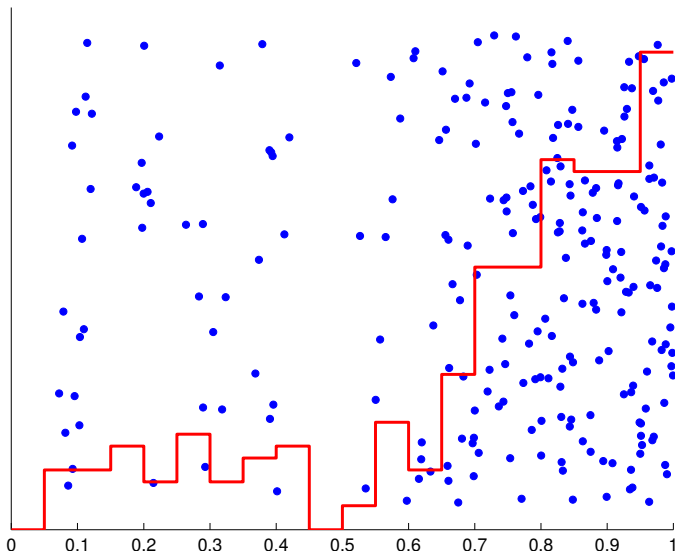




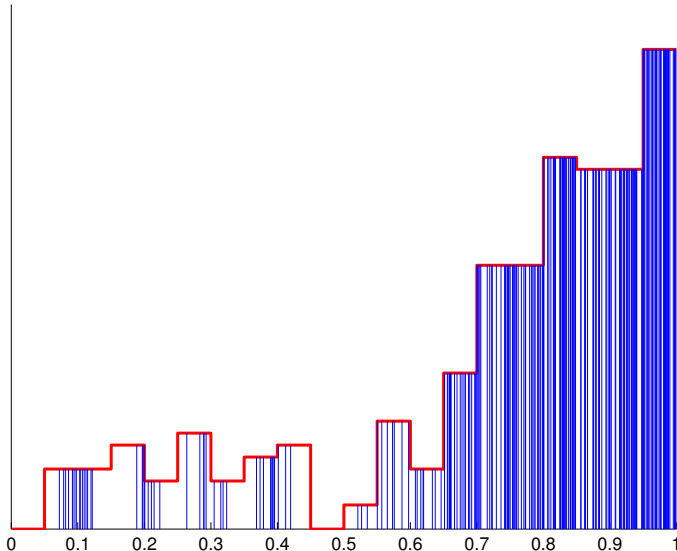
# Principe de la validation : échantillon d'entraînement



# Principe de la validation : échantillon de validation



# Principe de la validation : échantillon de validation



# Validation croisée

$$\text{Entraînement } \underbrace{D_n^{(e)}}_{\xi_1, \dots, \xi_{n_e}} \Rightarrow \hat{s}_m^{(e)} = \hat{s}_m(D_n^{(e)}) \quad \text{Validation } \underbrace{D_n^{(v)}}_{\xi_{n_e+1}, \dots, \xi_n} \Rightarrow \text{évaluer le risque}$$

# Validation croisée

$\underbrace{\xi_1, \dots, \xi_{n_e}}_{\text{Entraînement}} D_n^{(e)} \Rightarrow \hat{s}_m^{(e)} = \hat{s}_m(D_n^{(e)})$ 
 $\underbrace{\xi_{n_e+1}, \dots, \xi_n}_{\text{Validation}} D_n^{(v)} \Rightarrow \text{évaluer le risque}$

- Estimateur « Hold-out » du risque :

$$P_n^{(v)} \gamma(\hat{s}_m^{(e)}) = \frac{1}{n_v} \sum_{\xi \in D_n^{(v)}} \gamma(\hat{s}_m^{(e)}; \xi)$$

$$n_v = |D_n^{(v)}| = n - n_e$$

# Validation croisée

Entraînement  $\underbrace{\xi_1, \dots, \xi_{n_e}}_{D_n^{(e)}} \Rightarrow \hat{s}_m^{(e)} = \hat{s}_m(D_n^{(e)})$     Validation  $\underbrace{\xi_{n_e+1}, \dots, \xi_n}_{D_n^{(v)}} \Rightarrow$  évaluer le risque

- Estimateur « Hold-out » du risque :

$$P_n^{(v)} \gamma \left( \hat{s}_m^{(e)} \right) = \frac{1}{n_v} \sum_{\xi \in D_n^{(v)}} \gamma \left( \hat{s}_m^{(e)}; \xi \right) \quad n_v = |D_n^{(v)}| = n - n_e$$

- Validation croisée : moyenne d'estimateurs « hold-out »

$$\hat{\mathcal{R}}^{vc} \left( \hat{s}_m; D_n; (I_j^{(e)})_{1 \leq j \leq B} \right) = \frac{1}{B} \sum_{j=1}^B P_n^{(v,j)} \gamma \left( \hat{s}_m^{(e,j)} \right) \quad D_n^{(e,j)} = (\xi_i)_{i \in I_j^{(e)}}$$

# Validation croisée

Entraînement  $\underbrace{\xi_1, \dots, \xi_{n_e}}_{D_n^{(e)}} \Rightarrow \hat{s}_m^{(e)} = \hat{s}_m(D_n^{(e)})$     Validation  $\underbrace{\xi_{n_e+1}, \dots, \xi_n}_{D_n^{(v)}} \Rightarrow$  évaluer le risque

- Estimateur « Hold-out » du risque :

$$P_n^{(v)} \gamma(\hat{s}_m^{(e)}) = \frac{1}{n_v} \sum_{\xi \in D_n^{(v)}} \gamma(\hat{s}_m^{(e)}; \xi) \quad n_v = |D_n^{(v)}| = n - n_e$$

- Validation croisée : moyenne d'estimateurs « hold-out »

$$\hat{\mathcal{R}}^{vc}(\hat{s}_m; D_n; (I_j^{(e)})_{1 \leq j \leq B}) = \frac{1}{B} \sum_{j=1}^B P_n^{(v,j)} \gamma(\hat{s}_m^{(e,j)}) \quad D_n^{(e,j)} = (\xi_i)_{i \in I_j^{(e)}}$$

- Sélection de modèles :

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \hat{\mathcal{R}}^{vc}(\hat{s}_m; D_n) \right\} .$$

## Validation croisée « V-fold »

$\mathcal{B} = (B_j)_{1 \leq j \leq V}$  partition de  $\{1, \dots, n\}$

$$\underbrace{(\xi_i)_{i \in B_1}}_{\text{validation}}, \underbrace{(\xi_i)_{i \in B_2}, \dots, (\xi_i)_{i \in B_{V-1}}, (\xi_i)_{i \in B_V}}_{\text{entraînement}} \Rightarrow P_n^{(1)} \gamma(\widehat{S}_m^{(-1)})$$



## Validation croisée « V-fold »

$\mathcal{B} = (B_j)_{1 \leq j \leq V}$  partition de  $\{1, \dots, n\}$

$$\begin{array}{l}
 \underbrace{(\xi_i)_{i \in B_1}}_{\text{validation}}, \underbrace{(\xi_i)_{i \in B_2}, \dots, (\xi_i)_{i \in B_{V-1}}, (\xi_i)_{i \in B_V}}_{\text{entraînement}} \Rightarrow P_n^{(1)} \gamma(\widehat{S}_m^{(-1)}) \\
 \underbrace{(\xi_i)_{i \in B_1}}_{\text{entraîn.}}, \underbrace{(\xi_i)_{i \in B_2}}_{\text{validation}}, \underbrace{(\xi_i)_{i \in B_{V-1}}, (\xi_i)_{i \in B_V}}_{\text{entraînement}} \Rightarrow P_n^{(2)} \gamma(\widehat{S}_m^{(-2)})
 \end{array}$$

## Validation croisée « V-fold »

$\mathcal{B} = (B_j)_{1 \leq j \leq V}$  partition de  $\{1, \dots, n\}$

$$\begin{array}{l}
 \underbrace{(\xi_i)_{i \in B_1}}_{\text{validation}}, \underbrace{(\xi_i)_{i \in B_2}, \dots, (\xi_i)_{i \in B_{V-1}}, (\xi_i)_{i \in B_V}}_{\text{entraînement}} \Rightarrow P_n^{(1)} \gamma(\widehat{S}_m^{(-1)}) \\
 \underbrace{(\xi_i)_{i \in B_1}, (\xi_i)_{i \in B_2}}_{\text{entraîn.}}, \underbrace{(\xi_i)_{i \in B_{V-1}}, (\xi_i)_{i \in B_V}}_{\text{entraînement}}, \dots \Rightarrow P_n^{(2)} \gamma(\widehat{S}_m^{(-2)}) \\
 \vdots \\
 \underbrace{(\xi_i)_{i \in B_1}, (\xi_i)_{i \in B_2}, \dots, (\xi_i)_{i \in B_{V-1}}}_{\text{entraînement}}, \underbrace{(\xi_i)_{i \in B_V}}_{\text{validation}} \Rightarrow P_n^{(V)} \gamma(\widehat{S}_m^{(-V)})
 \end{array}$$

## Validation croisée « V-fold »

$\mathcal{B} = (B_j)_{1 \leq j \leq V}$  partition de  $\{1, \dots, n\}$

$\underbrace{(\xi_i)_{i \in B_1}}_{\text{validation}}, \underbrace{(\xi_i)_{i \in B_2}, \dots, (\xi_i)_{i \in B_{V-1}}, (\xi_i)_{i \in B_V}}_{\text{entraînement}} \Rightarrow P_n^{(1)} \gamma(\widehat{S}_m^{(-1)})$

$\underbrace{(\xi_i)_{i \in B_1}}_{\text{entraîn.}}, \underbrace{(\xi_i)_{i \in B_2}}_{\text{validation}}, \underbrace{(\xi_i)_{i \in B_{V-1}}, (\xi_i)_{i \in B_V}}_{\text{entraînement}} \Rightarrow P_n^{(2)} \gamma(\widehat{S}_m^{(-2)})$

⋮

$\underbrace{(\xi_i)_{i \in B_1}, (\xi_i)_{i \in B_2}, \dots, (\xi_i)_{i \in B_{V-1}}}_{\text{entraînement}}, \underbrace{(\xi_i)_{i \in B_V}}_{\text{validation}} \Rightarrow P_n^{(V)} \gamma(\widehat{S}_m^{(-V)})$

$$\Rightarrow \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) = \frac{1}{V} \sum_{j=1}^V P_n^{(j)} \gamma(\widehat{S}_m^{(-j)}) \quad \widehat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m) \right\}$$

# Analyse au premier ordre : lemme

- Objectif :

$$\text{minimiser } \mathcal{R}(m) = P(\hat{s}_m(D_n)) \text{ sur } m \in \mathcal{M}$$

- Méthode :

$$\text{minimiser } \mathcal{C}(m) = \hat{\mathcal{R}}^{\text{vf}}(\hat{s}_m; D_n; \mathcal{B}) \text{ sur } m \in \mathcal{M}$$

$$\Rightarrow \hat{m}_{\mathcal{C}} \in \underset{m \in \mathcal{M}}{\operatorname{argmin}} \{ \mathcal{C}(m) \}$$

# Analyse au premier ordre : lemme

- Objectif :

$$\text{minimiser } \mathcal{R}(m) = P(\hat{s}_m(D_n)) \text{ sur } m \in \mathcal{M}$$

- Méthode :

$$\text{minimiser } \mathcal{C}(m) = \hat{\mathcal{R}}^{\text{vf}}(\hat{s}_m; D_n; \mathcal{B}) \text{ sur } m \in \mathcal{M}$$

$$\Rightarrow \hat{m}_{\mathcal{C}} \in \underset{m \in \mathcal{M}}{\operatorname{argmin}} \{ \mathcal{C}(m) \}$$

## Lemme

Si  $\forall m \in \mathcal{M}, \quad -B(m) \leq \mathcal{C}(m) - \mathcal{R}(m) \leq A(m),$

alors,  $\forall \hat{m}_{\mathcal{C}} \in \underset{m \in \mathcal{M}}{\operatorname{argmin}} \{ \mathcal{C}(m) \},$

$$\mathcal{R}(\hat{m}_{\mathcal{C}}) - B(\hat{m}_{\mathcal{C}}) \leq \inf_{m \in \mathcal{M}} \{ \mathcal{R}(m) + A(m) \} .$$

# Optimalité au premier ordre

- Principe d'estimation sans biais du risque (Mallows, Akaike, 1973) : choisir  $\mathcal{C}$  tel que

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\mathcal{C}(m)] = \mathbb{E}[\mathcal{R}(m)] \quad .$$

# Optimalité au premier ordre

- Principe d'estimation sans biais du risque (Mallows, Akaike, 1973) : choisir  $\mathcal{C}$  tel que

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\mathcal{C}(m)] = \mathbb{E}[\mathcal{R}(m)] \quad .$$

- Sous réserve d'inégalités de concentration (uniformes sur  $m \in \mathcal{M}$ ), avec grande probabilité,

$$\forall m \in \mathcal{M}, \quad -\delta_n \mathcal{R}(m) \leq \mathcal{C}(m) - \mathcal{R}(m) \leq \delta_n \mathcal{R}(m)$$

avec  $\delta_n \in ]0, 1[$ .

# Optimalité au premier ordre

- Principe d'estimation sans biais du risque (Mallows, Akaike, 1973) : choisir  $\mathcal{C}$  tel que

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\mathcal{C}(m)] = \mathbb{E}[\mathcal{R}(m)] \quad .$$

- Sous réserve d'inégalités de concentration (uniformes sur  $m \in \mathcal{M}$ ), avec grande probabilité,

$$\forall m \in \mathcal{M}, \quad -\delta_n \mathcal{R}(m) \leq \mathcal{C}(m) - \mathcal{R}(m) \leq \delta_n \mathcal{R}(m)$$

avec  $\delta_n \in ]0, 1[$ .

⇒ d'après le lemme,

$$\forall \hat{m}_{\mathcal{C}} \in \underset{m \in \mathcal{M}}{\operatorname{argmin}} \{ \mathcal{C}(m) \}, \quad \mathcal{R}(\hat{m}_{\mathcal{C}}) \leq \frac{1 + \delta_n}{1 - \delta_n} \inf_{m \in \mathcal{M}} \{ \mathcal{R}(m) \} \quad .$$



# Optimalité au premier ordre

- Principe d'estimation sans biais du risque (Mallows, Akaike, 1973) : choisir  $\mathcal{C}$  tel que

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\mathcal{C}(m)] = \mathbb{E}[\mathcal{R}(m)] \quad .$$

- Sous réserve d'inégalités de concentration (uniformes sur  $m \in \mathcal{M}$ ), avec grande probabilité,

$$\forall m \in \mathcal{M}, \quad -\delta_n \mathcal{R}(m) \leq \mathcal{C}(m) - \mathcal{R}(m) \leq \delta_n \mathcal{R}(m)$$

avec  $\delta_n \in ]0, 1[$ .

⇒ d'après le lemme,

$$\forall \hat{m}_{\mathcal{C}} \in \operatorname{argmin}_{m \in \mathcal{M}} \{\mathcal{C}(m)\}, \quad \mathcal{R}(\hat{m}_{\mathcal{C}}) \leq \frac{1 + \delta_n}{1 - \delta_n} \inf_{m \in \mathcal{M}} \{\mathcal{R}(m)\} \quad .$$

- Optimal au premier ordre si  $\delta_n \rightarrow 0$ .**

# Au premier ordre : biais de la validation croisée

- Hypothèse :  $\text{Card}(B_j) = n/V$  pour tout  $j$ .
- Calcul d'espérances (moindres carrés, densité ou régression) :

$$\mathbb{E}[P\gamma(\hat{s}_m(D_n))] \approx \alpha(m) + \frac{\beta(m)}{n}$$

$$\Rightarrow \mathbb{E}[\hat{\mathcal{R}}^{\text{vf}}(\hat{s}_m; D_n; \mathcal{B})] = \mathbb{E}[P_n^{(j)}\gamma(\hat{s}_m^{(-j)})] = \mathbb{E}[P\gamma(\hat{s}_m^{(-j)})]$$

$$\approx \alpha(m) + \frac{V}{V-1} \frac{\beta(m)}{n}$$

# Au premier ordre : biais de la validation croisée

- Hypothèse :  $\text{Card}(B_j) = n/V$  pour tout  $j$ .
- Calcul d'espérances (moindres carrés, densité ou régression) :

$$\begin{aligned} \mathbb{E}[P\gamma(\hat{s}_m(D_n))] &\approx \alpha(m) + \frac{\beta(m)}{n} \\ \Rightarrow \mathbb{E}[\hat{\mathcal{R}}^{\text{vf}}(\hat{s}_m; D_n; \mathcal{B})] &= \mathbb{E}[P_n^{(j)}\gamma(\hat{s}_m^{(-j)})] = \mathbb{E}[P\gamma(\hat{s}_m^{(-j)})] \\ &\approx \alpha(m) + \frac{V}{V-1} \frac{\beta(m)}{n} \end{aligned}$$

⇒ **biais**, décroissant avec  $V$ , tend vers zéro quand  $V \rightarrow +\infty$

# Au premier ordre : biais de la validation croisée

- Hypothèse :  $\text{Card}(B_j) = n/V$  pour tout  $j$ .
- Calcul d'espérances (moindres carrés, densité ou régression) :

$$\begin{aligned} \mathbb{E}[P\gamma(\hat{s}_m(D_n))] &\approx \alpha(m) + \frac{\beta(m)}{n} \\ \Rightarrow \mathbb{E}[\hat{\mathcal{R}}^{\text{vf}}(\hat{s}_m; D_n; \mathcal{B})] &= \mathbb{E}[P_n^{(j)}\gamma(\hat{s}_m^{(-j)})] = \mathbb{E}[P\gamma(\hat{s}_m^{(-j)})] \\ &\approx \alpha(m) + \frac{V}{V-1} \frac{\beta(m)}{n} \end{aligned}$$

⇒ **biais**, décroissant avec  $V$ , tend vers zéro quand  $V \rightarrow +\infty$

⇒ **sous-optimalité** de la validation croisée «  $V$ -fold » à  $V$  fixé  
(A. 2008, régressogrammes; valable plus largement)

## Correction du biais et pénalisation « V-fold »

- Validation croisée « V-fold » corrigée (Burman, 1989) :

$$\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m; D_n; \mathcal{B}) := \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) + P_{n\gamma}(\widehat{s}_m) - \frac{1}{V} \sum_{j=1}^V P_{n\gamma}(\widehat{s}_m^{(-j)})$$

# Correction du biais et pénalisation « V-fold »

- Validation croisée « V-fold » corrigée (Burman, 1989) :

$$\begin{aligned}\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m; D_n; \mathcal{B}) &:= \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) + P_n\gamma(\widehat{s}_m) - \frac{1}{V} \sum_{j=1}^V P_n\gamma(\widehat{s}_m^{(-j)}) \\ &= P_n\gamma(\widehat{s}_m) + \underbrace{\text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B})}_{\text{pénalité V-fold (A. 2008)}}\end{aligned}$$

## Correction du biais et pénalisation « V-fold »

- Validation croisée « V-fold » corrigée (Burman, 1989) :

$$\begin{aligned}\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m; D_n; \mathcal{B}) &:= \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) + P_n\gamma(\widehat{s}_m) - \frac{1}{V} \sum_{j=1}^V P_n\gamma(\widehat{s}_m^{(-j)}) \\ &= P_n\gamma(\widehat{s}_m) + \underbrace{\text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B})}_{\text{pénalité V-fold (A. 2008)}}\end{aligned}$$

- Estimation de densité, moindres carrés (A. & Lerasle, 2014) :

$$\widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) = P_n\gamma(\widehat{s}_m(D_n)) + \underbrace{\left(1 + \frac{1}{2(V-1)}\right)}_{\text{surpénalisation}} \text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B})$$

$$\widehat{\mathcal{R}}^{\text{lp}}(\widehat{s}_m; D_n; \mathcal{B}) = P_n\gamma(\widehat{s}_m(D_n)) + \underbrace{\left(1 + \frac{1}{2\left(\frac{n}{p} - 1\right)}\right)}_{\text{surpénalisation}} \text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B}_{\text{loo}})$$

## Inégalités oracle optimales pour la pénalisation « V-fold »

## Théorème

Avec probabilité  $1 - n^{-2}$ ,  $\forall \delta > 0$ ,

$$\forall \hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \{ P_n \gamma(\hat{s}_m(D_n)) + \operatorname{pen}_{\text{VF}}(\hat{s}_m; D_n; \mathcal{B}) \},$$

$$\ell(s^*, \hat{s}_{\hat{m}}) \leq (1 + \delta) \inf_{m \in \mathcal{M}} \{ \ell(s^*, \hat{s}_m) \} + \frac{L [\log(\operatorname{Card}(\mathcal{M})) \vee \log(n)]^\alpha}{\delta^\beta n}$$

$\Rightarrow$  Optimal au premier ordre si  $\operatorname{Card}(\mathcal{M}) \leq an^b$



## Inégalités oracle optimales pour la pénalisation « V-fold »

## Théorème

Avec probabilité  $1 - n^{-2}$ ,  $\forall \delta > 0$ ,

$$\forall \hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \{ P_n \gamma(\hat{s}_m(D_n)) + \operatorname{pen}_{\text{VF}}(\hat{s}_m; D_n; \mathcal{B}) \},$$

$$\ell(s^*, \hat{s}_{\hat{m}}) \leq (1 + \delta) \inf_{m \in \mathcal{M}} \{ \ell(s^*, \hat{s}_m) \} + \frac{L [\log(\operatorname{Card}(\mathcal{M})) \vee \log(n)]^\alpha}{\delta^\beta n}$$

$\Rightarrow$  Optimal au premier ordre si  $\operatorname{Card}(\mathcal{M}) \leq an^b$

Valable sous des hypothèses raisonnablement faibles pour :

- Les **régressogrammes** en régression hétéroscédastique (A. 2008, 2009)
- L'**estimation de densité par moindres carrés** (A. & Lerasle, 2014; Celisse, 2014)

## Inégalités oracle optimales pour la pénalisation « V-fold »

## Théorème

Avec probabilité  $1 - n^{-2}$ ,  $\forall \delta > 0$ ,

$$\forall \hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \{ P_n \gamma(\hat{s}_m(D_n)) + \operatorname{pen}_{\text{VF}}(\hat{s}_m; D_n; \mathcal{B}) \},$$

$$\ell(s^*, \hat{s}_{\hat{m}}) \leq (1 + \delta) \inf_{m \in \mathcal{M}} \{ \ell(s^*, \hat{s}_m) \} + \frac{L [\log(\operatorname{Card}(\mathcal{M})) \vee \log(n)]^\alpha}{\delta^\beta n}$$

$\Rightarrow$  Optimal au premier ordre si  $\operatorname{Card}(\mathcal{M}) \leq an^b$

Valable sous des hypothèses raisonnablement faibles pour :

- Les **régressogrammes** en régression hétéroscédastique (A. 2008, 2009)
  - L'**estimation de densité par moindres carrés** (A. & Lerasle, 2014; Celisse, 2014)
- + Inégalité-oracle sous-optimale pour la validation croisée « V-fold » (constante  $1 + \frac{1}{V-1}$ ).

# Analyse au second ordre ?

- Comment comparer  $\hat{m}_{\mathcal{C}_1}$  et  $\hat{m}_{\mathcal{C}_2}$  lorsque

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\mathcal{C}_1(m)] = \mathbb{E}[\mathcal{C}_2(m)] \quad ?$$

# Analyse au second ordre ?

- Comment comparer  $\hat{m}_{C_1}$  et  $\hat{m}_{C_2}$  lorsque

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[C_1(m)] = \mathbb{E}[C_2(m)] \quad ?$$

- Tenir compte de la variance  $\text{var}(C_i(m))$  ?

# Analyse au second ordre ?

- Comment comparer  $\hat{m}_{C_1}$  et  $\hat{m}_{C_2}$  lorsque

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\mathcal{C}_1(m)] = \mathbb{E}[\mathcal{C}_2(m)] \quad ?$$

- Tenir compte de la variance  $\text{var}(\mathcal{C}_i(m))$  ?

- Variance de quelle quantité ?

Pour toute variable  $Z$ ,  $\hat{m}_C \in \operatorname{argmin}_{m \in \mathcal{M}} \{\mathcal{C}(m) + Z\}$   
mais  $\text{var}(\mathcal{C}(m) + Z)$  dépend de  $Z$ ...

# Analyse au second ordre ?

- Comment comparer  $\hat{m}_{C_1}$  et  $\hat{m}_{C_2}$  lorsque

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\mathcal{C}_1(m)] = \mathbb{E}[\mathcal{C}_2(m)] \quad ?$$

- Tenir compte de la variance  $\text{var}(\mathcal{C}_i(m))$  ?
- Variance de quelle quantité ?

Pour toute variable  $Z$ ,  $\hat{m}_C \in \operatorname{argmin}_{m \in \mathcal{M}} \{\mathcal{C}(m) + Z\}$   
mais  $\text{var}(\mathcal{C}(m) + Z)$  dépend de  $Z$ ...

- Ce qui compte :

$$\forall m, m' \in \mathcal{M}, \quad \text{sign}(\mathcal{C}(m) - \mathcal{C}(m')) = \text{sign}(\mathcal{R}(m) - \mathcal{R}(m'))$$

# Analyse au second ordre ?

- Comment comparer  $\hat{m}_{C_1}$  et  $\hat{m}_{C_2}$  lorsque

$$\forall m \in \mathcal{M}, \quad \mathbb{E}[\mathcal{C}_1(m)] = \mathbb{E}[\mathcal{C}_2(m)] \quad ?$$

- Tenir compte de la variance  $\text{var}(\mathcal{C}_i(m))$  ?
- Variance de quelle quantité ?

Pour toute variable  $Z$ ,  $\hat{m}_C \in \operatorname{argmin}_{m \in \mathcal{M}} \{\mathcal{C}(m) + Z\}$   
 mais  $\text{var}(\mathcal{C}(m) + Z)$  dépend de  $Z$ ...

- Ce qui compte :

$$\forall m, m' \in \mathcal{M}, \quad \text{sign}(\mathcal{C}(m) - \mathcal{C}(m')) = \text{sign}(\mathcal{R}(m) - \mathcal{R}(m'))$$

⇒ variance des incréments

$$\text{var}(\mathcal{C}(m) - \mathcal{C}(m')) .$$

# Variance et sélection de modèles : heuristique

- $\forall m \notin \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\ell(s^*, \hat{s}_m)]$ , on veut minimiser

$$\mathbb{P}(\hat{m}_C = m)$$



# Variance et sélection de modèles : heuristique

- $\forall m \notin \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\ell(s^*, \hat{s}_m)]$ , on veut minimiser

$$\begin{aligned} \mathbb{P}(\hat{m}_C = m) &= \mathbb{P}(\forall m' \in \mathcal{M}, \mathcal{C}(m) - \mathcal{C}(m') < 0) \\ &\leq \min_{m' \neq m} \mathbb{P}(\mathcal{C}(m) - \mathcal{C}(m') < 0) \end{aligned}$$

# Variance et sélection de modèles : heuristique

- $\forall m \notin \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\ell(s^*, \hat{s}_m)]$ , on veut minimiser

$$\mathbb{P}(\hat{m}_C = m) = \mathbb{P}(\forall m' \in \mathcal{M}, \mathcal{C}(m) - \mathcal{C}(m') < 0)$$

$$\leq \min_{m' \neq m} \mathbb{P}(\mathcal{C}(m) - \mathcal{C}(m') < 0)$$

$$\approx \min_{m' \neq m} \mathbb{P}\left(\mathbb{E}[\mathcal{C}(m) - \mathcal{C}(m')] - \mathcal{N}\sqrt{\operatorname{var}(\mathcal{C}(m) - \mathcal{C}(m'))} < 0\right)$$

# Variance et sélection de modèles : heuristique

- $\forall m \notin \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\ell(s^*, \hat{s}_m)]$ , on veut minimiser

$$\mathbb{P}(\hat{m}_C = m) = \mathbb{P}(\forall m' \in \mathcal{M}, \mathcal{C}(m) - \mathcal{C}(m') < 0)$$

$$\leq \min_{m' \neq m} \mathbb{P}(\mathcal{C}(m) - \mathcal{C}(m') < 0)$$

$$\approx \min_{m' \neq m} \mathbb{P}\left(\mathbb{E}[\mathcal{C}(m) - \mathcal{C}(m')] - \mathcal{N} \sqrt{\operatorname{var}(\mathcal{C}(m) - \mathcal{C}(m'))} < 0\right)$$

$$= \bar{\Phi}\left(\max_{m' \neq m} \frac{\mathbb{E}[\mathcal{C}(m) - \mathcal{C}(m')]}{\sqrt{\operatorname{var}(\mathcal{C}(m) - \mathcal{C}(m'))}}\right) \quad \text{où } \bar{\Phi}(t) = \mathbb{P}(\mathcal{N} > t)$$

# Variance et sélection de modèles : heuristique

- $\forall m \notin \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\ell(s^*, \hat{s}_m)]$ , on veut minimiser

$$\begin{aligned} \mathbb{P}(\hat{m}_C = m) &= \mathbb{P}(\forall m' \in \mathcal{M}, \mathcal{C}(m) - \mathcal{C}(m') < 0) \\ &\leq \min_{m' \neq m} \mathbb{P}(\mathcal{C}(m) - \mathcal{C}(m') < 0) \\ &\approx \min_{m' \neq m} \mathbb{P}\left(\mathbb{E}[\mathcal{C}(m) - \mathcal{C}(m')] - \mathcal{N}\sqrt{\operatorname{var}(\mathcal{C}(m) - \mathcal{C}(m'))} < 0\right) \\ &= \bar{\Phi}\left(\max_{m' \neq m} \frac{\mathbb{E}[\mathcal{C}(m) - \mathcal{C}(m')]}{\sqrt{\operatorname{var}(\mathcal{C}(m) - \mathcal{C}(m'))}}\right) \quad \text{où } \bar{\Phi}(t) = \mathbb{P}(\mathcal{N} > t) \end{aligned}$$

- **Hypothèses** :  $\forall m \in \mathcal{M}, \mathbb{E}[\mathcal{C}_1(m)] = \mathbb{E}[\mathcal{C}_2(m)]$  et  
 $\forall i, \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\mathcal{C}_i(m)] = \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\mathcal{R}(m)]$   
 $\Rightarrow$  on veut minimiser  $\operatorname{var}(\mathcal{C}_i(m) - \mathcal{C}_i(m'))$

# Variance et sélection de modèles (densité, moindres carrés)

$$\Delta(m, m', V) = \widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m) - \widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_{m'})$$

Théorème (A. & Lerasle, 2014)

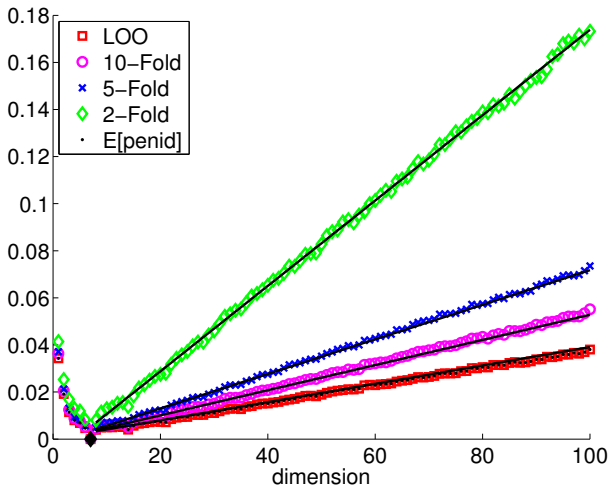
$$\begin{aligned} \text{var}(\Delta(m, m', V)) &= 4 \left( 1 + \frac{2}{n} + \frac{1}{n^2} \right) \frac{\text{var}_P(s_m^* - s_{m'}^*)}{n} \\ &\quad + 2 \left( 1 + \frac{4}{V-1} - \frac{1}{n} \right) \underbrace{\frac{B(m, m')}{n^2}}_{\geq 0} \end{aligned}$$

Si de plus  $S_m \subset S_{m'}$  sont deux modèles d'histogrammes réguliers de pas  $d_m^{-1}$ ,  $d_{m'}^{-1}$ , alors

$$B(m, m') \propto \|s_m^* - s_{m'}^*\| d_m .$$

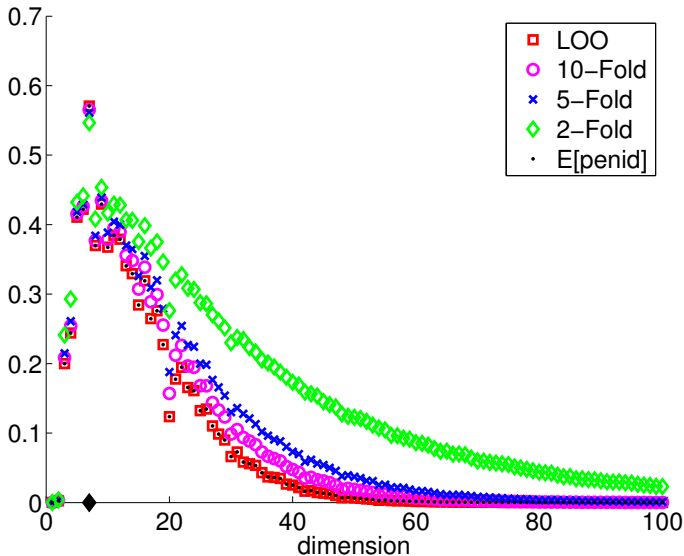
Les deux termes sont du même ordre si  $\|s_m^* - s_{m'}^*\| \approx d_m/n$ .

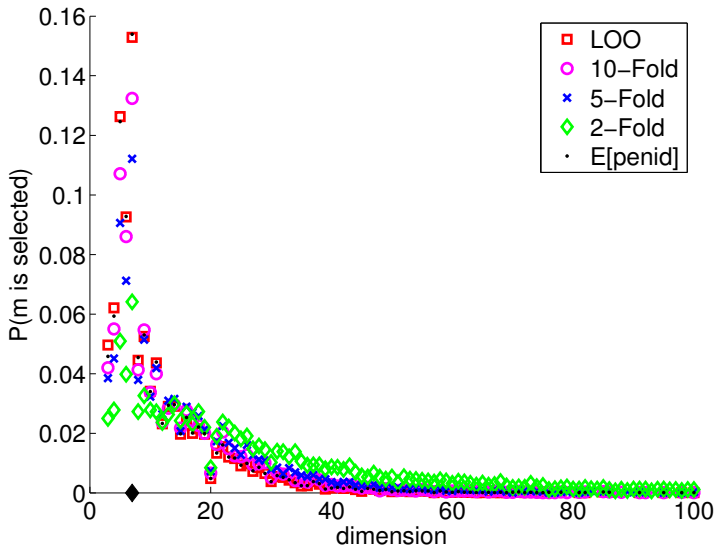
# Variance de $\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m) - \widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_{m^*})$ vs. $(d_m, V)$



$$\text{var}(\Delta(m, m', V)) \approx n^{-2} \left[ 29 \left( 1 + \frac{0.8}{V-1} \right) + 3.7 \left( 1 + \frac{3.8}{V-1} \right) (d_m - d_{m^*}) \right]$$

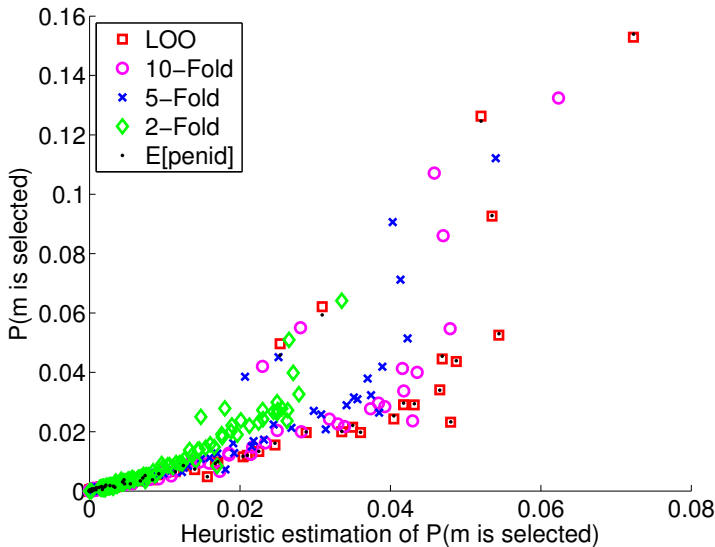
$$\overline{\Phi}(\max_{m' \neq m} \mathbb{E}[\mathcal{C}(m) - \mathcal{C}(m')] / \sqrt{\text{var}(\mathcal{C}(m) - \mathcal{C}(m'))})$$



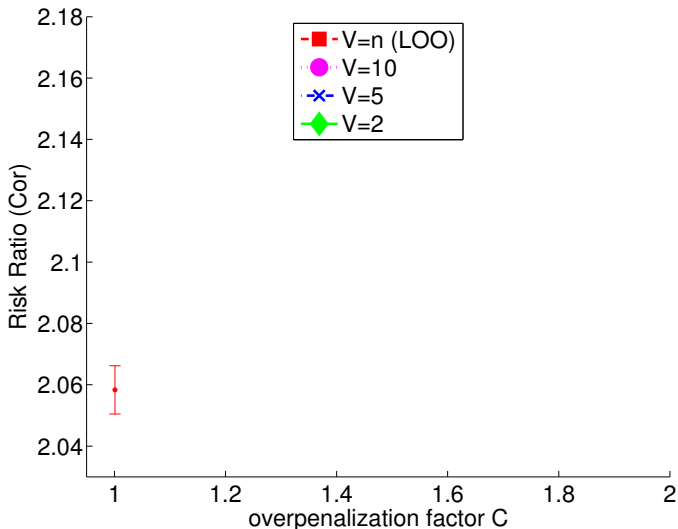
Évaluation de l'heuristique :  $\mathbb{P}(\hat{m}_C = m)$ 



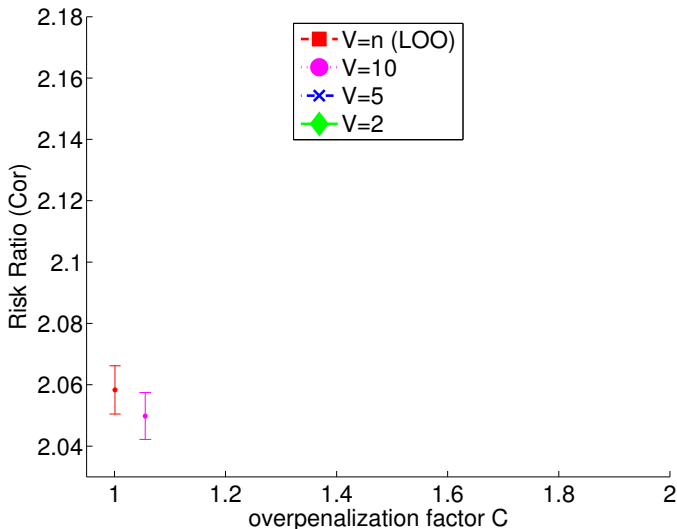
# Évaluation de l'heuristique (2)



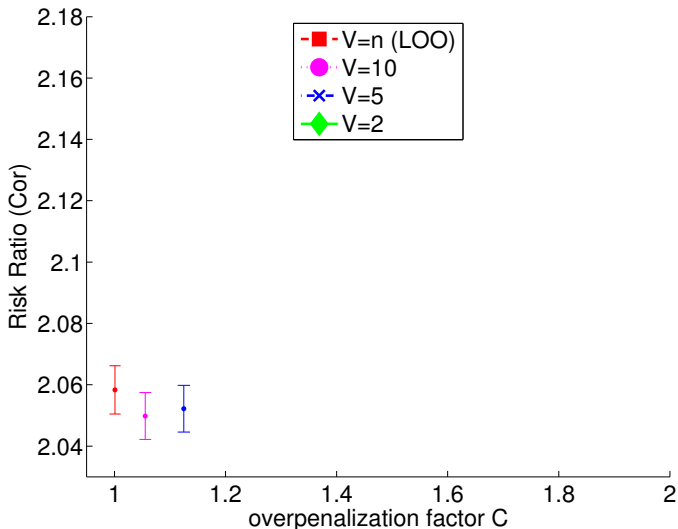
# Simulations : validation croisée « V-fold »



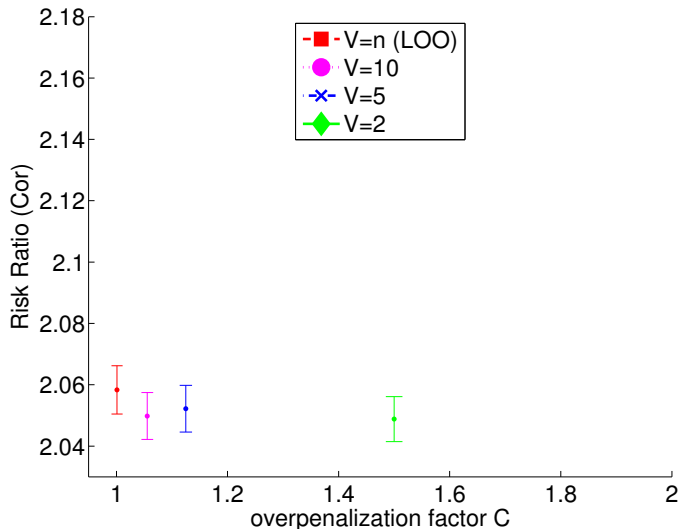
## Simulations : validation croisée « V-fold »

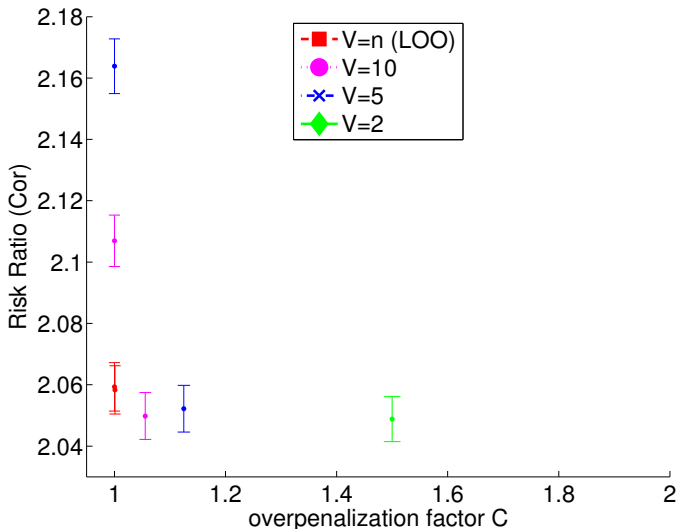


## Simulations : validation croisée « V-fold »

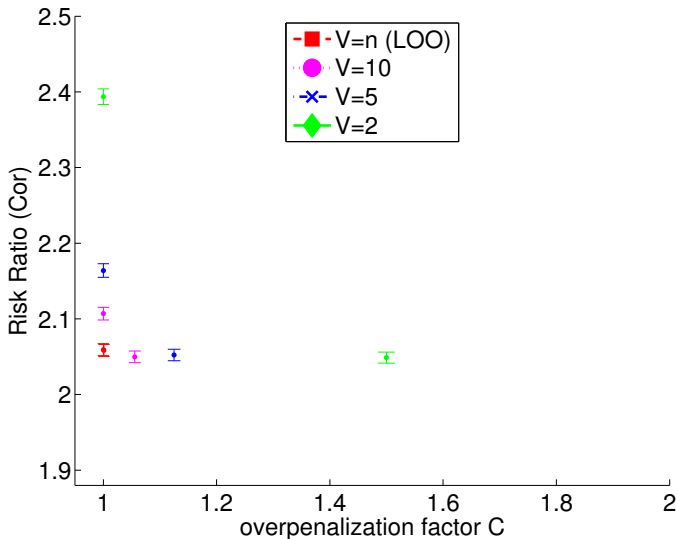


## Simulations : validation croisée « V-fold »

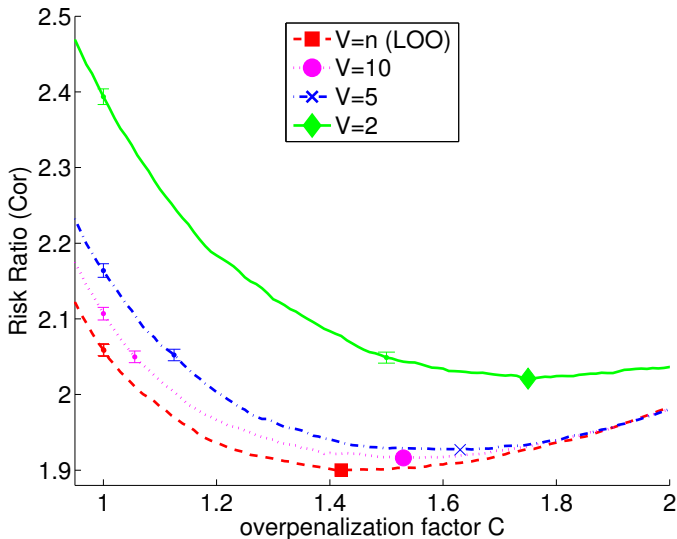


Simulations : pénalisation  $\ll V\text{-fold} \gg$ 

## Simulations : pénalisation « V-fold »

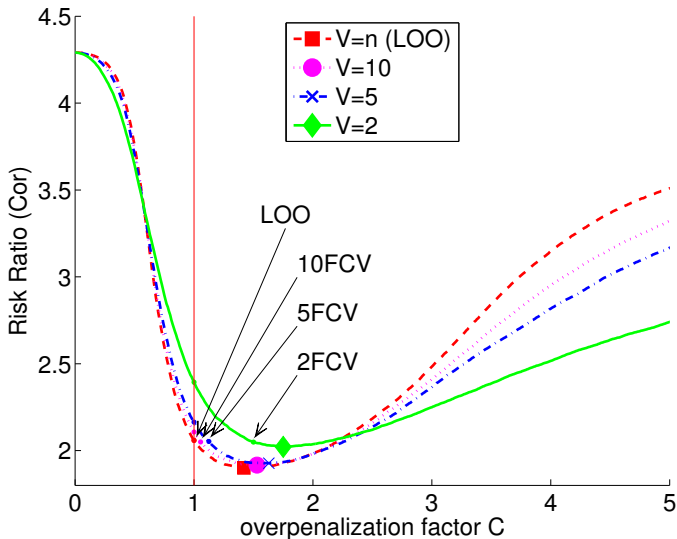


# Simulations : surpénalisation

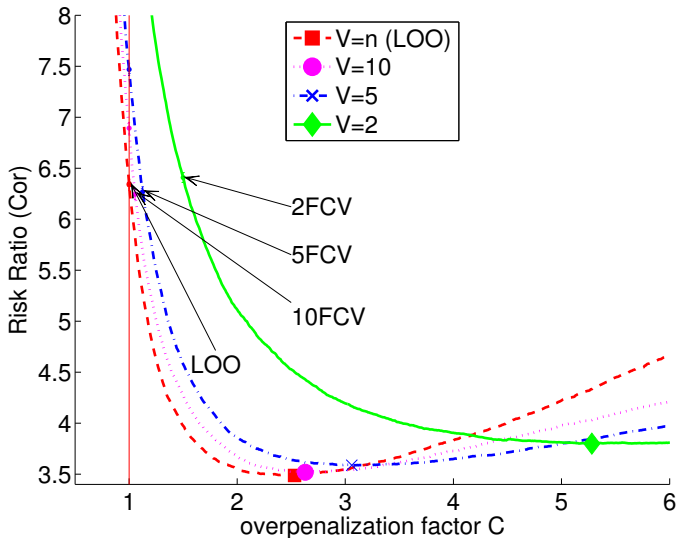




## Simulations : conclusion



## Simulations : cadre « pratiquement paramétrique »



## « V-fold » et sélection de modèles : conclusion

- Temps de calcul :  $\mathcal{O}(V)$  en général

# « V-fold » et sélection de modèles : conclusion

- **Temps de calcul** :  $\mathcal{O}(V)$  en général
- **Validation croisée « V-fold »** :
  - **Biais** : diminue avec  $V$  / peut être supprimé
  - **Variance** : diminue avec  $V$  / quasi minimal pour  $V \in [5, 10]$

⇒ performance optimale lorsque  $V$  est maximal, **quasi optimale pour  $V = 10$** ...

# « V-fold » et sélection de modèles : conclusion

- Temps de calcul :  $\mathcal{O}(V)$  en général
  - Validation croisée « V-fold » :
    - Biais : diminue avec  $V$  / peut être supprimé
    - Variance : diminue avec  $V$  / quasi minimal pour  $V \in [5, 10]$
- ⇒ performance optimale lorsque  $V$  est maximal, **quasi optimale pour  $V = 10$** ...
- ... **si** le facteur de surpénalisation optimal  $C^* \approx 1$  (**nombreux cas possibles**).

## « V-fold » et sélection de modèles : conclusion

- **Temps de calcul** :  $\mathcal{O}(V)$  en général
- **Validation croisée « V-fold »** :
  - Biais : diminue avec  $V$  / peut être supprimé
  - Variance : diminue avec  $V$  / quasi minimal pour  $V \in [5, 10]$
 ⇒ performance optimale lorsque  $V$  est maximal, quasi optimale pour  $V = 10$ ...  
 ... si le facteur de surpénalisation optimal  $C^* \approx 1$  (nombreux cas possibles).
- **Pénalisation « V-fold »** :
  - **Découplage** entre biais et variance ⇒ plus simple à comprendre et utiliser.
  - Biais : **choisi directement** à travers  $C$ , sans contrainte.
  - Variance : diminue avec  $V$  / quasi minimale pour  $V \in [5, 10]$ .

## Généralité de ces résultats

- Valable en régression par moindres carrés et en estimation de densité par moindres carrés ou noyaux (travail en cours avec M. Lerasle et N. Magalhães).

# Généralité de ces résultats

- Valable en régression par moindres carrés et en estimation de densité par moindres carrés ou noyaux (travail en cours avec M. Lerasle et N. Magalhães).
- Correction du biais / pénalisation «  $V$ -fold » : valable lorsque

$$\mathbb{E} \left[ (P - P_n) \gamma(\hat{s}_m) \right] \approx \frac{\gamma(m)}{n} .$$

Sinon : «  $V$ -fold » répété ou VC Monte-Carlo avec  $n_e$  bien choisi.



# Généralité de ces résultats

- Valable en régression par moindres carrés et en estimation de densité par moindres carrés ou noyaux (travail en cours avec M. Lerasle et N. Magalhães).
- Correction du biais / pénalisation «  $V$ -fold » : valable lorsque

$$\mathbb{E}\left[(P - P_n)\gamma(\hat{s}_m)\right] \approx \frac{\gamma(m)}{n} .$$

Sinon : «  $V$ -fold » répété ou VC Monte-Carlo avec  $n_e$  bien choisi.

- **Variance : d'autres comportements sont possibles (expérimentalement).**

# Généralité de ces résultats

- Valable en régression par moindres carrés et en estimation de densité par moindres carrés ou noyaux (travail en cours avec M. Lerasle et N. Magalhães).
- Correction du biais / pénalisation « V-fold » : valable lorsque

$$\mathbb{E}\left[(P - P_n)\gamma(\hat{s}_m)\right] \approx \frac{\gamma(m)}{n} .$$

Sinon : « V-fold » répété ou VC Monte-Carlo avec  $n_e$  bien choisi.

- Variance : d'autres comportements sont possibles (expérimentalement).
- **Tout peut se tester sur des données synthétiques : tracer**

$$n \rightarrow \mathbb{E}\left[P\gamma(\hat{s}_m(D_n))\right] \quad \text{et} \quad m \rightarrow \text{var}\left(\hat{\mathcal{R}}^{\text{vc}}(\hat{s}_m) - \hat{\mathcal{R}}^{\text{vc}}(\hat{s}_{m^*})\right) .$$