

Contributions à la théorie statistique de l'apprentissage: sélection d'estimateurs et détection de ruptures

Sylvain Arlot

¹CNRS

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

Soutenance d'habilitation à diriger des recherches, 3 décembre
2014

Parcours

- Thèse et monitorat à l'Université Paris-Sud (2004 à 2008)
 - Séjour à Berkeley (fév.-mars 2008, P. Bartlett)

Parcours

- Thèse et monitorat à l'Université Paris-Sud (2004 à 2008)
 - Séjour à Berkeley (fév.-mars 2008, P. Bartlett)
- Chargé de recherches CNRS, affecté au Département d'Informatique de l'École normale supérieure
 - Équipe Willow (2008 à 2010)
 - Équipe Sierra (2011 à aujourd'hui)

Pour une théorie utile en pratique

- Comprendre pourquoi certaines procédures sont performantes

Pour une théorie utile en pratique

- Comprendre pourquoi certaines procédures sont performantes
- Pourquoi certaines procédures fonctionnent-elles mieux ?
Compromis entre complexité algorithmique et performance statistique ?

Pour une théorie utile en pratique

- Comprendre **pourquoi certaines procédures sont performantes**
- Pourquoi certaines procédures fonctionnent-elles mieux ?
Compromis entre complexité algorithmique et performance statistique ?
- **Corriger les défauts** de méthodes couramment utilisées

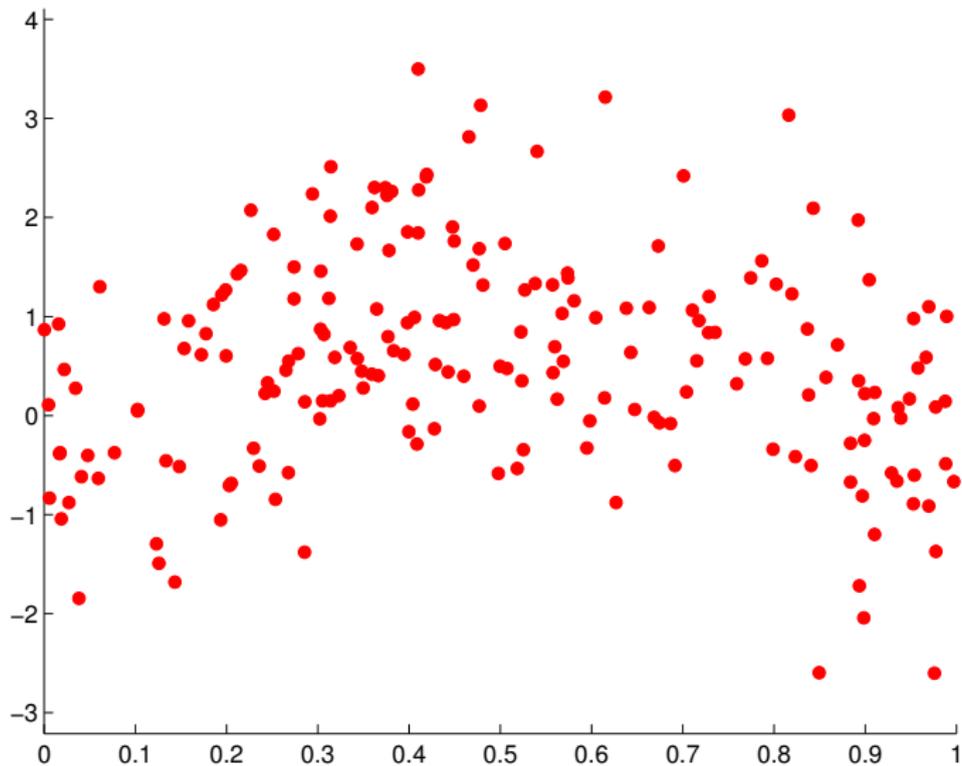
Pour une théorie utile en pratique

- Comprendre **pourquoi certaines procédures sont performantes**
- Pourquoi certaines procédures fonctionnent-elles mieux ?
Compromis entre complexité algorithmique et performance statistique ?
- **Corriger les défauts** de méthodes couramment utilisées
- Proposer de **nouvelles méthodes** sur des bases théoriques e.g., heuristique de pente (Birgé & Massart, 2001)

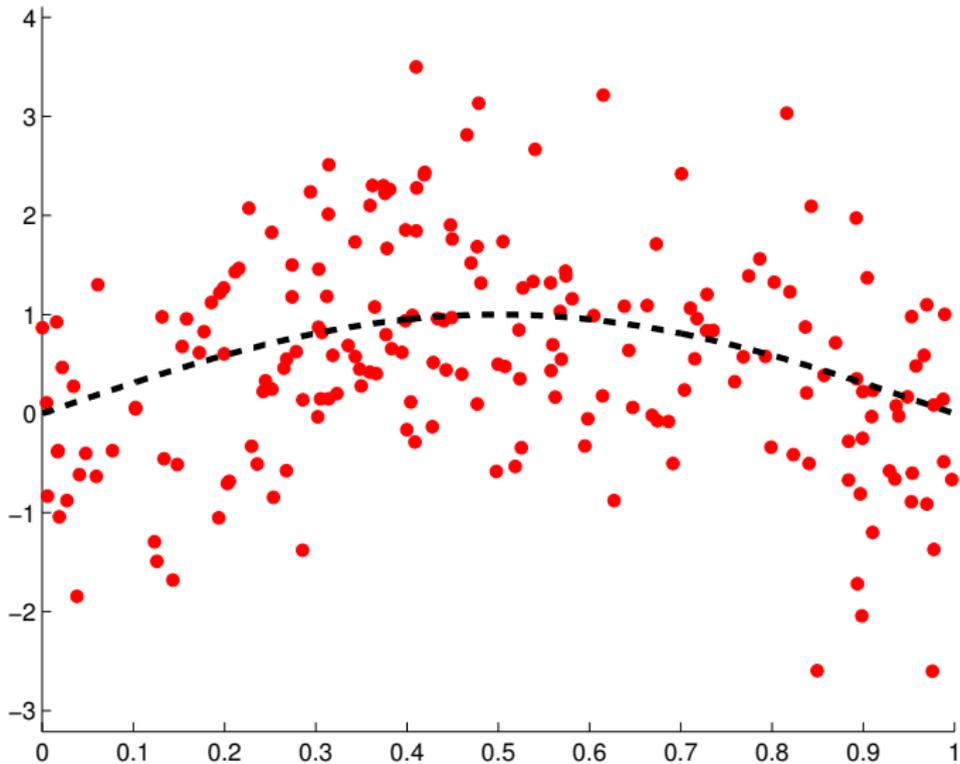
Plan

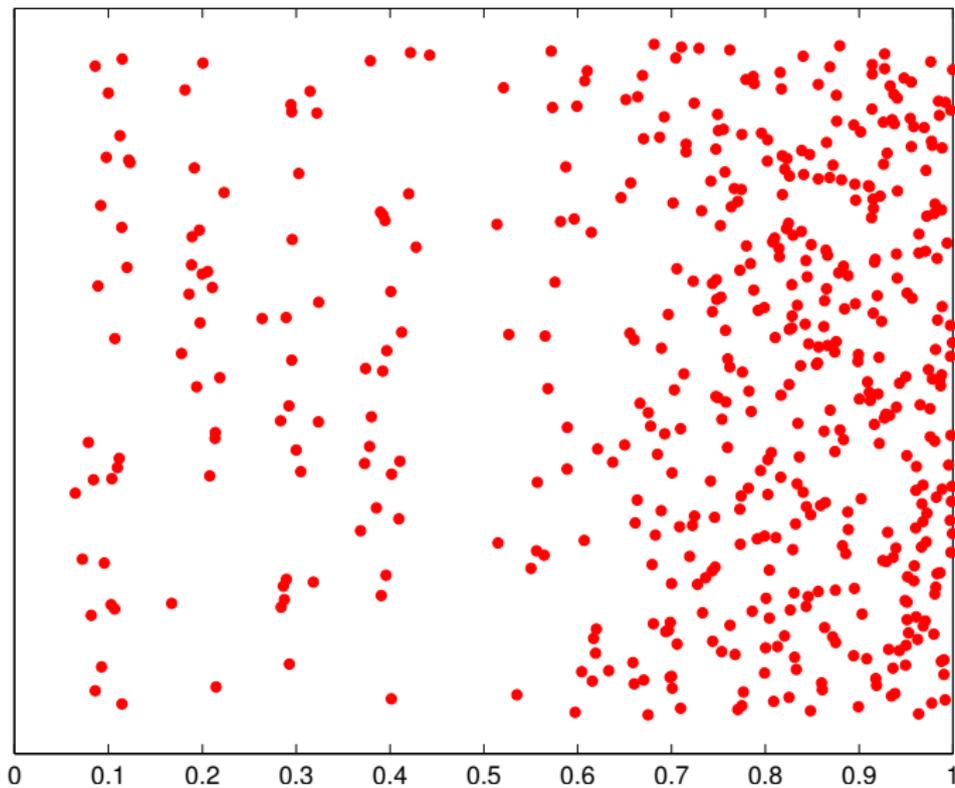
- 1 Sélection d'estimateurs
- 2 Validation croisée
- 3 Pénalités minimales
- 4 Détection de ruptures
- 5 Conclusion

Régression : données $(X_1, Y_1), \dots, (X_n, Y_n)$

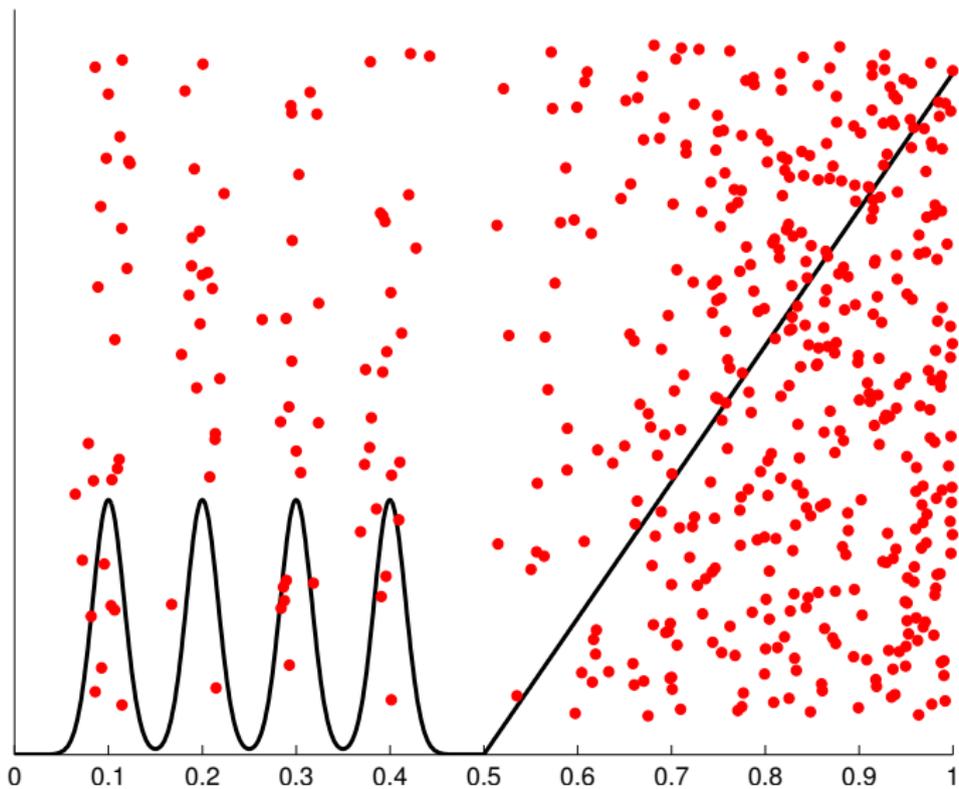


But : reconstruire le signal



Estimation de densité : données ξ_1, \dots, ξ_n 

But : estimer la densité s^* des observations ξ_i



Cadre général

- Données $\xi_1, \dots, \xi_n \in \Xi$ i.i.d. de loi P
prédiction : $\xi_j = (X_j, Y_j) \in \mathcal{X} \times \mathcal{Y}$

Cadre général

- Données $\xi_1, \dots, \xi_n \in \Xi$ i.i.d. de loi P
prédiction : $\xi_i = (X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$
- But : estimer une caractéristique $s^* \in \mathcal{S}$ de P
densité, fonction de régression, meilleur prédicteur, etc.

Cadre général

- Données $\xi_1, \dots, \xi_n \in \Xi$ i.i.d. de loi P
prédiction : $\xi_i = (X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$
- But : estimer une caractéristique $s^* \in \mathbb{S}$ de P
densité, fonction de régression, meilleur prédicteur, etc.
- **Fonction de contraste $\gamma : \mathbb{S} \times \Xi \rightarrow \mathbb{R}$ telle que**

$$s^* \in \underset{t \in \mathbb{S}}{\operatorname{argmin}} \{ P\gamma(t) \} \quad \text{avec} \quad P\gamma(t) := \mathbb{E}_{\xi \sim P} [\gamma(t; \xi)]$$

Cadre général

- Données $\xi_1, \dots, \xi_n \in \Xi$ i.i.d. de loi P
prédiction : $\xi_i = (X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$
- But : estimer une caractéristique $s^* \in \mathbb{S}$ de P
densité, fonction de régression, meilleur prédicteur, etc.
- Fonction de contraste $\gamma : \mathbb{S} \times \Xi \rightarrow \mathbb{R}$ telle que

$$s^* \in \operatorname{argmin}_{t \in \mathbb{S}} \{P\gamma(t)\} \quad \text{avec} \quad P\gamma(t) := \mathbb{E}_{\xi \sim P} [\gamma(t; \xi)]$$

- Perte relative

$$\ell(s^*, t) := P\gamma(t) - P\gamma(s^*) \geq 0$$

Exemples

- Prédiction : $\xi_i = (X_i, Y_i)$

$X_{n+1} \rightsquigarrow$ « prédire » Y_{n+1} avec $t(X_{n+1})$?

$\gamma(t; (x, y))$ mesure la « distance » entre $t(x)$ et y

Exemples

- Prédiction : $\xi_i = (X_i, Y_i)$

$X_{n+1} \rightsquigarrow$ « prédire » Y_{n+1} avec $t(X_{n+1})$?

$\gamma(t; (x, y))$ mesure la « distance » entre $t(x)$ et y

- Régression ($\mathcal{Y} = \mathbb{R}$), moindres carrés :

$$\gamma(t; (x, y)) = (t(x) - y)^2 \quad s^*(X) = \mathbb{E}[Y|X]$$

Exemples

- Prédiction : $\xi_i = (X_i, Y_i)$

$X_{n+1} \rightsquigarrow$ « prédire » Y_{n+1} avec $t(X_{n+1})$?

$\gamma(t; (x, y))$ mesure la « distance » entre $t(x)$ et y

- Régression ($\mathcal{Y} = \mathbb{R}$), moindres carrés :

$$\gamma(t; (x, y)) = (t(x) - y)^2 \quad s^*(X) = \mathbb{E}[Y|X]$$

- Classification binaire ($\mathcal{Y} = \{0, 1\}$), perte 0-1 :

$$\gamma(t; (x, y)) = \mathbb{1}_{t(x) \neq y}$$

Exemples

- Prédiction : $\xi_i = (X_i, Y_i)$

$X_{n+1} \rightsquigarrow \ll \text{prédire} \gg Y_{n+1}$ avec $t(X_{n+1})$?

$\gamma(t; (x, y))$ mesure la « distance » entre $t(x)$ et y

- Régression ($\mathcal{Y} = \mathbb{R}$), moindres carrés :

$$\gamma(t; (x, y)) = (t(x) - y)^2 \quad s^*(X) = \mathbb{E}[Y|X]$$

- Classification binaire ($\mathcal{Y} = \{0, 1\}$), perte 0–1 :

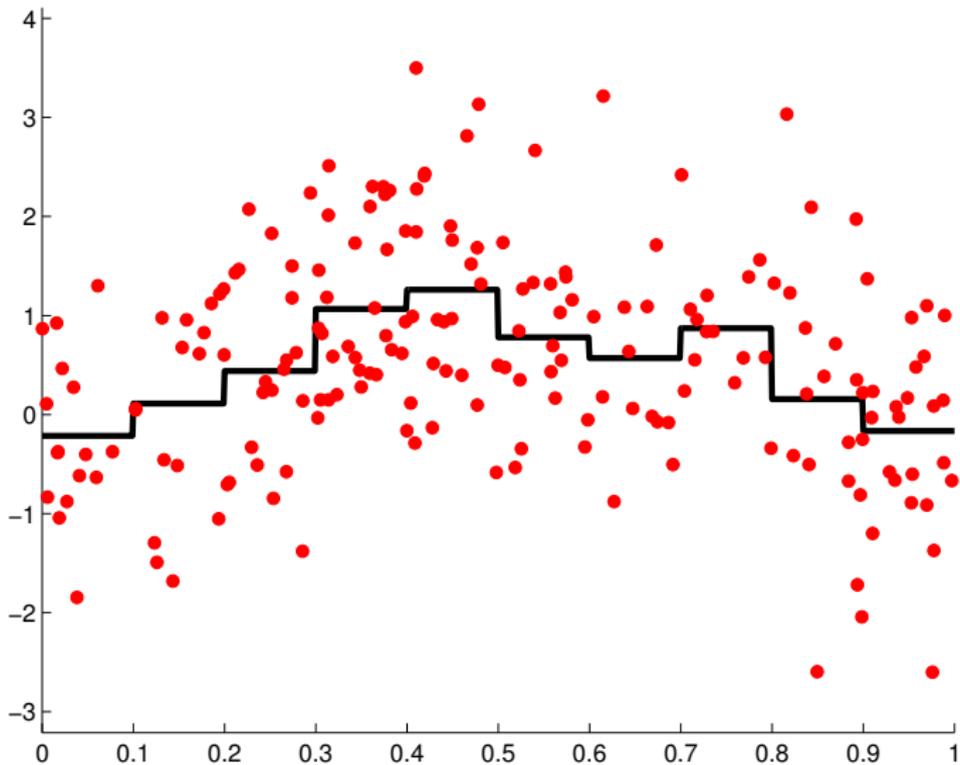
$$\gamma(t; (x, y)) = \mathbb{1}_{t(x) \neq y}$$

- Estimation de densité (mesure de référence μ) :

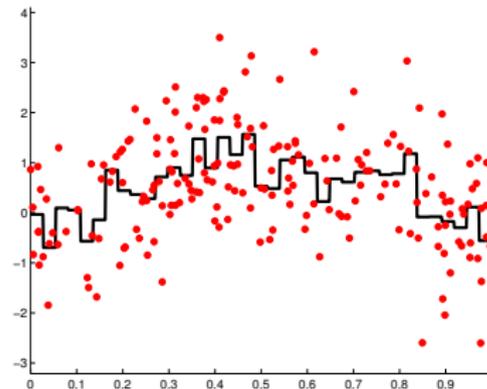
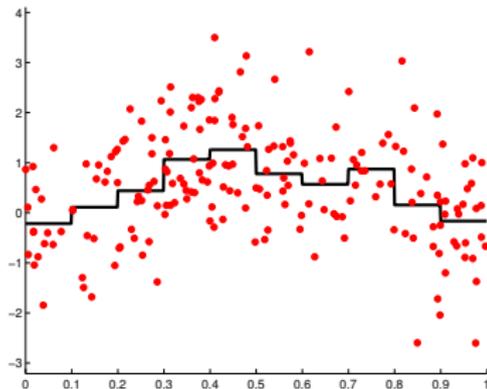
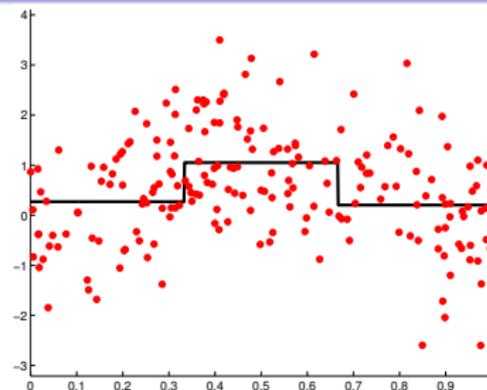
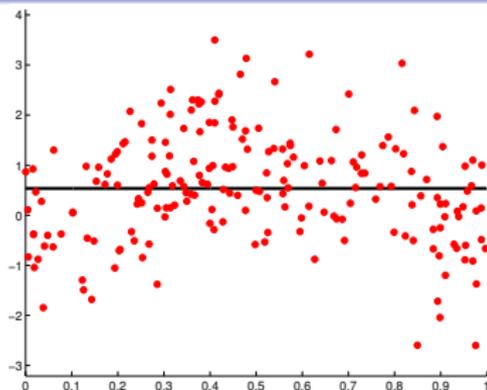
moindres carrés : $\gamma(t; \xi) = \|t\|_{L^2(\mu)}^2 - 2t(\xi)$

log-vraisemblance : $\gamma(t; \xi) = -\log(t(\xi))$

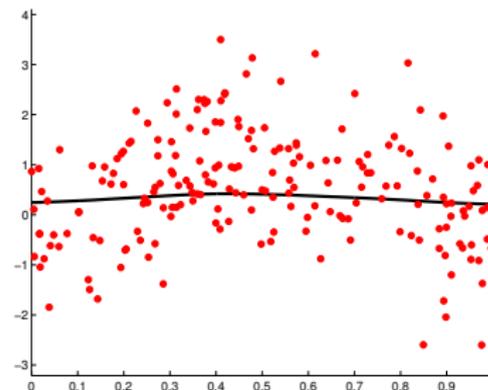
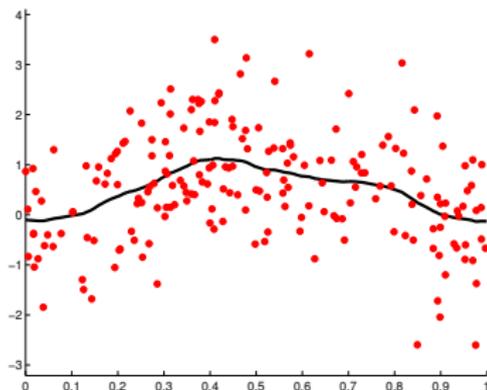
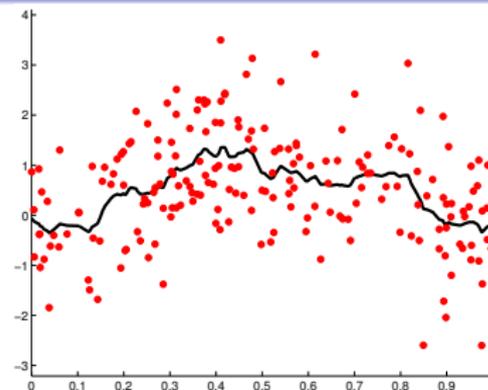
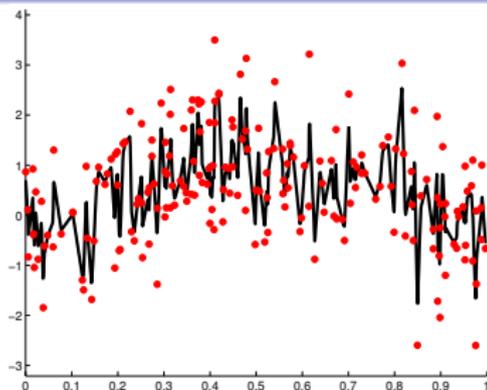
Estimateurs : un régressogramme



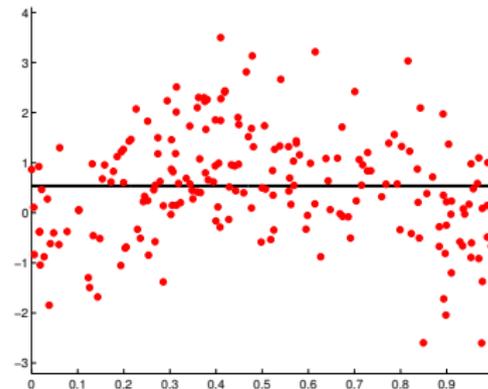
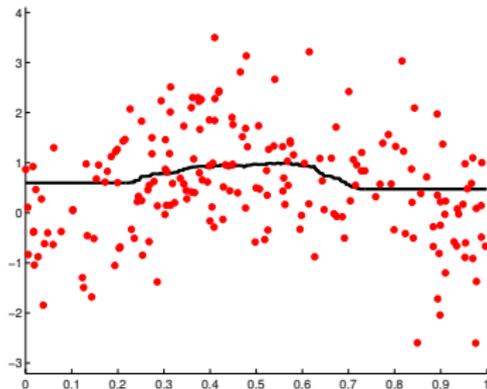
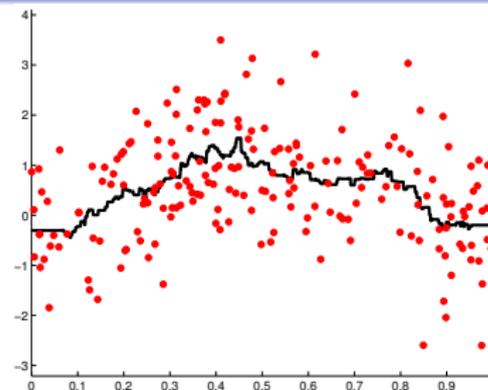
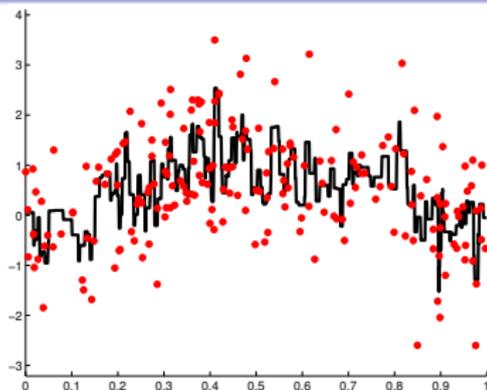
Sélection d'estimateurs : régressogrammes réguliers



Sélection d'estimateurs : régression ridge à noyau



Sélection d'estimateurs : k plus proches voisins



Sélection d'estimateurs

- **Estimateur/Algorithme d'apprentissage** : $\hat{s} : D_n \mapsto \hat{s}(D_n) \in \mathbb{S}$

Sélection d'estimateurs

- Estimateur/Algorithme d'apprentissage : $\hat{s} : D_n \mapsto \hat{s}(D_n) \in \mathbb{S}$
- Exemple : **estimateur des moindres carrés** sur un modèle $S_m \subset \mathbb{S}$

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

Exemple de modèle : histogrammes

Sélection d'estimateurs

- Estimateur/Algorithme d'apprentissage : $\hat{s} : D_n \mapsto \hat{s}(D_n) \in \mathbb{S}$
- Exemple : estimateur des moindres carrés sur un modèle $S_m \subset \mathbb{S}$

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

Exemple de modèle : histogrammes

- Famille d'estimateurs $(\hat{s}_m)_{m \in \mathcal{M}} \Rightarrow$ choisir $\hat{m} = \hat{m}(D_n)$?
e.g., famille de modèles $(S_m)_{m \in \mathcal{M}} \Rightarrow$ famille d'estimateurs des moindres carrés

Sélection d'estimateurs

- Estimateur/Algorithme d'apprentissage : $\hat{s} : D_n \mapsto \hat{s}(D_n) \in \mathbb{S}$
- Exemple : estimateur des moindres carrés sur un modèle $S_m \subset \mathbb{S}$

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

Exemple de modèle : histogrammes

- Famille d'estimateurs $(\hat{s}_m)_{m \in \mathcal{M}} \Rightarrow$ choisir $\hat{m} = \hat{m}(D_n)$?
e.g., famille de modèles $(S_m)_{m \in \mathcal{M}} \Rightarrow$ famille d'estimateurs des moindres carrés
- Objectif : minimiser le risque, *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

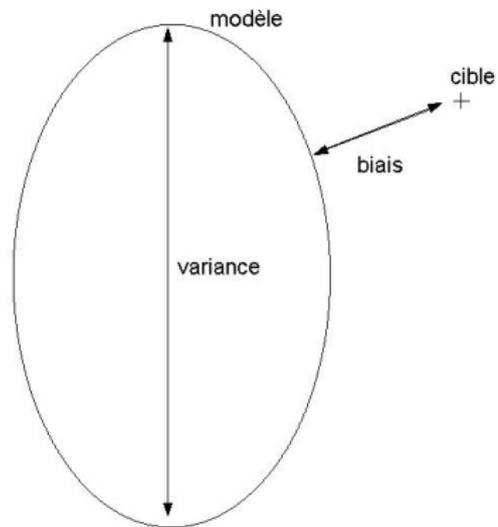
$$\mathbb{E}[\ell(s^*, \hat{s}_m)] = \text{Biais} + \text{Variance}$$

Biais ou Erreur d'approximation

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

Variance ou Erreur d'estimation

Régression, moindres carrés : $\frac{\sigma^2 \dim(S_m)}{n}$



Compromis biais-variance

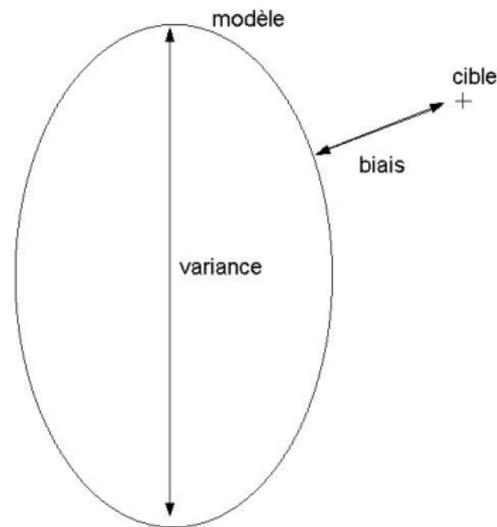
$$\mathbb{E}[\ell(s^*, \hat{s}_m)] = \text{Biais} + \text{Variance}$$

Biais ou Erreur d'approximation

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

Variance ou Erreur d'estimation

Régression, moindres carrés : $\frac{\sigma^2 \dim(S_m)}{n}$



Compromis biais-variance

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

Un problème général

- Objectif :

minimiser $\mathcal{R}(m)$ sur $m \in \mathcal{M}$

Un problème général

- Objectif :

minimiser $\mathcal{R}(m)$ sur $m \in \mathcal{M}$

- Méthode :

minimiser $\mathcal{C}(m)$ sur $m \in \mathcal{M} \Rightarrow \hat{m}_{\mathcal{C}}$

Un problème général

- Objectif :

minimiser $\mathcal{R}(m)$ sur $m \in \mathcal{M}$

- Méthode :

minimiser $\mathcal{C}(m)$ sur $m \in \mathcal{M} \Rightarrow \hat{m}_{\mathcal{C}}$

- Exemples :

- sélection d'estimateurs :

$$\mathcal{R}(m) = \ell(s^*, \hat{s}_m)$$

Un problème général

- Objectif :

minimiser $\mathcal{R}(m)$ sur $m \in \mathcal{M}$

- Méthode :

minimiser $\mathcal{C}(m)$ sur $m \in \mathcal{M} \Rightarrow \hat{m}_{\mathcal{C}}$

- Exemples :

- sélection d'estimateurs :

$$\mathcal{R}(m) = \ell(s^*, \hat{s}_m)$$

- estimateur par minimum de contraste sur S :

$$\mathcal{M} = S \subset \mathbb{S} \quad \mathcal{R}(t) = \ell(s^*, t) \quad \mathcal{C}(t) = P_n \gamma(t)$$

Un problème général

- Objectif :

minimiser $\mathcal{R}(m)$ sur $m \in \mathcal{M}$

- Méthode :

minimiser $\mathcal{C}(m)$ sur $m \in \mathcal{M} \Rightarrow \hat{m}_{\mathcal{C}}$

- Exemples :

- sélection d'estimateurs :

$$\mathcal{R}(m) = \ell(s^*, \hat{s}_m)$$

- estimateur par minimum de contraste sur S :

$$\mathcal{M} = S \subset \mathbb{S} \quad \mathcal{R}(t) = \ell(s^*, t) \quad \mathcal{C}(t) = P_n \gamma(t)$$

- relaxations (convexes) en optimisation

Analyse du problème : un lemme

Lemme

Si

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

alors,

$$\begin{aligned} \forall \hat{m}_C \in \operatorname{argmin}_{m \in \mathcal{M}} \{C(m)\}, \\ \mathcal{R}(\hat{m}_C) - B(\hat{m}_C) \leq \inf_{m \in \mathcal{M}} \{\mathcal{R}(m) + A(m)\}. \end{aligned} \quad (1)$$

Analyse du problème : un lemme

Lemme

Si

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

alors,

$$\begin{aligned} \forall \hat{m}_C \in \operatorname{argmin}_{m \in \mathcal{M}} \{C(m)\}, \\ \mathcal{R}(\hat{m}_C) - B(\hat{m}_C) \leq \inf_{m \in \mathcal{M}} \{ \mathcal{R}(m) + A(m) \}. \end{aligned} \quad (1)$$

En fait, (1) a lieu dès que

$$\begin{aligned} \forall m, m' \in \mathcal{M}, \quad (C(m) - \mathcal{R}(m)) - (C(m') - \mathcal{R}(m')) \\ \leq A(m) + B(m'). \end{aligned}$$

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 $\delta_n \rightarrow 0$.

Une autre utilisation du lemme

- Choisir \mathcal{C} qui est une borne supérieure sur \mathcal{R} , uniformément sur $m \in \mathcal{M}$, c'est-à-dire tel que

$$\forall m \in \mathcal{M}, \quad \mathcal{C}(m) \geq \mathcal{R}(m) .$$

Une autre utilisation du lemme

- Choisir \mathcal{C} qui est une **borne supérieure sur \mathcal{R} , uniformément sur $m \in \mathcal{M}$** , c'est-à-dire tel que

$$\forall m \in \mathcal{M}, \quad \mathcal{C}(m) \geq \mathcal{R}(m) .$$

⇒ d'après le lemme,

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

Une autre utilisation du lemme

- Choisir \mathcal{C} qui est une **borne supérieure sur \mathcal{R} , uniformément sur $m \in \mathcal{M}$** , c'est-à-dire tel que

$$\forall m \in \mathcal{M}, \quad \mathcal{C}(m) \geq \mathcal{R}(m) .$$

⇒ d'après le lemme,

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

- Exemples :
 - **grandes collections d'estimateurs (sélection de variables, détection de ruptures, etc.)**

Une autre utilisation du lemme

- Choisir \mathcal{C} qui est une **borne supérieure sur \mathcal{R} , uniformément sur $m \in \mathcal{M}$** , c'est-à-dire tel que

$$\forall m \in \mathcal{M}, \quad \mathcal{C}(m) \geq \mathcal{R}(m) .$$

⇒ d'après le lemme,

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

- Exemples :
 - grandes collections d'estimateurs (sélection de variables, détection de ruptures, etc.)
 - **minimisation du risque structurel (Vapnik)**

Une autre utilisation du lemme

- Choisir \mathcal{C} qui est une **borne supérieure sur \mathcal{R} , uniformément sur $m \in \mathcal{M}$** , c'est-à-dire tel que

$$\forall m \in \mathcal{M}, \quad \mathcal{C}(m) \geq \mathcal{R}(m) .$$

⇒ d'après le lemme,

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

- Exemples :
 - grandes collections d'estimateurs (sélection de variables, détection de ruptures, etc.)
 - minimisation du risque structurel (Vapnik)
 - **relaxations en optimisation**

Analyse au second ordre ?

- Comment comparer \mathcal{C}_1 et \mathcal{C}_2 tels que

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

Analyse au second ordre ?

- Comment comparer \mathcal{C}_1 et \mathcal{C}_2 tels que

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

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

Analyse au second ordre ?

- Comment comparer \mathcal{C}_1 et \mathcal{C}_2 tels que

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

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

- Variance de quelle quantité ?

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

Analyse au second ordre ?

- Comment comparer \mathcal{C}_1 et \mathcal{C}_2 tels que

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

- 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 ...

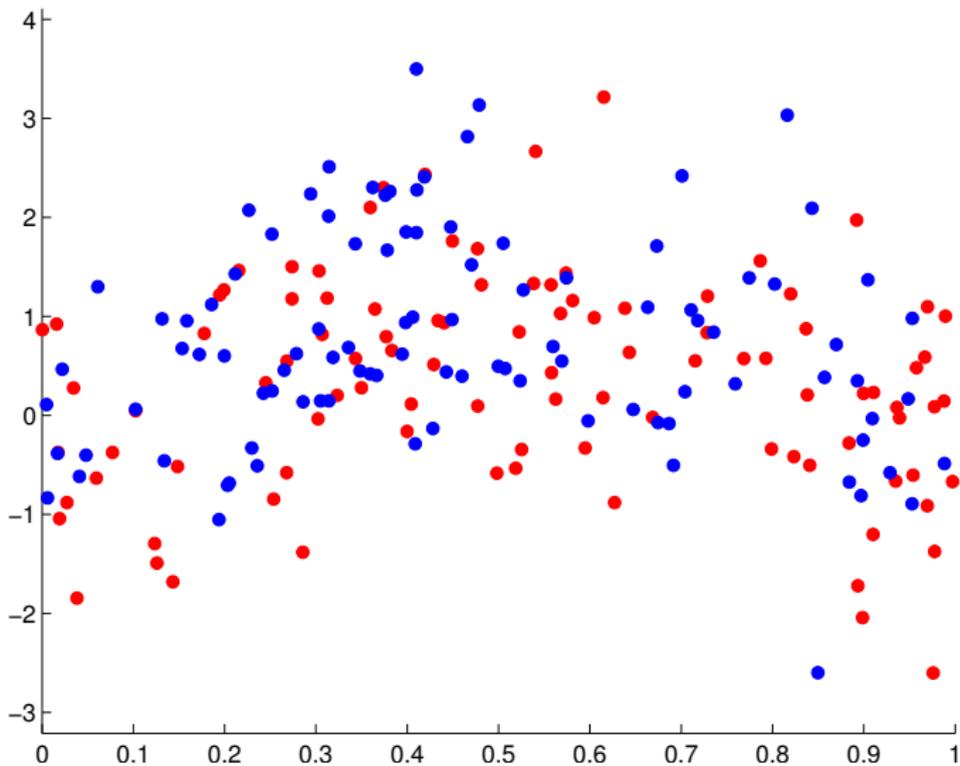
⇒ variance des incréments

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

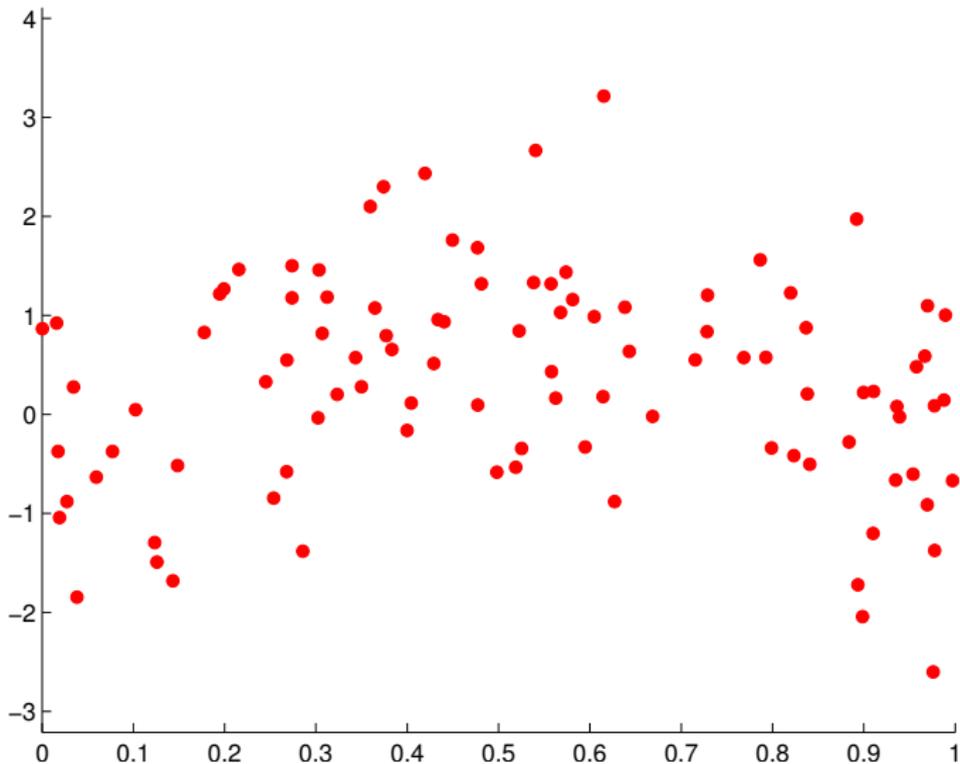
Plan

- 1 Sélection d'estimateurs
- 2 Validation croisée
- 3 Pénalités minimales
- 4 Détection de ruptures
- 5 Conclusion

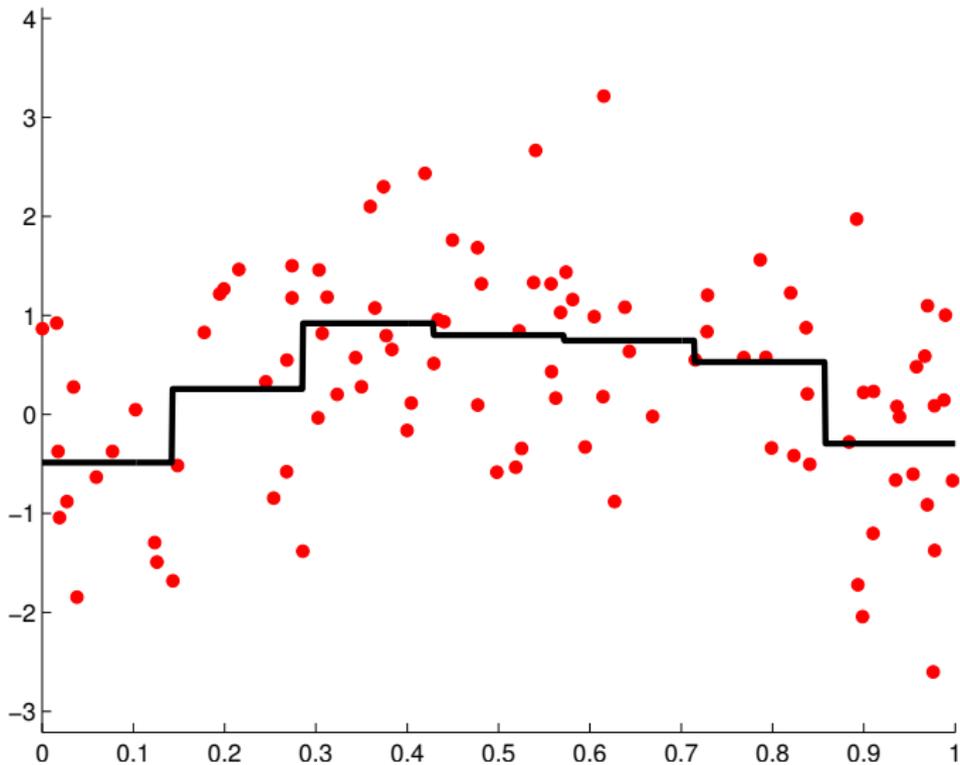
Principe de la validation simple



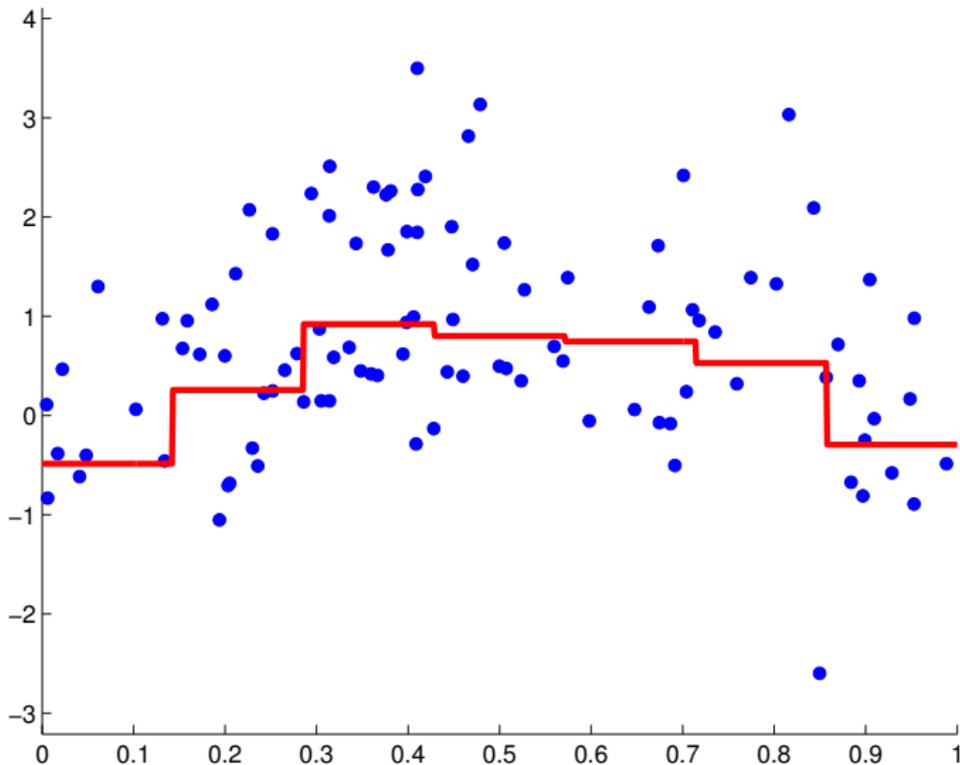
Principe de la validation : échantillon d'entraînement



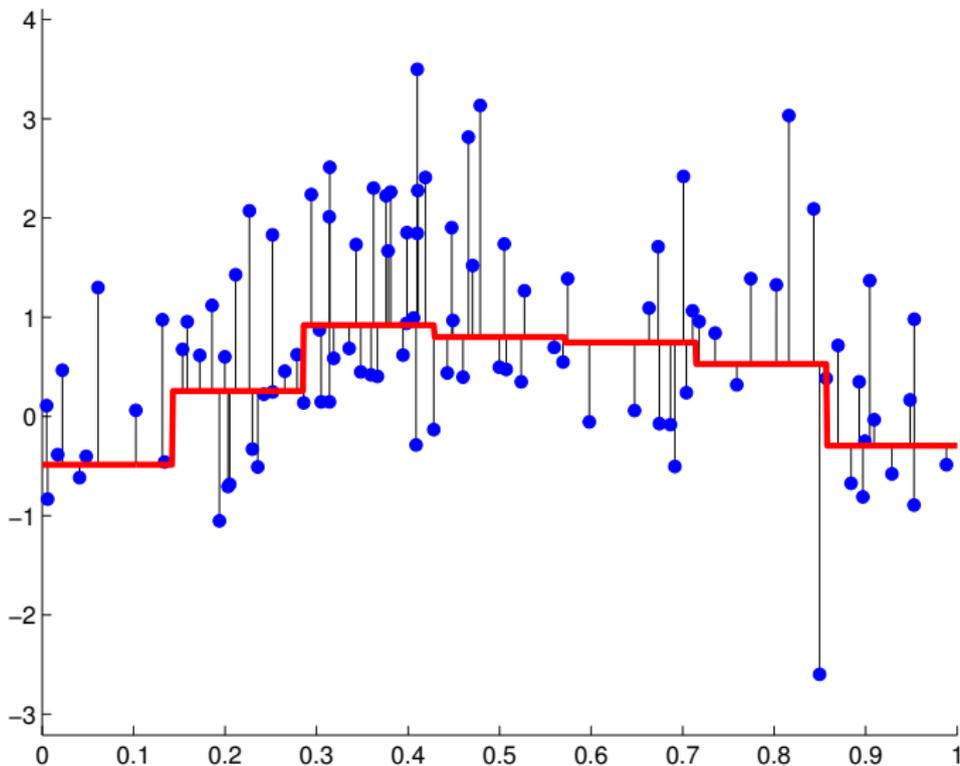
Principe de la validation : échantillon d'entraînement



Principe de la validation : échantillon de validation



Principe de la validation : échantillon de validation



Validation croisée « V-fold »

$$\underbrace{\xi_1, \dots, \xi_q}_{\text{Entraînement}}, \underbrace{\xi_{q+1}, \dots, \xi_n}_{\text{Validation}}$$

Entraînement Validation

Validation croisée « V-fold » :

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

$$\Rightarrow \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m; D_n; \mathcal{B}) = \frac{1}{V} \sum_{j=1}^V P_n^j \gamma(\widehat{s}_m^{(-j)})$$

$$\widehat{m} \in \arg \min_{m \in \mathcal{M}} \left\{ \widehat{\mathcal{R}}^{\text{vf}}(\widehat{s}_m) \right\}$$

A. & Celisse. A survey of cross-validation procedures for model selection. *Statistics Surveys*, 4:40–79, 2010.

Au premier ordre : biais de la validation croisée

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

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

Au premier ordre : biais de la validation croisée

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

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

\Rightarrow **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 :

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

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

⇒ **sous-optimalité** de la validation croisée « V -fold » à V fixé
(prouvé pour les régressogrammes, valable plus largement)

A. V -fold cross-validation improved : V -fold penalization, 2008. arXiv:0802.0566v2

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) :

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

- Heuristique de rééchantillonnage (Efron, 1983), sous-échantillonnage « V-fold » et pénalisation
⇒ pénalité « V-fold »

$$\text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B}) := \frac{V-1}{V} \sum_{j=1}^V (P_n - P_n^{(-j)}) \gamma(\widehat{s}_m^{(-j)})$$

A. V-fold cross-validation improved : V-fold penalization, 2008. arXiv:0802.0566v2

Correction du biais et pénalisation « V-fold »

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

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

- Heuristique de rééchantillonnage (Efron, 1983), sous-échantillonnage « V-fold » et pénalisation
⇒ pénalité « V-fold »

$$\text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B}) := \frac{V-1}{V} \sum_{j=1}^V (P_n - P_n^{(-j)}) \gamma(\widehat{s}_m^{(-j)})$$

A. V-fold cross-validation improved : V-fold penalization, 2008. arXiv:0802.0566v2

- $\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m; D_n; \mathcal{B}) = P_n \gamma(\widehat{s}_m(D_n)) + \text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B})$

Correction du biais et pénalisation « V-fold »

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

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

- Heuristique de rééchantillonnage (Efron, 1983), sous-échantillonnage « V-fold » et pénalisation
⇒ **pénalité « V-fold »**

$$\text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B}) := \frac{V-1}{V} \sum_{j=1}^V (P_n - P_n^{(-j)}) \gamma(\widehat{s}_m^{(-j)})$$

A. V-fold cross-validation improved : V-fold penalization, 2008. arXiv:0802.0566v2

- $\widehat{\mathcal{R}}^{\text{vf,corr}}(\widehat{s}_m; D_n; \mathcal{B}) = P_n \gamma(\widehat{s}_m(D_n)) + \text{pen}_{\text{VF}}(\widehat{s}_m; D_n; \mathcal{B})$
- **Sans biais** si $\mathbb{E}[(P - P_n)(\widehat{s}_m(D_n))] = \gamma(m)/n$

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

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

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

⇒ 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. V-fold cross-validation improved : V-fold penalization, 2008. arXiv:0802.0566v2

A. Model selection by resampling penalization. *Electronic Journal of Statistics*, 3:557–624, 2009.

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. V-fold cross-validation improved : V-fold penalization, 2008. arXiv:0802.0566v2

A. Model selection by resampling penalization. *Electronic Journal of Statistics*, 3:557–624, 2009.

- L'**estimation de densité par moindres carrés**

A. & Lerasle. Why $V = 5$ is enough in V-fold cross-validation, 2014. arXiv:1210.5830v2.

A. & Celisse. A survey of cross-validation procedures for model selection. *Statistics Surveys*, 4:40–79, 2010. 27/68

Variance et sélection de modèles : heuristique

- Heuristique : pour tout $m \notin \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\ell(s^*, \hat{s}_m)]$, on veut minimiser $\mathbb{P}(\hat{m}_C = m)$

Variance et sélection de modèles : heuristique

- **Heuristique** : pour tout $m \notin \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\ell(s^*, \hat{s}_m)]$, on veut minimiser $\mathbb{P}(\hat{m}_{\mathcal{C}} = m)$

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

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

Variance et sélection de modèles : heuristique

- **Heuristique** : pour tout $m \notin \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\ell(s^*, \hat{s}_m)]$, 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

- **Heuristique** : pour tout $m \notin \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\ell(s^*, \hat{s}_m)]$, on veut minimiser $\mathbb{P}(\hat{m}_{\mathcal{C}} = m)$

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

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

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

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

Variance et sélection de modèles : heuristique

- Heuristique : pour tout $m \notin \operatorname{argmin}_{m \in \mathcal{M}} \mathbb{E}[\ell(s^*, \hat{s}_m)]$, on veut minimiser $\mathbb{P}(\hat{m}_{\mathcal{C}} = m)$

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

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

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

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

- Biais constant parmi $(\hat{\mathcal{R}}^{\text{vf,corr}})_{2 \leq v \leq n} \Rightarrow$ comparer

$$\operatorname{var}\left(\hat{\mathcal{R}}^{\text{vf,corr}}(\hat{s}_m; D_n; \mathcal{B}) - \hat{\mathcal{R}}^{\text{vf,corr}}(\hat{s}_{m'}; D_n; \mathcal{B})\right)$$

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

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

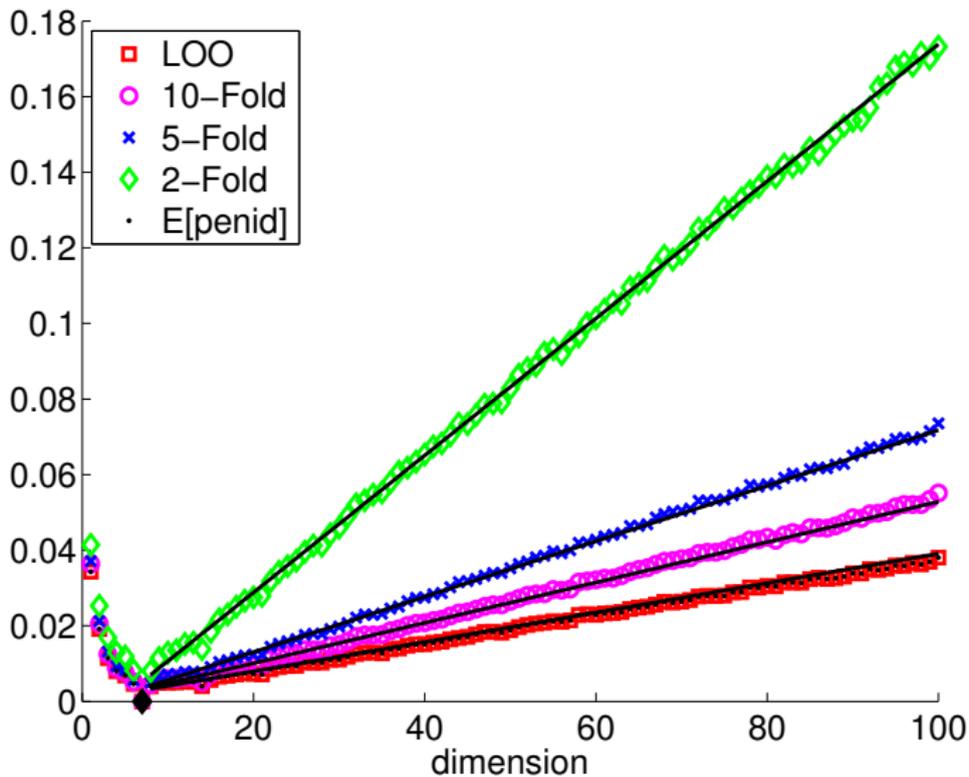
Si de plus $S_m \subset S_{m'}$ sont deux modèles d'histogrammes réguliers de pas d_m^{-1} , $d_{m'}^{-1}$, alors

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

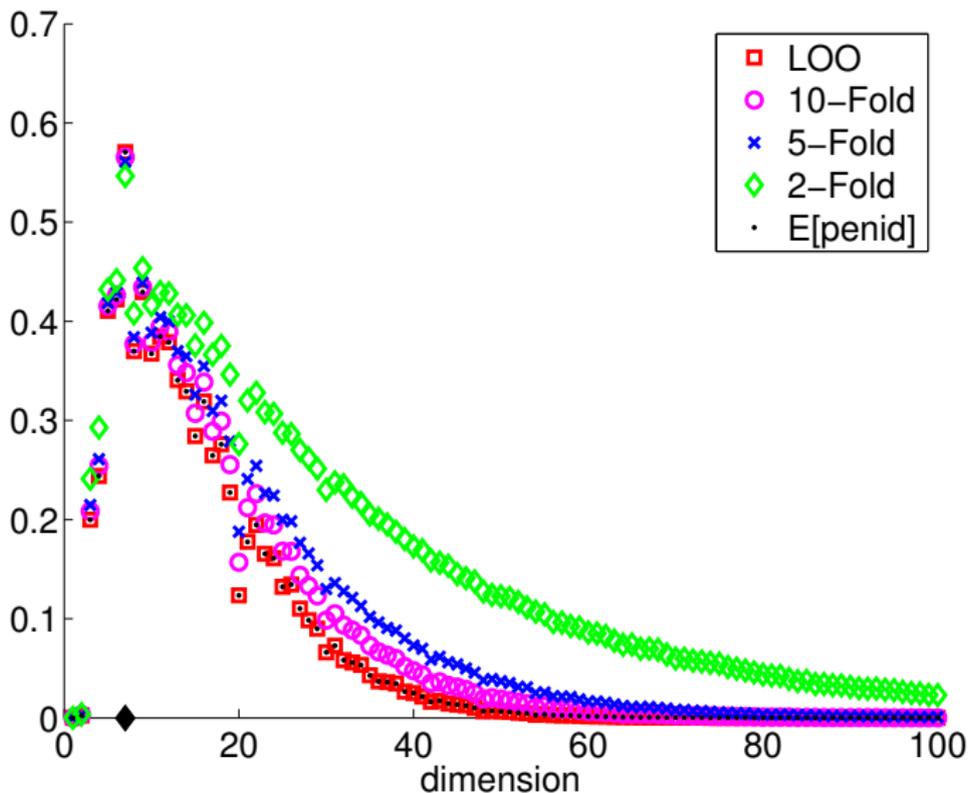
A. & Lerasle. Why $V = 5$ is enough in V -fold cross-validation, 2014. arXiv:1210.5830v2.

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

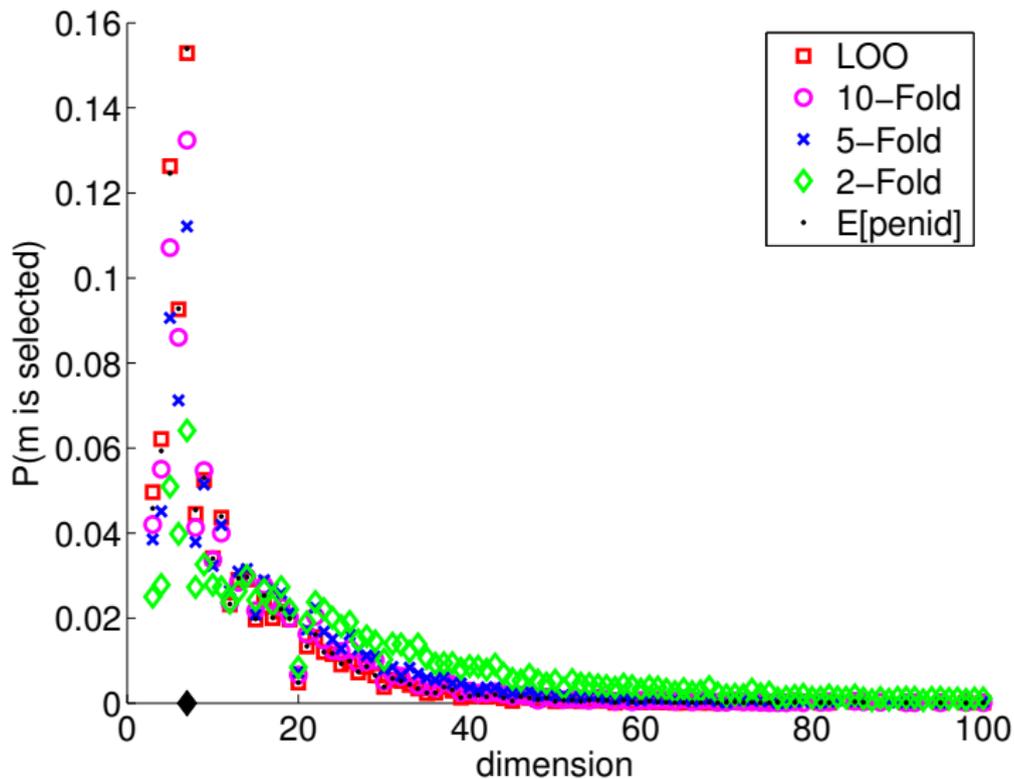
Variance de $\mathcal{C}(m) - \mathcal{C}(m^*)$ en fonction de d_m et V



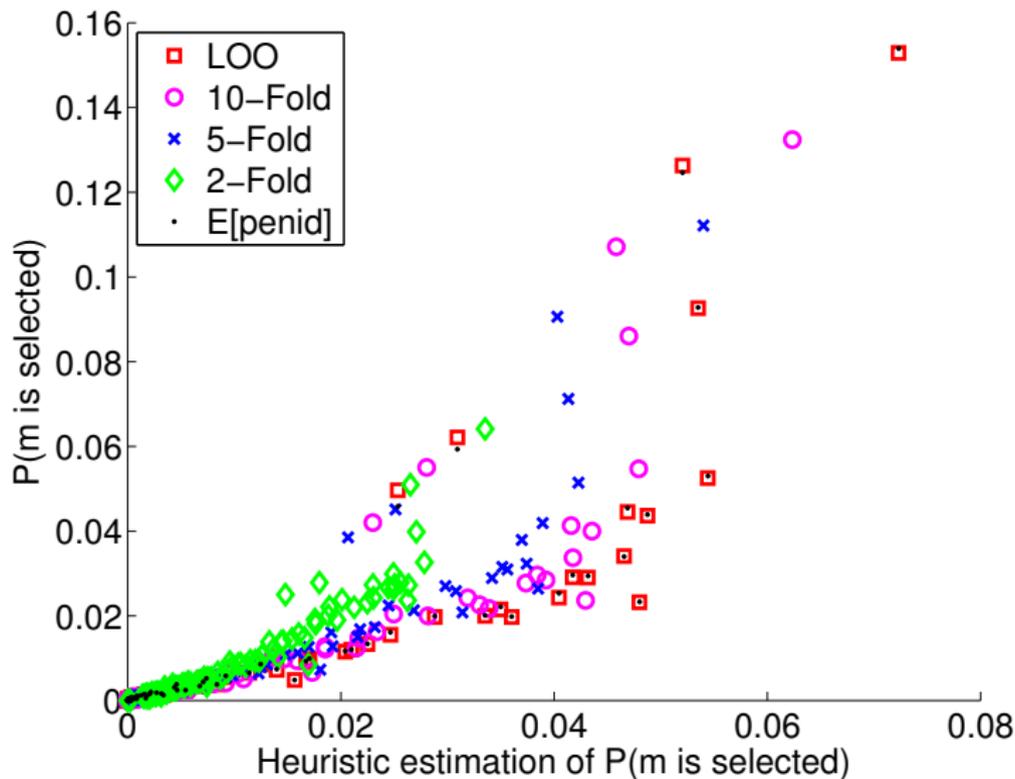
$$\bar{\Phi}(\max_{m' \neq m} \mathbb{E}[\mathcal{C}(m) - \mathcal{C}(m')] / \sqrt{\text{var}(\mathcal{C}(m) - \mathcal{C}(m'))})$$



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



Évaluation de l'heuristique, densité (2)



Conclusion : choix de V pour le « V -fold »

- **Biais** : fonction décroissante de V
peut être supprimé (pénalisation)
 - **Variance** : fonction décroissante de V
proche de son minimum pour $V = 5$ ou 10
- ⇒ meilleure performance pour les plus grands V
(pas nécessairement vrai en général en apprentissage)

Conclusion : choix de V pour le « V -fold »

- **Biais** : fonction décroissante de V
peut être supprimé (pénalisation)
 - **Variance** : fonction décroissante de V
proche de son minimum pour $V = 5$ ou 10
- ⇒ meilleure performance pour les plus grands V
(pas nécessairement vrai en général en apprentissage)
- **Complexité algorithmique** : $\mathcal{O}(V)$ en général

Plan

- 1 Sélection d'estimateurs
- 2 Validation croisée
- 3 Pénalités minimales**
- 4 Détection de ruptures
- 5 Conclusion

Motivation (1) : calibration de pénalités

Pénalisation :

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \{ P_n \gamma(\hat{s}_m) + \text{pen}(m) \}$$

⇒ pénalité idéale $\text{pen}_{\text{id}}(m) = (P - P_n) \gamma(\hat{s}_m)$.

Motivation (1) : calibration de pénalités

Pénalisation :

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \{ P_n \gamma(\hat{s}_m) + \operatorname{pen}(m) \}$$

⇒ pénalité idéale $\operatorname{pen}_{\text{id}}(m) = (P - P_n) \gamma(\hat{s}_m)$.

- C_p et C_L (Mallows, 1973) :

$$\frac{2\sigma^2 D_m}{n}$$

$$\frac{2\sigma^2 \operatorname{tr}(A_m)}{n}$$

Motivation (1) : calibration de pénalités

Pénalisation :

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \{ P_n \gamma(\hat{s}_m) + \operatorname{pen}(m) \}$$

⇒ pénalité idéale $\operatorname{pen}_{\text{id}}(m) = (P - P_n) \gamma(\hat{s}_m)$.

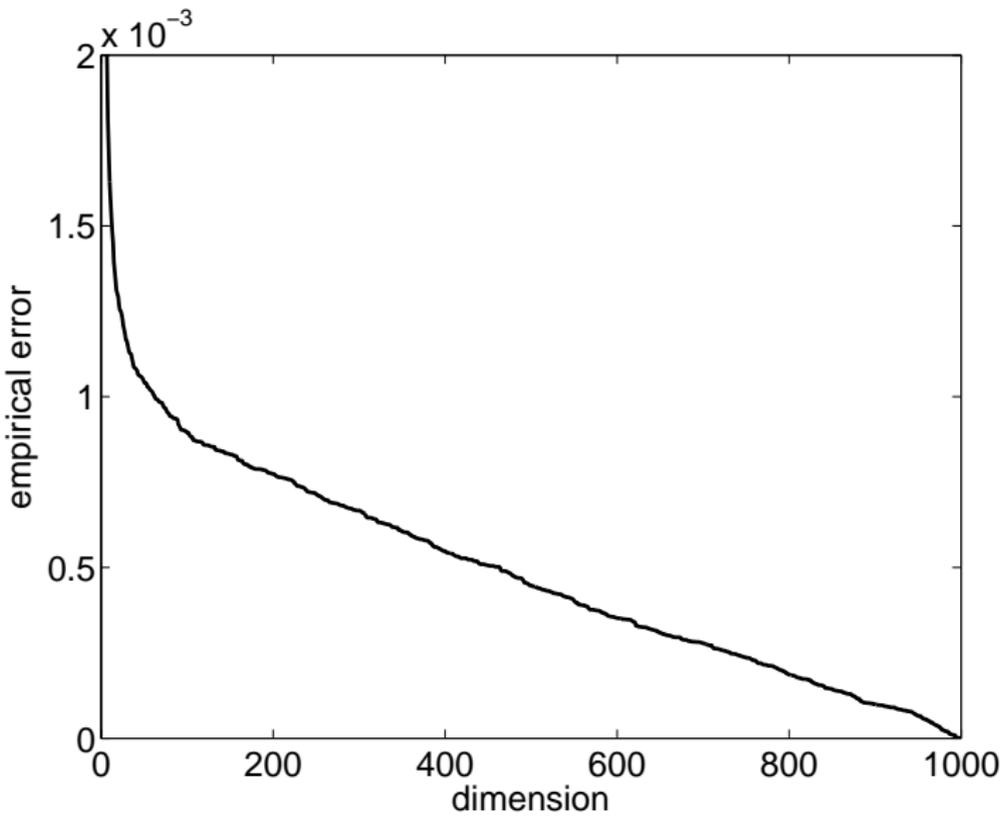
- C_p et C_L (Mallows, 1973) :

$$\frac{2\sigma^2 D_m}{n}$$

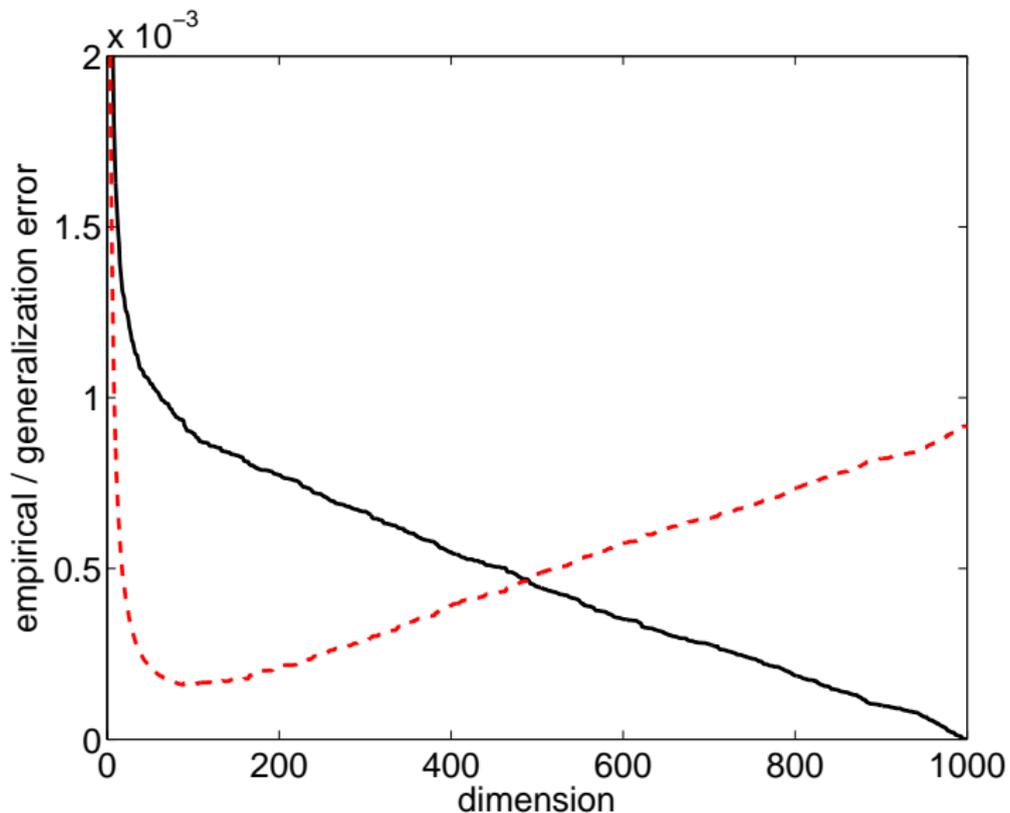
$$\frac{2\sigma^2 \operatorname{tr}(A_m)}{n}$$

- Pénalités justifiées **asymptotiquement** : AIC, BIC
- Pénalités par rééchantillonnage
- Nombreuses pénalités faisant intervenir une ou plusieurs constante dont la valeur optimale est inconnue (complexités de Rademacher, détection de ruptures, etc.)

Motivation (2) : « L-curve » et heuristique de coude



Motivation (2) : « L-curve » et heuristique de coude



Motivation (3) : niveau minimal de pénalisation

Régression, moindres carrés :

- Pénalité optimale (C_p , Mallows, 1973) :

$$\frac{2\sigma^2 D_m}{n}$$

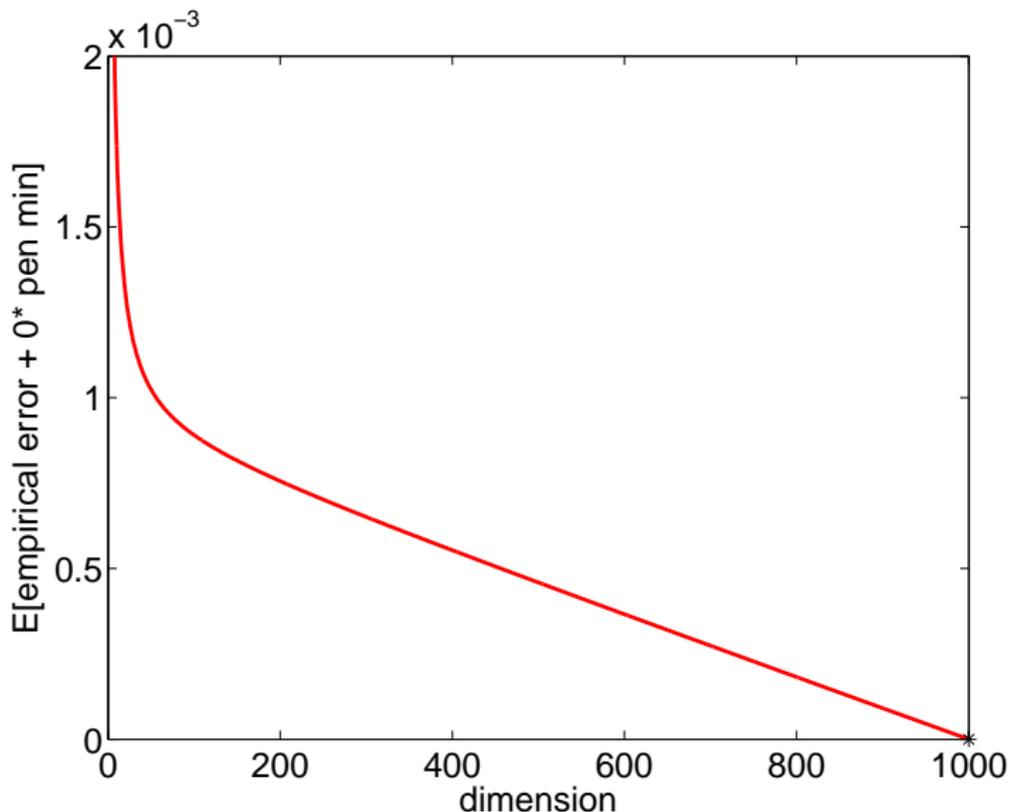
- Inégalité oracle démontrée pour

$$\frac{K\sigma^2 D_m}{n}$$

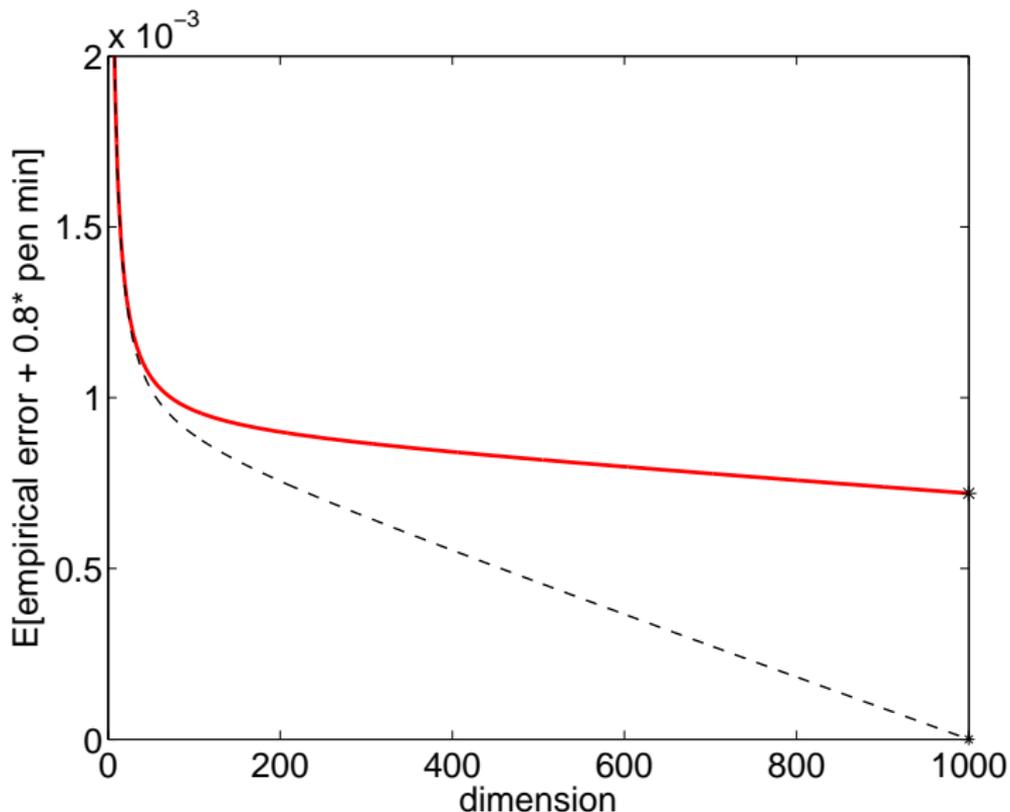
pour tout $K > 1$ (Birgé & Massart, 2001)

⇒ Que se passe-t-il pour $K \leq 1$?

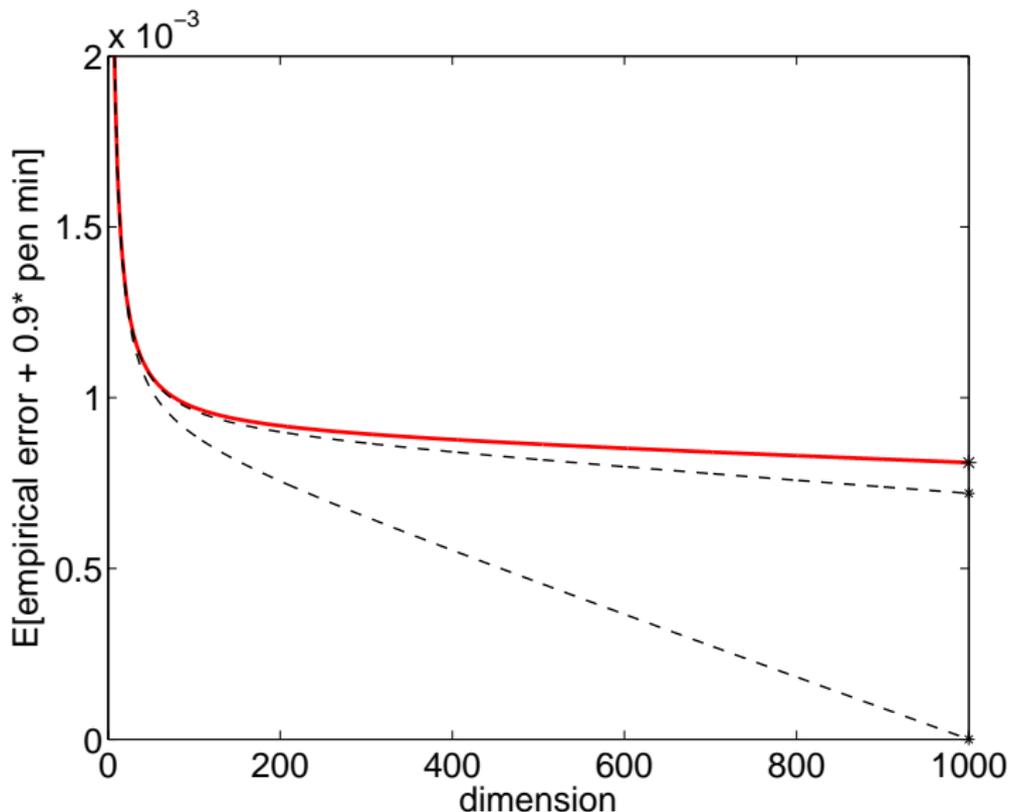
$\mathbb{E}[\text{Risque empirique}] + 0 \times \sigma^2 D_m n^{-1}$ (moindres carrés)



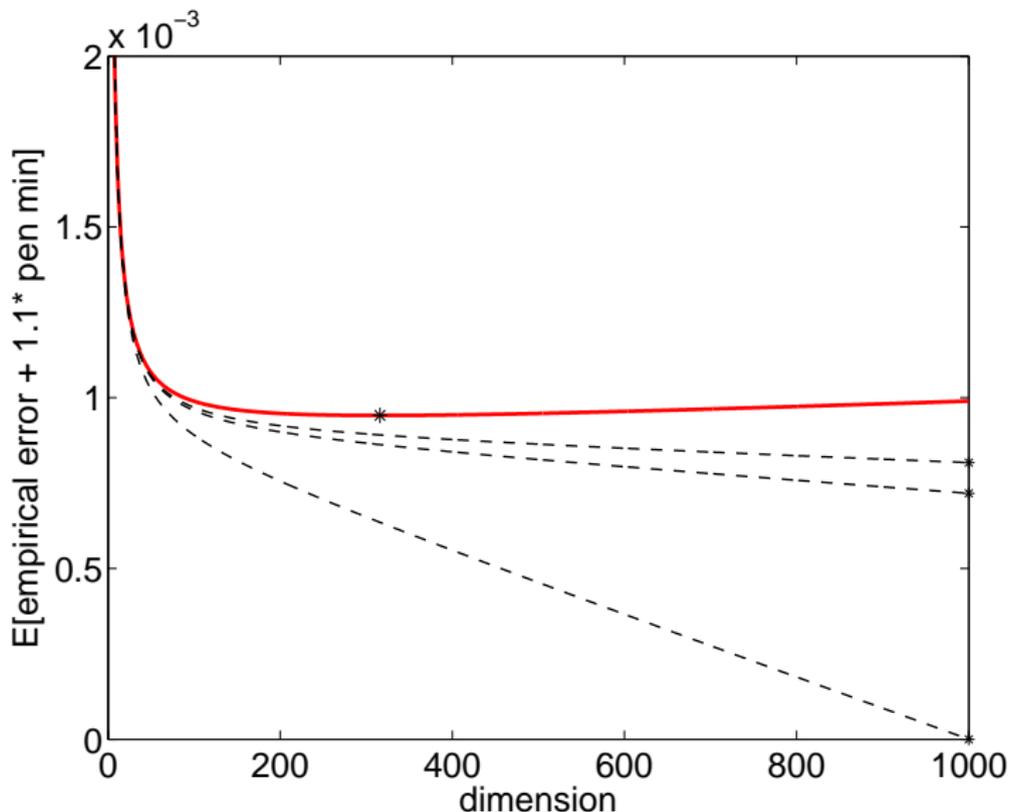
$$\mathbb{E}[\text{Risque empirique}] + 0.8 \times \sigma^2 D_m n^{-1} \text{ (moindres carrés)}$$

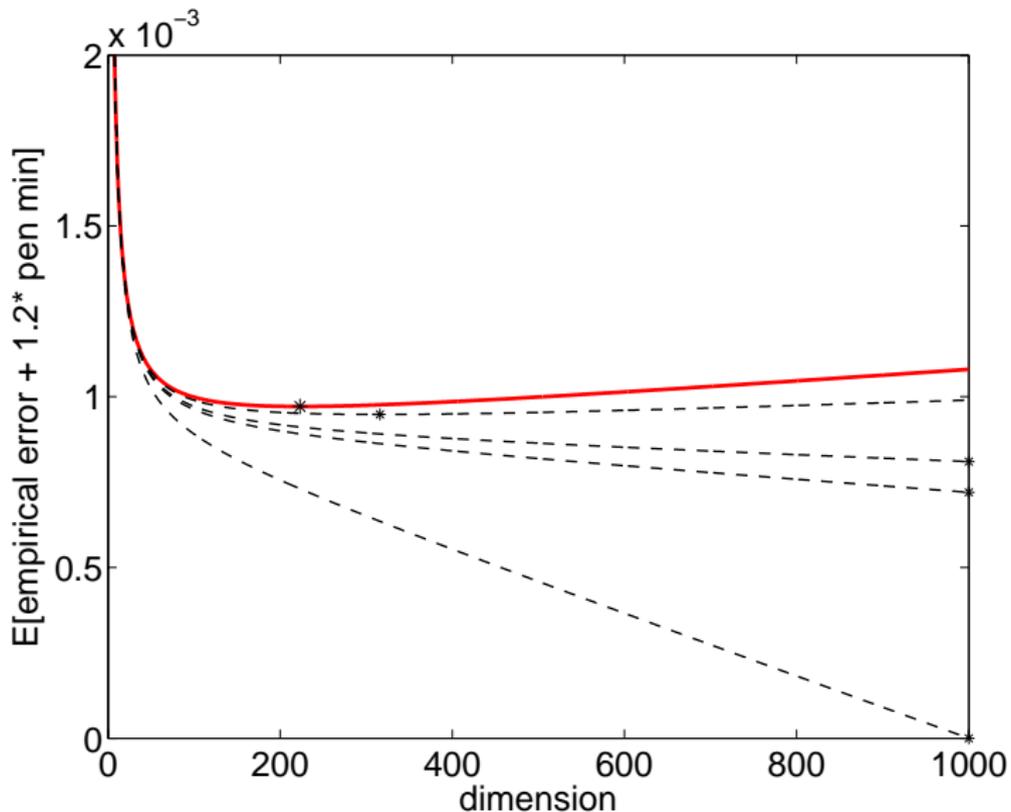


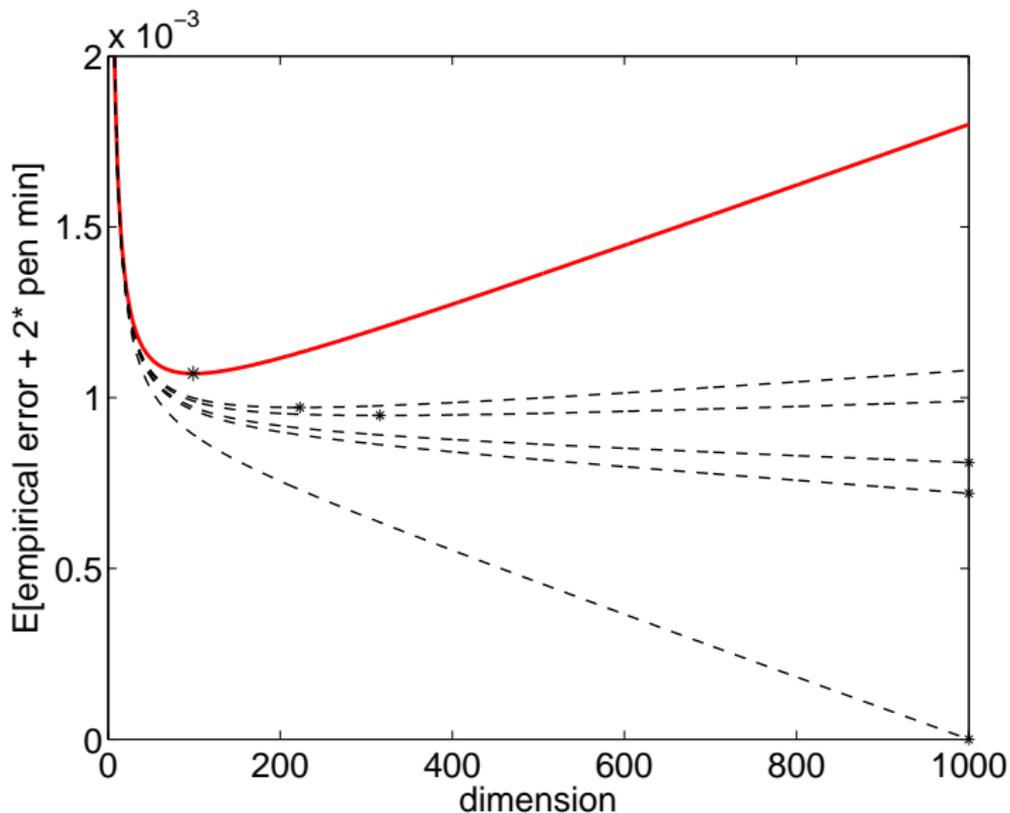
$$\mathbb{E}[\text{Risque empirique}] + 0.9 \times \sigma^2 D_m n^{-1} \text{ (moindres carrés)}$$



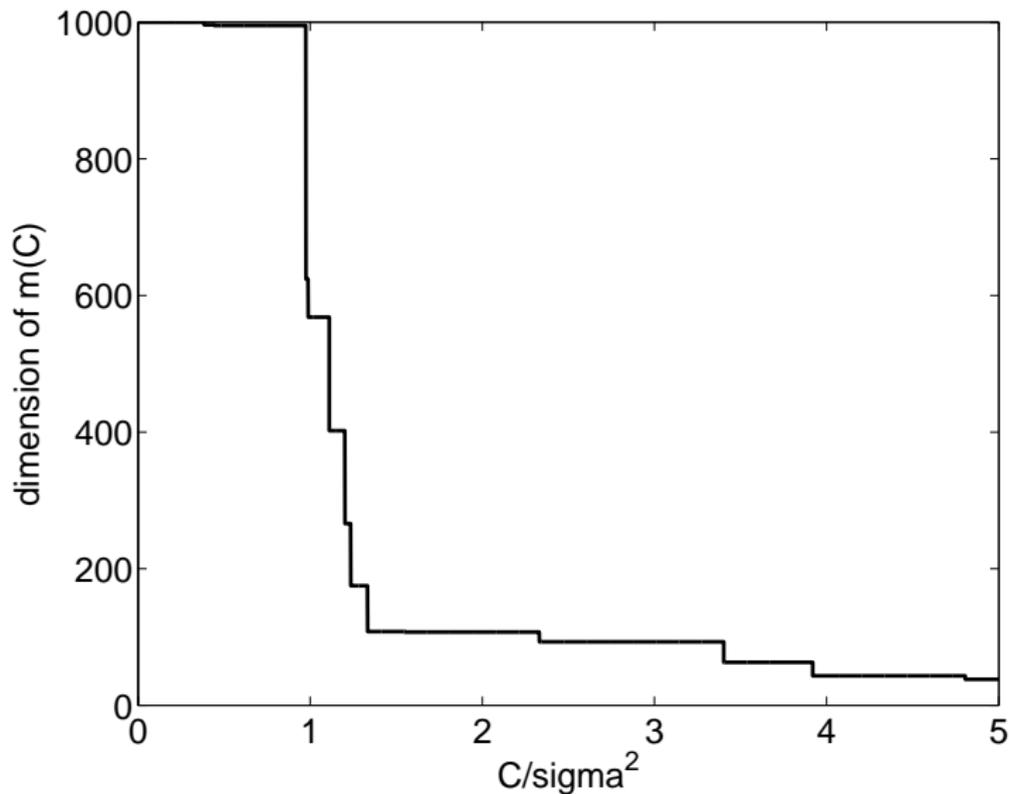
$\mathbb{E}[\text{Risque empirique}] + 1.1 \times \sigma^2 D_m n^{-1}$ (moindres carrés)



$\mathbb{E}[\text{Risque empirique}] + 1.2 \times \sigma^2 D_m n^{-1}$ (moindres carrés)

$$\mathbb{E}[\text{Risque empirique}] + 2 \times \sigma^2 D_m n^{-1} \text{ (moindres carrés)}$$


Moindres carrés : Saut de dimension



Moindres carrés : algorithme (Birgé & Massart 2007)

- ① pour tout $C > 0$, calculer

$$\hat{m}(C) \in \operatorname{argmin}_{m \in \mathcal{M}_n} \left\{ P_n \gamma(\hat{s}_m) + C \frac{D_m}{n} \right\}$$

- ② trouver \hat{C}_{\min} tel que $D_{\hat{m}(C)}$ est « très grande » lorsque $C < \hat{C}_{\min}$ et « raisonnablement petite » lorsque $C > \hat{C}_{\min}$

- ③ choisir $\hat{m} = \hat{m}(2\hat{C}_{\min})$

Mise en œuvre pratique : package CAPUSHE (Baudry, Maugis & Michel, 2011)

<http://www.math.univ-toulouse.fr/~maugis/CAPUSHE.html>

Théorème (1) : Saut de dimension / Pénalité minimale

Théorème (Birgé & Massart 2007, reformulé)

Hypothèses :

- *plan d'expérience déterministe, bruit Gaussien, variance σ^2 constante, $\exists m_0 \in \mathcal{M}$, $S_{m_0} = \mathbb{S}$,*
- $\inf_{m \in \mathcal{M}} \{ \mathbb{E}[\ell(s^*, \hat{s}_m)] \} \leq \sigma^2 \delta_n$, $\delta_n \leq 1/20$.

Alors, $\forall \gamma > 0$, $n \geq n_0(\gamma)$, avec prob. au moins $1 - 4 \text{Card}(\mathcal{M})n^{-\gamma}$,

$$\forall C < (1 - \eta_n^-) \sigma^2, \quad D_{\hat{m}(C)} \geq \frac{9n}{10}$$

$$\forall C > (1 + \eta_n^+) \sigma^2, \quad D_{\hat{m}(C)} \leq \frac{n}{10}$$

$$\text{avec } \eta_n^- = 81 \sqrt{\frac{\gamma \log(n)}{n}}, \quad \eta_n^+ = \eta_n^- + 20\delta_n.$$

Théorème (2) : Inégalité-oracle

Théorème (Birgé & Massart 2007, reformulé)

Sous les hypothèses précédentes, avec probabilité au moins $1 - 4 \text{Card}(\mathcal{M})n^{-\gamma}$, pour tout

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ P_n \gamma(\hat{s}_m) + 2 \hat{C}_{\min} \frac{D_m}{n} \right\},$$

si $n \geq n_0(\gamma)$, pour tout $\eta \geq 2 \max\{\eta_n^-, \eta_n^+\}$,

$$\ell(s^*, \hat{s}_{\hat{m}}) \leq (1 + 3\eta) \inf_{m \in \mathcal{M}} \{\ell(s^*, \hat{s}_m)\} + \frac{880\sigma^2\gamma \log(n)}{\eta n}$$

A. Minimal penalties and the slope heuristics : a survey, 2014.

Estimateurs linéaires

Définition :

$$(\widehat{s}_m(X_i))_{1 \leq i \leq n} = A_m(Y_i)_{1 \leq i \leq n}$$

Estimateurs linéaires

Définition :

$$(\widehat{s}_m(X_i))_{1 \leq i \leq n} = A_m(Y_i)_{1 \leq i \leq n}$$

Exemples :

- moindres carrés

Estimateurs linéaires

Définition :

$$(\widehat{s}_m(X_i))_{1 \leq i \leq n} = A_m(Y_i)_{1 \leq i \leq n}$$

Exemples :

- moindres carrés
- k plus proches voisins
- Nadaraya-Watson
- splines (Wahba, 1990), régression ridge à noyaux :

$$\widehat{s}_\lambda \in \operatorname{argmin}_{t \in \mathcal{H}} \{ P_n \gamma(t) + \lambda \|t\|_{\mathcal{H}}^2 \} .$$

Généralisation aux estimateurs linéaires ?

Moindres carrés

$$\text{pen}_{\text{Cp}}(m) = \frac{2\sigma^2 D_m}{n}$$

$$\operatorname{argmin}_{m \in \mathcal{M}} \left\{ P_n \gamma(\hat{s}_m) + c \frac{D_m}{n} \right\}$$

$\Rightarrow D_{\hat{m}(c)} \ll \text{saute} \gg$ en $\hat{C}_{\min} \approx \sigma^2$

\Rightarrow choix optimal avec $\hat{m}(2\hat{C}_{\min})$

Généralisation aux estimateurs linéaires ?

Moindres carrés

$$\text{pen}_{\text{Cp}}(m) = \frac{2\sigma^2 D_m}{n}$$

$$\operatorname{argmin}_{m \in \mathcal{M}} \left\{ P_n \gamma(\hat{s}_m) + c \frac{D_m}{n} \right\}$$

$\Rightarrow D_{\hat{m}(c)} \ll \text{saute} \gg$ en $\hat{C}_{\min} \approx \sigma^2$

\Rightarrow choix optimal avec $\hat{m}(2\hat{C}_{\min})$

Estimateurs linéaires

$$\text{pen}_{\text{CL}}(m) = \frac{2\sigma^2 \operatorname{tr}(A_m)}{n}$$

Généralisation aux estimateurs linéaires ?

Moindres carrés

$$\text{pen}_{\text{Cp}}(m) = \frac{2\sigma^2 D_m}{n}$$

$$\operatorname{argmin}_{m \in \mathcal{M}} \left\{ P_n \gamma(\hat{s}_m) + c \frac{D_m}{n} \right\}$$

$\Rightarrow D_{\hat{m}(c)} \ll \text{saute} \gg$ en $\hat{C}_{\min} \approx \sigma^2$

\Rightarrow choix optimal avec $\hat{m}(2\hat{C}_{\min})$

Estimateurs linéaires

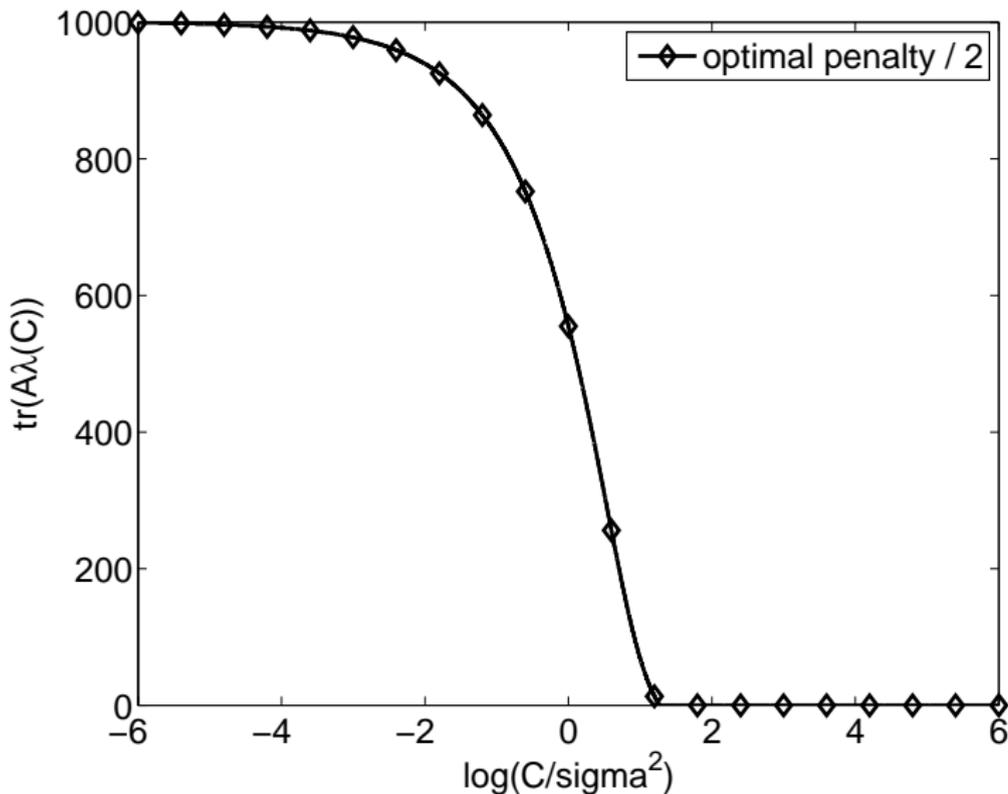
$$\text{pen}_{\text{CL}}(m) = \frac{2\sigma^2 \operatorname{tr}(A_m)}{n}$$

$$\operatorname{argmin}_{m \in \mathcal{M}} \left\{ P_n \gamma(\hat{s}_m) + c \frac{\operatorname{tr}(A_m)}{n} \right\}$$

$\operatorname{tr}(A_{\hat{m}(c)})$ saute-t-il en $\hat{C}_{\min} \approx \sigma^2$?

choix optimal avec $\hat{m}(2\hat{C}_{\min})$?

Pas de saut de dimension avec une pénalité $\propto \text{tr}(A_m)$



Pénalité minimale pour les estimateurs linéaires

$$\mathbb{E}[\ell(s^*, \hat{s}_m)] = \text{Err. approx.} + \underbrace{\frac{\text{tr}(A_m^\top A_m) \sigma^2}{n}}_{\text{Err. estim.}}$$

Pénalité minimale pour les estimateurs linéaires

$$\mathbb{E}[\ell(\mathbf{s}^*, \hat{\mathbf{s}}_m)] = \text{Err. approx.} + \underbrace{\frac{\text{tr}(A_m^\top A_m) \sigma^2}{n}}_{\text{Err. estim.}}$$

$$\mathbb{E}[P_n \gamma(\hat{\mathbf{s}}_m)] = \sigma^2 + \text{Err. approx.} - \frac{(2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)) \sigma^2}{n}$$

Pénalité minimale pour les estimateurs linéaires

$$\mathbb{E}[\ell(\mathbf{s}^*, \hat{\mathbf{s}}_m)] = \text{Err. approx.} + \underbrace{\frac{\text{tr}(A_m^\top A_m) \sigma^2}{n}}_{\text{Err. estim.}}$$

$$\mathbb{E}[P_n \gamma(\hat{\mathbf{s}}_m)] = \sigma^2 + \text{Err. approx.} - \frac{(2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)) \sigma^2}{n}$$

$$\Rightarrow \text{pénalité optimale} \quad \frac{(2 \text{tr}(A_m)) \sigma^2}{n}$$

Pénalité minimale pour les estimateurs linéaires

$$\mathbb{E}[\ell(\mathbf{s}^*, \hat{\mathbf{s}}_m)] = \text{Err. approx.} + \underbrace{\frac{\text{tr}(A_m^\top A_m) \sigma^2}{n}}_{\text{Err. estim.}}$$

$$\mathbb{E}[P_n \gamma(\hat{\mathbf{s}}_m)] = \sigma^2 + \text{Err. approx.} - \underbrace{\frac{(2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)) \sigma^2}{n}}_{\text{pénalité minimale}}$$

$$\Rightarrow \text{pénalité optimale} \quad \frac{(2 \text{tr}(A_m)) \sigma^2}{n}$$

Pénalité minimale pour les estimateurs linéaires

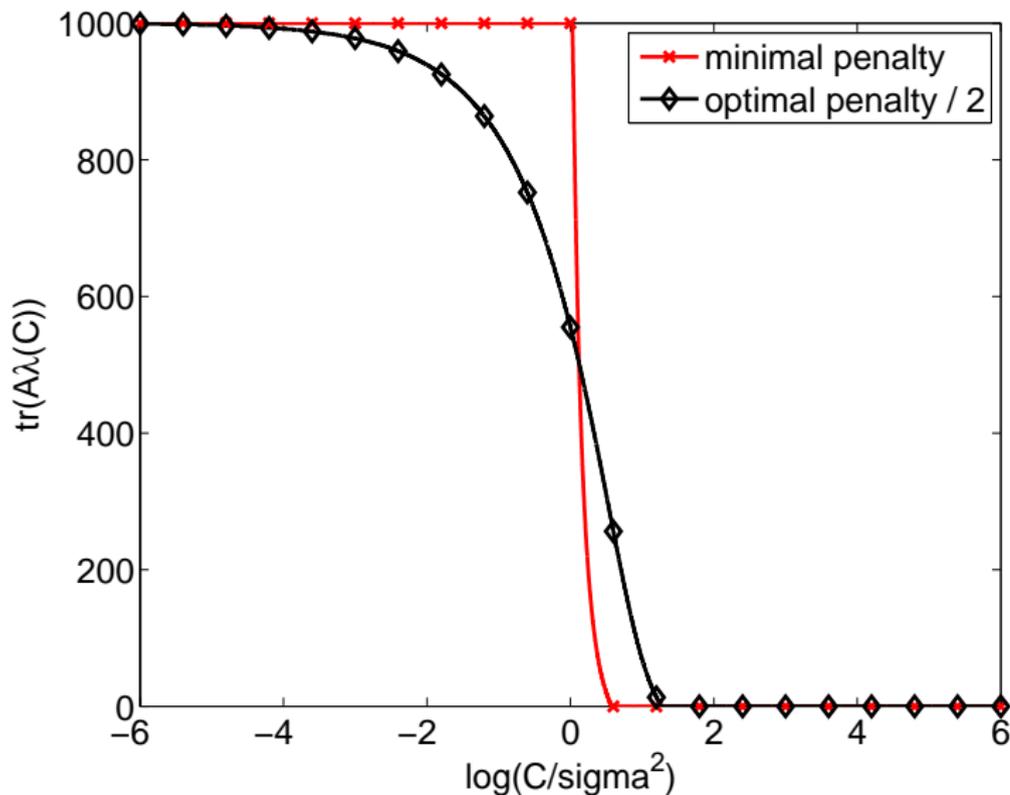
$$\mathbb{E}[\ell(\mathbf{s}^*, \hat{\mathbf{s}}_m)] = \text{Err. approx.} + \underbrace{\frac{\text{tr}(A_m^\top A_m) \sigma^2}{n}}_{\text{Err. estim.}}$$

$$\mathbb{E}[P_n \gamma(\hat{\mathbf{s}}_m)] = \sigma^2 + \text{Err. approx.} - \underbrace{\frac{(2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)) \sigma^2}{n}}_{\text{pénalité minimale}}$$

$$\Rightarrow \text{pénalité optimale} \quad \frac{(2 \text{tr}(A_m)) \sigma^2}{n}$$

$$\hat{m}(C) \in \underset{m \in \mathcal{M}}{\text{argmin}} \left\{ P_n \gamma(\hat{\mathbf{s}}_m) + C \times \frac{2 \text{tr}(A_m) - \text{tr}(A_m^\top A_m)}{n} \right\}$$

Saut de « dimension » (régression ridge)



Algorithme de calibration de pénalités

- ① pour tout $C > 0$, calculer

$$\hat{m}_{\min}(C) \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ P_n \gamma(\hat{s}_m) + \frac{C (2 \operatorname{tr}(A_m) - \operatorname{tr}(A_m^\top A_m))}{n} \right\}$$

- ② trouver \hat{C}_{\min} tel que $\operatorname{tr}(A_{\hat{m}_{\min}(C)})$ est « très grande » lorsque $C < \hat{C}_{\min}$ et « raisonnablement petite » lorsque $C > \hat{C}_{\min}$,

- ③ choisir

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ P_n \gamma(\hat{s}_m) + \frac{2\hat{C}_{\min} \operatorname{tr}(A_m)}{n} \right\}$$

⇒ Inégalité-oracle optimale (au premier ordre) :

A. & Bach. Data-driven calibration of linear estimators with minimal penalties. In *Advances in Neural Information Processing Systems 22*, pages 46–54, 2009.

Pénalités minimales : bilan

- Moindres carrés et k -plus proches voisins :

$$\text{pen}_{\text{opt}} = 2 \text{pen}_{\text{min}}$$

mais pas pour les autres estimateurs linéaires (kernel ridge, Nadaraya-Watson, etc.) !

Pénalités minimales : bilan

- Moindres carrés et k -plus proches voisins :

$$\text{pen}_{\text{opt}} = 2 \text{pen}_{\text{min}}$$

mais pas pour les autres estimateurs linéaires (kernel ridge, Nadaraya-Watson, etc.) !

- Application à l'apprentissage multi-tâches, *via* l'estimation d'une matrice de covariance en régression multivariée

Solnon, A. & Bach. Multi-task regression using minimal penalties. *Journal of Machine Learning Research*, 13:2773–2812, 2012.

Pénalités minimales : bilan

- Moindres carrés et k -plus proches voisins :

$$\text{pen}_{\text{opt}} = 2 \text{pen}_{\text{min}}$$

mais pas pour les autres estimateurs linéaires (kernel ridge, Nadaraya-Watson, etc.) !

- Application à l'apprentissage multi-tâches, *via* l'estimation d'une matrice de covariance en régression multivariée

Solnon, A. & Bach. Multi-task regression using minimal penalties. *Journal of Machine Learning Research*, 13:2773–2812, 2012.

- Régression hétéroscédastique, régressogrammes :

$$\text{pen}_{\text{opt}} \approx 2 \text{pen}_{\text{min}}$$

A. & Massart. Data-driven calibration of penalties for least-squares regression. *Journal of Machine Learning Research*, 10:245–279, 2009.

Pénalités minimales : bilan

- Moindres carrés et k -plus proches voisins :

$$\text{pen}_{\text{opt}} = 2 \text{pen}_{\text{min}}$$

mais pas pour les autres estimateurs linéaires (kernel ridge, Nadaraya-Watson, etc.) !

- Application à l'apprentissage multi-tâches, *via* l'estimation d'une matrice de covariance en régression multivariée

Solnon, A. & Bach. Multi-task regression using minimal penalties. *Journal of Machine Learning Research*, 13:2773–2812, 2012.

- Régression hétéroscédastique, régressogrammes :

$$\text{pen}_{\text{opt}} \approx 2 \text{pen}_{\text{min}}$$

A. & Massart. Data-driven calibration of penalties for least-squares regression. *Journal of Machine Learning Research*, 10:245–279, 2009.

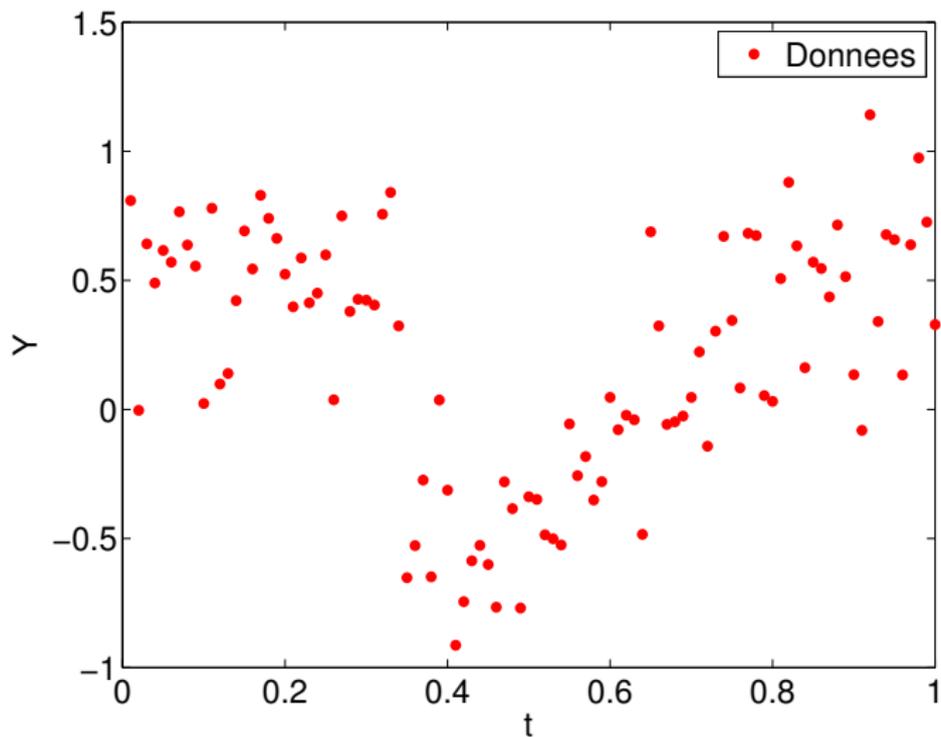
- Article de survol :

A. Minimal penalties and the slope heuristics : a survey, 2014. 51/68

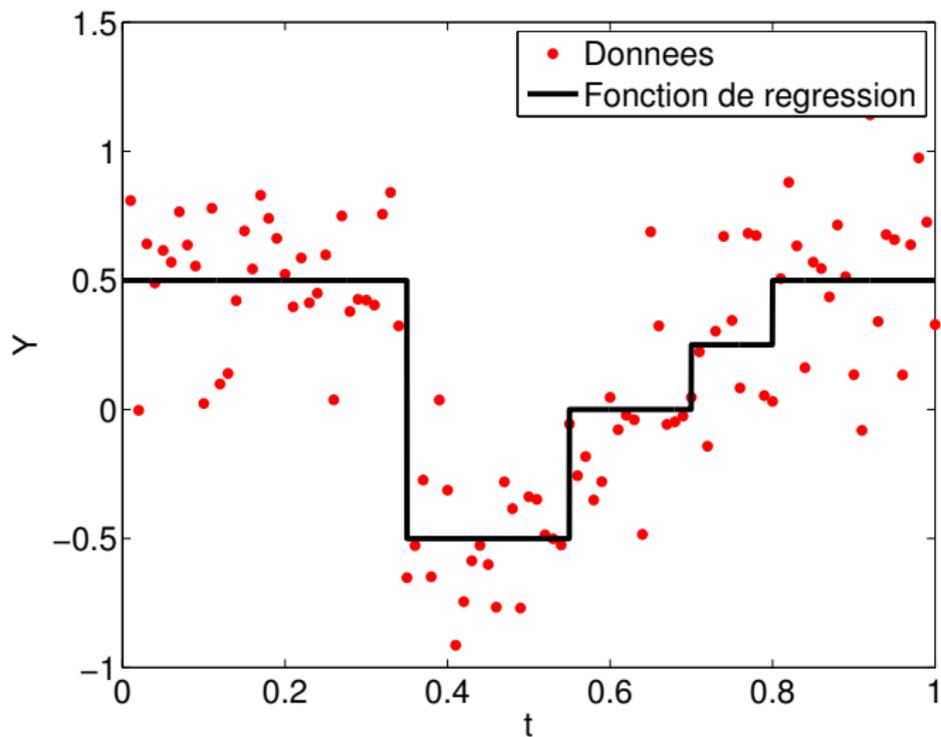
Plan

- 1 Sélection d'estimateurs
- 2 Validation croisée
- 3 Pénalités minimales
- 4 Détection de ruptures**
- 5 Conclusion

Détection de ruptures : données



Détection de ruptures : objectif



Détection de ruptures et sélection de modèles

$$Y_i = s^*(t_i) + \sigma(t_i) \varepsilon_i \quad \text{avec} \quad \mathbb{E}[\varepsilon_i] = 0 \quad \mathbb{E}[\varepsilon_i^2] = 1$$

- But : détecter les ruptures dans la moyenne s^* du signal Y

⇒ Sélection de modèles, collection des régressogrammes avec $\mathcal{M}_n = \mathfrak{P}_{\text{interv}}(\{t_1, \dots, t_n\})$ (ensemble des partitions en intervalles)

Approche classique (Lebarbier, 2005)

- Pénalité « Birgé-Massart » (suppose $\sigma(t_i) \equiv \sigma$) :

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}_n} \left\{ P_n \gamma(\hat{s}_m) + \frac{C\sigma^2 D_m}{n} \left(5 + 2 \log \left(\frac{n}{D_m} \right) \right) \right\}$$

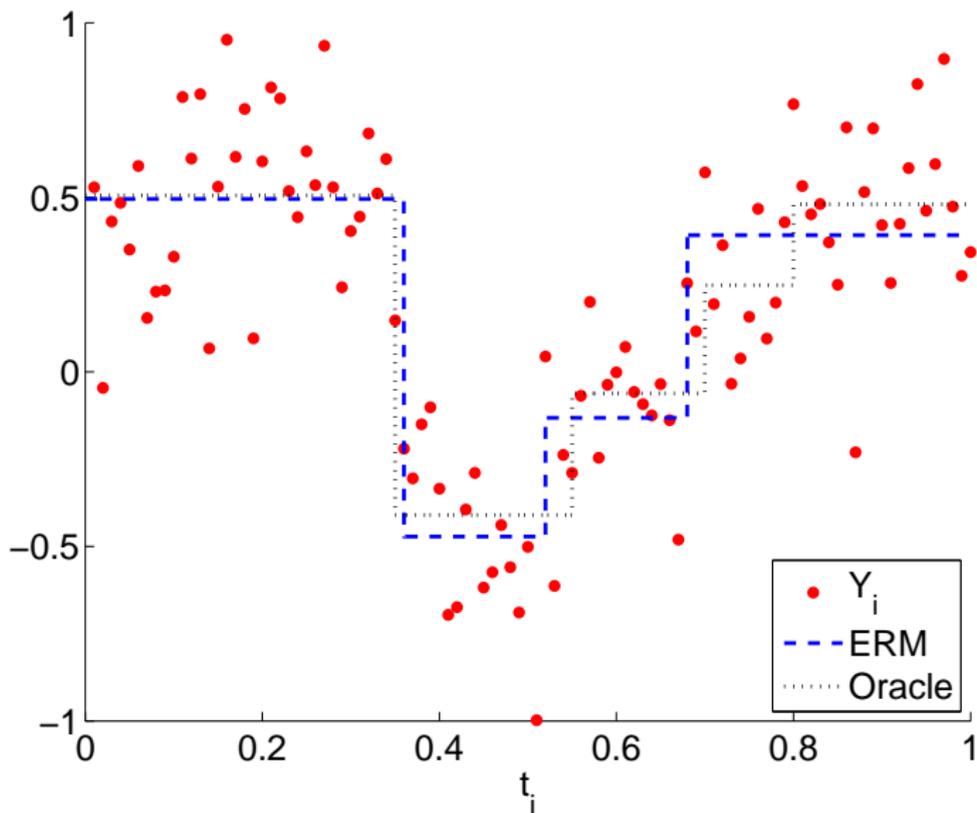
- Revient à agréger les modèles de même dimension :

$$\tilde{S}_D := \bigcup_{m \in \mathcal{M}_n, D_m = D} S_m$$

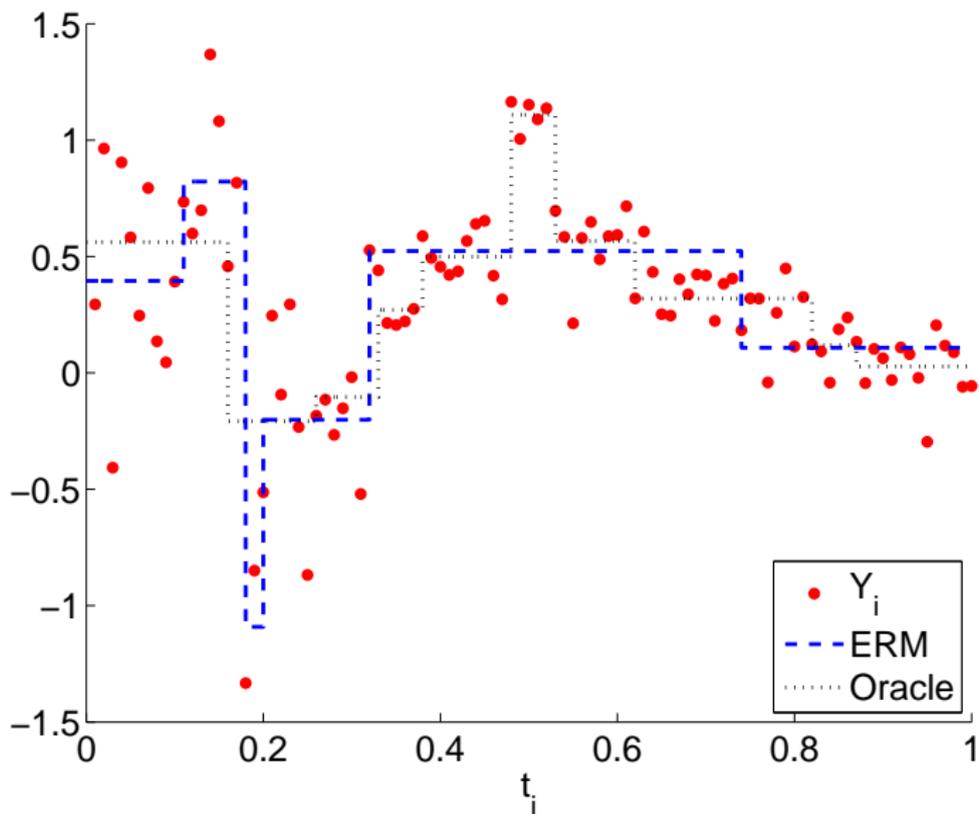
$$\hat{s}_D \in \operatorname{argmin}_{t \in \tilde{S}_D} \{ P_n \gamma(t) \} \quad \text{programmation dynamique}$$

$$\hat{D} \in \operatorname{argmin}_{1 \leq D \leq n} \left\{ P_n \gamma(\hat{s}_D) + \frac{C\sigma^2 D}{n} \left(5 + 2 \log \left(\frac{n}{D} \right) \right) \right\}$$

$D = 4$, homoscédastique ; $n = 100$, $\sigma = 0,25$



$D = 6$, hétéroscédastique ; $n = 100$, $\|\sigma\| = 0,30$



Détection de ruptures par validation croisée

- Échec de la minimisation du risque empirique dans le cas hétéroscédastique suggéré par

A. Choosing a penalty for model selection in heteroscedastic regression, 2010. arXiv:0812.3141v2.

Détection de ruptures par validation croisée

- Échec de la minimisation du risque empirique dans le cas hétéroscédastique suggéré par

A. Choosing a penalty for model selection in heteroscedastic regression, 2010. arXiv:0812.3141v2.

- La validation croisée s'adapte à l'hétéroscédasticité :

A. V-fold cross-validation improved : V-fold penalization, 2008. arXiv:0802.0566v2.

Détection de ruptures par validation croisée

- Échec de la minimisation du risque empirique dans le cas hétéroscédastique suggéré par

A. Choosing a penalty for model selection in heteroscedastic regression, 2010. arXiv:0812.3141v2.

- La validation croisée s'adapte à l'hétéroscédasticité :

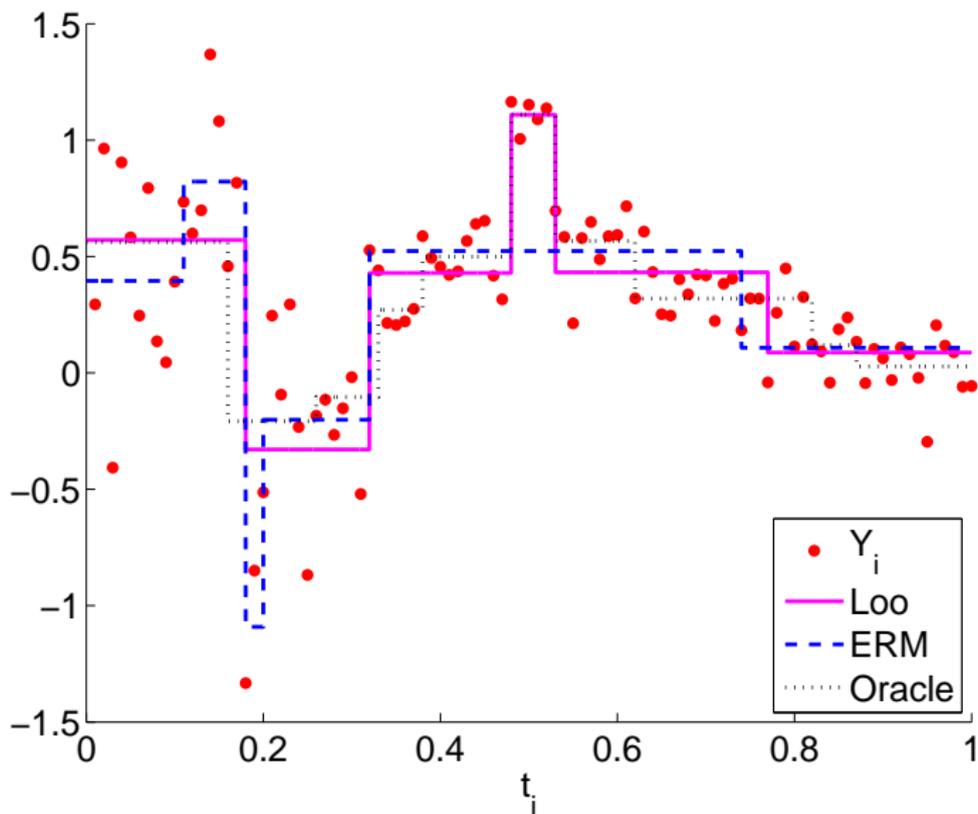
A. V-fold cross-validation improved : V-fold penalization, 2008. arXiv:0802.0566v2.

⇒ Algorithme :

$$\left(\widehat{S}_m\right)_{D_m=D} \rightsquigarrow \widehat{S}_D^{CV} = \widehat{S}_{\widehat{m}^{CV}}(D) \rightsquigarrow \widetilde{S} = \widehat{S}_{\widehat{D}^{CV}}^{CV}$$

A. & Celisse. Segmentation of the mean of heteroscedastic data via cross-validation. *Statistics and Computing*, pages 1–20, 2010.

$D = 6$, hétéroscédastique ; $n = 100$, $\|\sigma\| = 0,30$



Données de grande dimension ou complexes

$Y_1, \dots, Y_n \in \mathcal{Y}$ de loi constante par morceaux

- \mathcal{Y} de grande dimension (structure euclidienne inadaptée) ou qui n'est pas un espace vectoriel

Données de grande dimension ou complexes

$Y_1, \dots, Y_n \in \mathcal{Y}$ de loi constante par morceaux

- \mathcal{Y} de grande dimension (structure euclidienne inadaptée) ou qui n'est pas un espace vectoriel

⇒ Idée : considérer un noyau défini positif $k : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ et chercher les ruptures dans la série

$$\Phi(Y_1), \dots, \Phi(Y_n) \in \mathcal{H} \quad \text{où} \quad \Phi(y) = k(y, \cdot)$$

et \mathcal{H} est le RKHS associé à k .

Données de grande dimension ou complexes

$Y_1, \dots, Y_n \in \mathcal{Y}$ de loi constante par morceaux

- \mathcal{Y} de grande dimension (structure euclidienne inadaptée) ou qui n'est pas un espace vectoriel

⇒ Idée : considérer un noyau défini positif $k : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ et chercher les ruptures dans la série

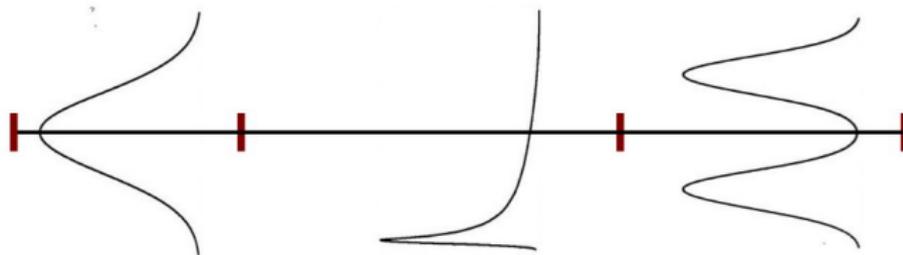
$$\Phi(Y_1), \dots, \Phi(Y_n) \in \mathcal{H} \quad \text{où} \quad \Phi(y) = k(y, \cdot)$$

et \mathcal{H} est le RKHS associé à k .

- Pénalité similaire à celle de Lebarbier (2005)
⇒ inégalité-oracle.

A., Celisse & Harchaoui. Kernel change-point detection, 2012. arXiv:1202.3878v1.

Détection de ruptures à noyaux : applications

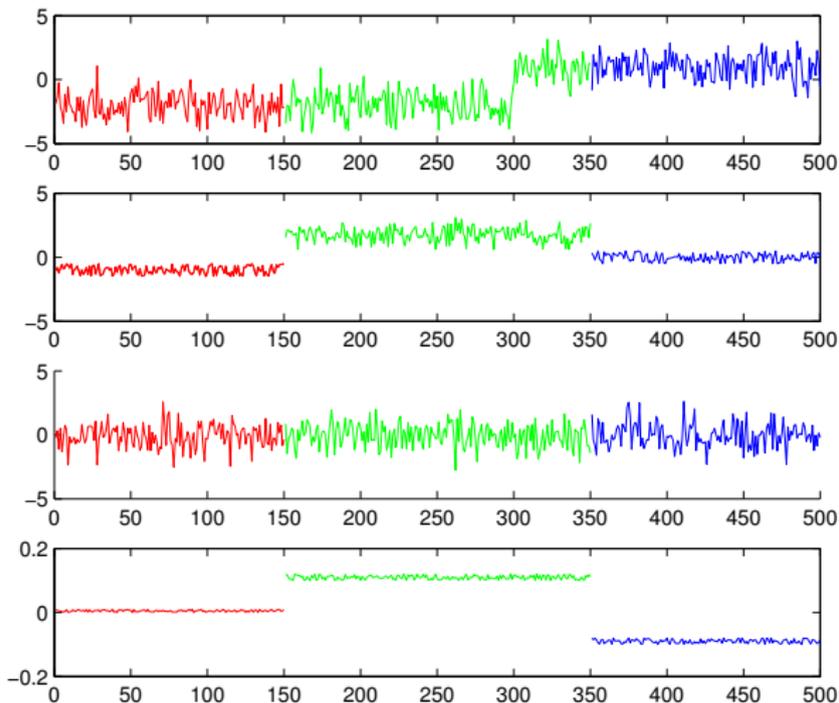


music

applause

speech

Apprentissage de métrique : objectif



Apprentissage de métrique

$$Y = (Y_1, \dots, Y_n) \in (\mathbb{R}^d)^n$$

$$\forall B \succeq 0, \quad \hat{m}_B(Y) \in \operatorname{argmin}_{m \in \mathcal{M}} \{P_n \gamma_B(\hat{s}_m) + \lambda D_m\} \Rightarrow \hat{B} ?$$

Apprentissage de métrique

$$Y = (Y_1, \dots, Y_n) \in (\mathbb{R}^d)^n$$

$$\forall B \succeq 0, \quad \hat{m}_B(Y) \in \underset{m \in \mathcal{M}}{\operatorname{argmin}} \{P_n \gamma_B(\hat{s}_m) + \lambda D_m\} \Rightarrow \hat{B} ?$$

- Hypothèse : on dispose de séries segmentées

$$(Y^1, m^1), \dots, (Y^N, m^N) .$$

Apprentissage de métrique

$$Y = (Y_1, \dots, Y_n) \in (\mathbb{R}^d)^n$$

$$\forall B \succeq 0, \quad \hat{m}_B(Y) \in \operatorname{argmin}_{m \in \mathcal{M}} \{P_n \gamma_B(\hat{s}_m) + \lambda D_m\} \Rightarrow \hat{B} ?$$

- Hypothèse : on dispose de séries segmentées

$$(Y^1, m^1), \dots, (Y^N, m^N) .$$

⇒ prédiction structurée : idéalement, on voudrait

$$\hat{B} \in \operatorname{argmin}_{B \succeq 0} \left\{ \frac{1}{N} \sum_{j=1}^N \gamma(m^j, \hat{m}_B(Y^j)) + \underbrace{\Omega(B)}_{\text{régularisation}} \right\}$$

Apprentissage de métrique

$$Y = (Y_1, \dots, Y_n) \in (\mathbb{R}^d)^n$$

$$\forall B \succeq 0, \quad \hat{m}_B(Y) \in \operatorname{argmin}_{m \in \mathcal{M}} \{P_n \gamma_B(\hat{s}_m) + \lambda D_m\} \Rightarrow \hat{B} ?$$

- Hypothèse : on dispose de séries segmentées

$$(Y^1, m^1), \dots, (Y^N, m^N) .$$

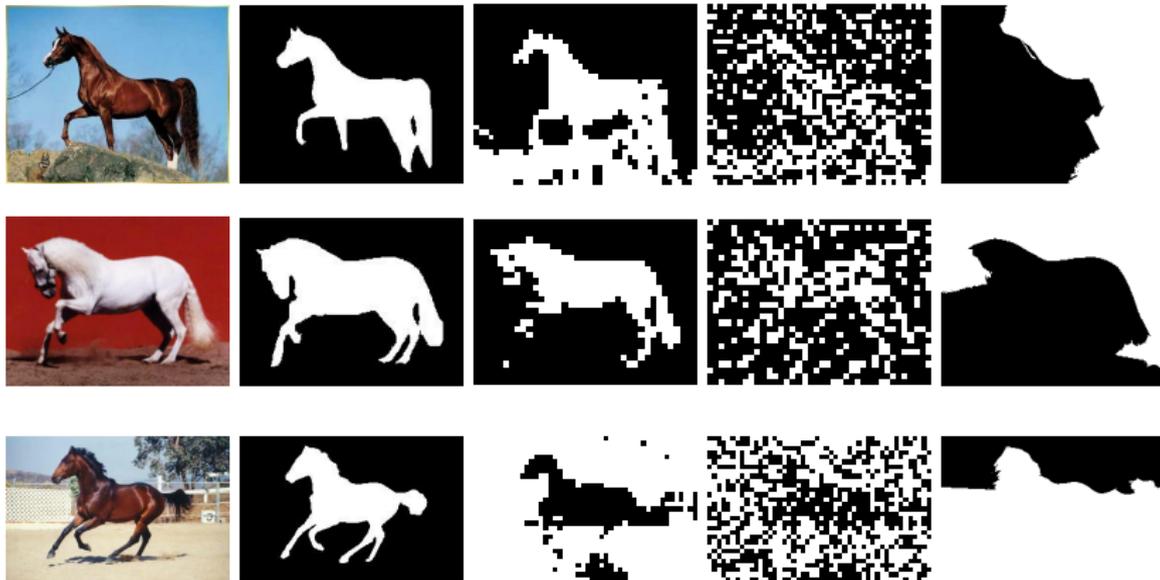
⇒ prédiction structurée : idéalement, on voudrait

$$\hat{B} \in \operatorname{argmin}_{B \succeq 0} \left\{ \frac{1}{N} \sum_{j=1}^N \gamma(m^j, \hat{m}_B(Y^j)) + \underbrace{\Omega(B)}_{\text{régularisation}} \right\}$$

⇒ **relaxation convexe (Tsochantaridis *et al.*, 2005)**

Lajugie, A. & Bach. Large-margin metric learning for constrained partitioning problems. In *International Conference on Machine Learning (ICML)*, volume 32, pages 297–305, 2014.

Apprentissage de métrique : résultats (segmentation 2d)



données

cible

notre méthode

sans prior

normalized cuts

Autres travaux

- Régions de confiance et tests multiples par **rééchantillonnage**

A., Blanchard & Roquain. Some nonasymptotic results on resampling in high dimension, I : Confidence regions. II : Multiple tests. *The Annals of Statistics*, 38(1):51–99, 2010.

- **Rééchantillonnage** et précision d'extrapolation dans une série temporelle, application en astronomie :

Desmars, A., Arlot, Lainey & Vienne. Estimating the accuracy of satellite ephemerides using the bootstrap method. *Astronomy and Astrophysics*, 499:321–330, 2009.

- Adaptation à la **condition de marge** en classification :

A. & Bartlett. Margin adaptive model selection in statistical learning. *Bernoulli*, 17(2):687–713, 2011.

- Erreur d'approximation de **forêts purement aléatoires** :

A. & Genuer. Analysis of purely random forests bias, 2014. arXiv:1407.3939v1.

- **Apprentissage de métrique** pour l'alignement dynamique de séquences :

Garreau, Lajugie, A. & Bach. Metric learning for temporal sequence alignment. In *Advances in Neural Information Processing Systems* 27, 2014.

Bilan

- Comprendre pourquoi certaines procédures sont performantes
⇒ heuristique de coude, forêts (purement) aléatoires
(Breiman, 2001).

Bilan

- Comprendre pourquoi certaines procédures sont performantes
⇒ heuristique de coude, forêts (purement) aléatoires (Breiman, 2001).
- Pourquoi certaines procédures fonctionnent-elles mieux ?
Compromis entre complexité algorithmique et performance statistique ?
⇒ choix de V pour la validation croisée « V -fold ».

Bilan

- Comprendre pourquoi certaines procédures sont performantes
⇒ heuristique de coude, forêts (purement) aléatoires (Breiman, 2001).
- Pourquoi certaines procédures fonctionnent-elles mieux ?
Compromis entre complexité algorithmique et performance statistique ?
⇒ choix de V pour la validation croisée « V -fold ».
- Corriger les défauts de méthodes couramment utilisées
⇒ correction du biais de la validation croisée « V -fold ».

Bilan

- Comprendre pourquoi certaines procédures sont performantes
⇒ **heuristique de coude, forêts (purement) aléatoires (Breiman, 2001).**
- Pourquoi certaines procédures fonctionnent-elles mieux ?
Compromis entre complexité algorithmique et performance statistique ?
⇒ **choix de V pour la validation croisée « V -fold ».**
- Corriger les défauts de méthodes couramment utilisées
⇒ **correction du biais de la validation croisée « V -fold ».**
- Proposer de nouvelles méthodes sur des bases théoriques
⇒ **généralisation de l'heuristique de pente aux estimateurs linéaires, détection de ruptures par validation croisée ou à noyaux, etc.**

Perspectives

- Comparaisons au second ordre
 - Rendre l'heuristique quantitative
 - Application à la surpénalisation ?

Perspectives

- Comparaisons au second ordre
 - Rendre l'heuristique quantitative
 - Application à la surpénalisation ?
- Validation croisée : extension à d'autres cadres
 - Correction du biais
 - Prise en compte de la variance (comportements différents ?)

Perspectives

- Comparaisons au second ordre
 - Rendre l'heuristique quantitative
 - Application à la surpénalisation ?
- Validation croisée : extension à d'autres cadres
 - Correction du biais
 - Prise en compte de la variance (comportements différents ?)
- Pénalités minimales
 - Grandes familles de modèles
 - Surpénalisation automatique ?

Perspectives

- Comparaisons au second ordre
 - Rendre l'heuristique quantitative
 - Application à la surpénalisation ?
- Validation croisée : extension à d'autres cadres
 - Correction du biais
 - Prise en compte de la variance (comportements différents ?)
- Pénalités minimales
 - Grandes familles de modèles
 - Surpénalisation automatique ?
- Détection de ruptures
 - Prise en compte de l'hétéroscédasticité dans le cadre hilbertien
 - Choix du noyau ?