

# Fondamentaux de l'apprentissage statistique

Sylvain Arlot

Laboratoire de Mathématiques d'Orsay, Université Paris-Sud

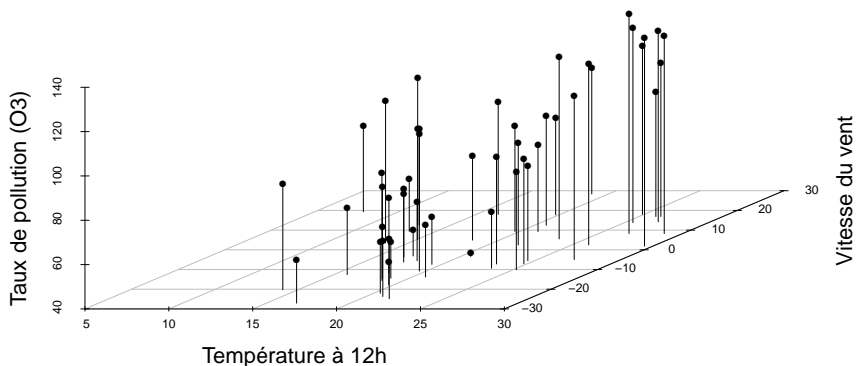
JES 2016, Fréjus

3 Octobre 2016

# Plan

- 1 **Prévision**
- 2 Régression et classification
- 3 Minimisation du risque empirique
- 4 Moyennes locales
- 5 On n'a rien sans rien
- 6 Conclusion : enjeux de l'apprentissage

## Exemple : taux de pollution (régression)



Nouvelle observation (T12, Vx)  $\Rightarrow$  taux de pollution ?

Figure : [Cornillon and Matzner-Løber, 2011]

## Reconnaissance de chiffres manuscrits (MNIST)

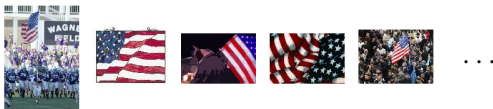


5 ⇒ ?

<http://yann.lecun.com/exdb/mnist/>

# Reconnaissance d'objets

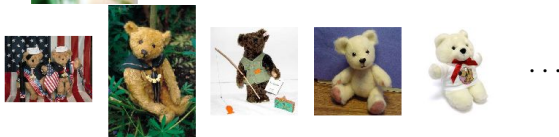
Drapeau américain :



Papillon :



Ours en peluche :



⇒ Drapeau américain ? Papillon ?  
Ours en peluche ? ...

[http://www.vision.caltech.edu/Image\\_Datasets/Caltech256/](http://www.vision.caltech.edu/Image_Datasets/Caltech256/) 4/63

# Nombreux autres exemples

- **Bioinformatique** :
  - données de séquençage  $\Rightarrow$  diagnostic et pronostic (cancer, ...)
  - médecine personnalisée
  - ...
- **Classification de texte** :
  - Détection de spams
  - Publicité en ligne
  - Classification automatique de documents
- Reconnaissance d'**actions humaines** dans des **vidéos**
- Reconnaissance de **parole**
- Évaluation des risques-clients (**prêt bancaire**)
- ...

# Formalisation : problème de prévision

- **Données** :  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$   
 $X_i \in \mathcal{X}$  : variable explicative  
 $Y_i \in \mathcal{Y}$  : variable d'intérêt  
Hypothèse :  $(X, Y), (X_1, Y_1), \dots, (X_n, Y_n), \dots$  i.i.d.  $\sim P$

# Formalisation : problème de prévision

- **Données** :  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$   
 $X_i \in \mathcal{X}$  : variable explicative  
 $Y_i \in \mathcal{Y}$  : variable d'intérêt  
Hypothèse :  $(X, Y), (X_1, Y_1), \dots, (X_n, Y_n), \dots$  i.i.d.  $\sim P$
- **Prédicteur** :  $f : \mathcal{X} \rightarrow \mathcal{Y}$   
( $\mathcal{F}$  : ensemble des prédicteurs)  
Nouvelle observation  $X_{n+1} \Rightarrow f(X_{n+1})$  « prévoit »  $Y_{n+1}$



# Formalisation : problème de prévision

- **Données** :  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$   
 $X_i \in \mathcal{X}$  : variable explicative  
 $Y_i \in \mathcal{Y}$  : variable d'intérêt  
Hypothèse :  $(X, Y), (X_1, Y_1), \dots, (X_n, Y_n), \dots$  i.i.d.  $\sim P$
- **Prédicteur** :  $f : \mathcal{X} \rightarrow \mathcal{Y}$   
( $\mathcal{F}$  : ensemble des prédicteurs)  
Nouvelle observation  $X_{n+1} \Rightarrow f(X_{n+1})$  « prévoit »  $Y_{n+1}$
- **Mesure de qualité** : fonction de coût  $c : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, +\infty[$   
Risque (erreur de prévision) :  $\mathcal{R}_P(f) = \mathbb{E} \left[ c(f(X), Y) \right]$

# Formalisation : problème de prévision

- Données** :  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$   
 $X_i \in \mathcal{X}$  : variable explicative  
 $Y_i \in \mathcal{Y}$  : variable d'intérêt  
 Hypothèse :  $(X, Y), (X_1, Y_1), \dots, (X_n, Y_n), \dots$  i.i.d.  $\sim P$
- Prédicteur** :  $f : \mathcal{X} \rightarrow \mathcal{Y}$   
 ( $\mathcal{F}$  : ensemble des prédicteurs)  
 Nouvelle observation  $X_{n+1} \Rightarrow f(X_{n+1})$  « prévoit »  $Y_{n+1}$
- Mesure de qualité** : fonction de coût  $c : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, +\infty[$   
 Risque (erreur de prévision) :  $\mathcal{R}_P(f) = \mathbb{E} \left[ c(f(X), Y) \right]$
- En résumé : avec  $D_n$  uniquement, on cherche un prédicteur  $f \in \mathcal{F}$  tel que  $\mathcal{R}_P(f)$  est minimal.

# Formalisation : problème de régression

- Données** :  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$   
 $X_i \in \mathcal{X}$  : variable explicative  
 $Y_i \in \mathcal{Y}$  : variable d'intérêt (Régression :  $\mathcal{Y} = \mathbb{R}$ )  
 Hypothèse :  $(X, Y), (X_1, Y_1), \dots, (X_n, Y_n), \dots$  i.i.d.  $\sim P$
- Prédicteur** :  $f : \mathcal{X} \rightarrow \mathcal{Y}$   
 ( $\mathcal{F}$  : ensemble des prédicteurs)  
 Nouvelle observation  $X_{n+1} \Rightarrow f(X_{n+1})$  « prévoit »  $Y_{n+1}$
- Mesure de qualité** : fonction de coût  $c : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, +\infty[$   
 Risque (erreur de prévision) :  $\mathcal{R}_P(f) = \mathbb{E} \left[ c(f(X), Y) \right]$   
 Ex. en régression :  $c(y, y') = (y - y')^2$
- En résumé : avec  $D_n$  uniquement, on cherche un prédicteur  $f \in \mathcal{F}$  tel que  $\mathcal{R}_P(f)$  est minimal.

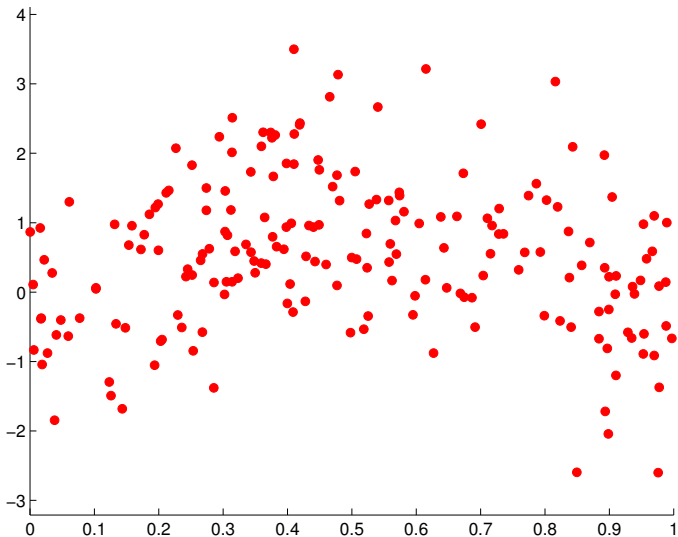
## Formalisation : problème de classification binaire supervisée

- **Données** :  $D_n = (X_i, Y_i)_{1 \leq i \leq n}$   
 $X_i \in \mathcal{X}$  : variable explicative  
 $Y_i \in \mathcal{Y}$  : variable d'intérêt (Classif. binaire :  $\mathcal{Y} = \{0, 1\}$ )  
 Hypothèse :  $(X, Y), (X_1, Y_1), \dots, (X_n, Y_n), \dots$  i.i.d.  $\sim P$
- **Prédicteur** :  $f : \mathcal{X} \rightarrow \mathcal{Y}$   
 ( $\mathcal{F}$  : ensemble des prédicteurs)  
 Nouvelle observation  $X_{n+1} \Rightarrow f(X_{n+1})$  « prévoit »  $Y_{n+1}$
- **Mesure de qualité** : fonction de coût  $c : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, +\infty[$   
 Risque (erreur de prévision) :  $\mathcal{R}_P(f) = \mathbb{E} \left[ c(f(X), Y) \right]$   
 Ex. en classification :  $c(y, y') = \mathbb{1}_{y \neq y'}$
- En résumé : avec  $D_n$  uniquement, on cherche un prédicteur  $f \in \mathcal{F}$  tel que  $\mathcal{R}_P(f)$  est minimal.

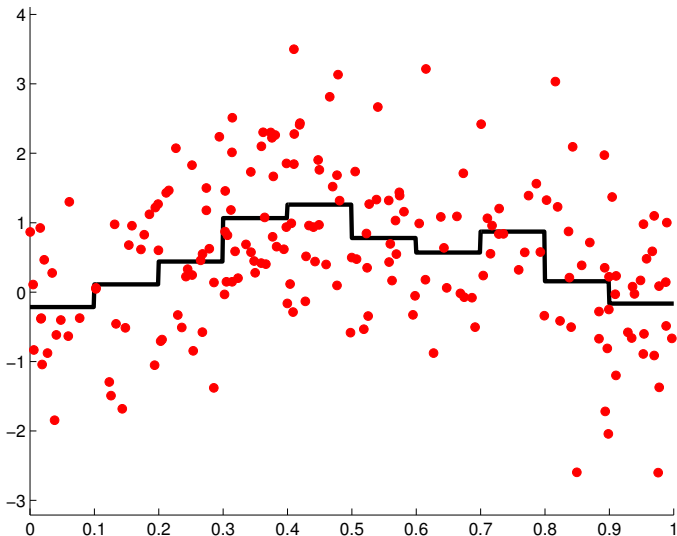
# Prévision idéale

- **Risque de Bayes** :  $\mathcal{R}_P^* := \inf_{f \in \mathcal{F}} \mathcal{R}_P(f)$
- **Prédicteur de Bayes** :  $f^* \in \operatorname{argmin}_{f \in \mathcal{F}} \{\mathcal{R}_P(f)\}$
- **Excès de risque** :  $\ell(f^*, f) = \mathcal{R}_P(f) - \mathcal{R}_P^* \geq 0$

# Et avec les données uniquement ?



# Règle d'apprentissage par partition (régression)



# Règles d'apprentissage

- Règle d'apprentissage :  $\hat{f} : \bigcup_{n \geq 1} (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{F}$   
Entrée : un échantillon  $D_n$  (de taille quelconque  $n \geq 1$ )  
Sortie : un prédicteur  $\hat{f}(D_n) : \mathcal{X} \rightarrow \mathcal{Y}$



# Règles d'apprentissage

- Règle d'apprentissage :  $\hat{f} : \bigcup_{n \geq 1} (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{F}$   
Entrée : un échantillon  $D_n$  (de taille quelconque  $n \geq 1$ )  
Sortie : un prédicteur  $\hat{f}(D_n) : \mathcal{X} \rightarrow \mathcal{Y}$

- **Risque conditionnel** :

$$\mathcal{R}_P(\hat{f}(D_n)) = \mathbb{E} \left[ c(\hat{f}(D_n; X), Y) \mid D_n \right]$$

- **Risque moyen** :

$$\mathbb{E} \left[ \mathcal{R}_P(\hat{f}(D_n)) \right] = \mathbb{E} \left[ c(\hat{f}(D_n; X), Y) \right]$$

# Règles d'apprentissage

- Règle d'apprentissage :  $\hat{f} : \bigcup_{n \geq 1} (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{F}$   
 Entrée : un échantillon  $D_n$  (de taille quelconque  $n \geq 1$ )  
 Sortie : un prédicteur  $\hat{f}(D_n) : \mathcal{X} \rightarrow \mathcal{Y}$

- Risque conditionnel :

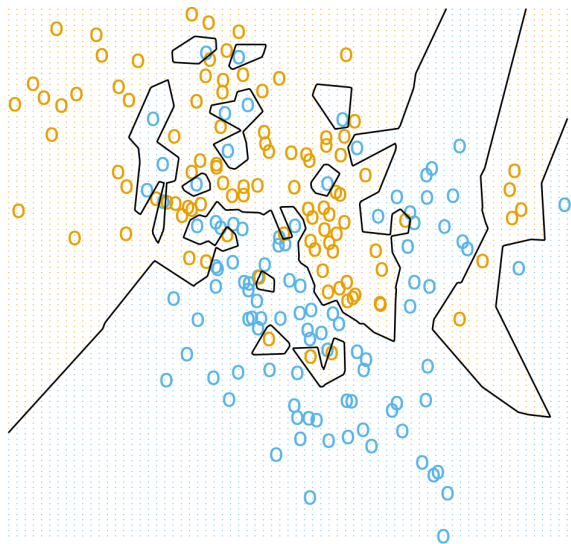
$$\mathcal{R}_P(\hat{f}(D_n)) = \mathbb{E} \left[ c(\hat{f}(D_n; X), Y) \mid D_n \right]$$

- Risque moyen :

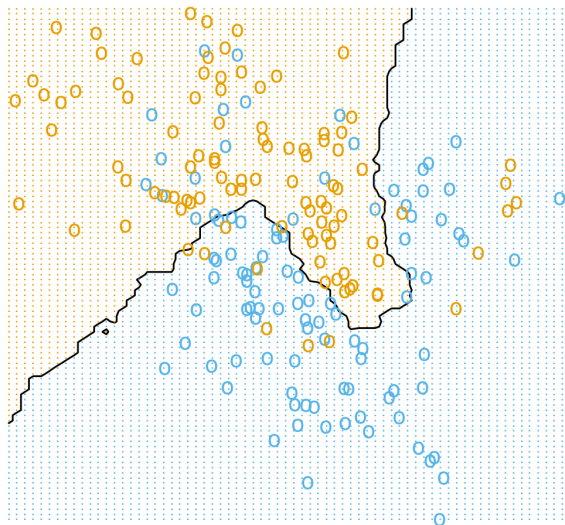
$$\mathbb{E} \left[ \mathcal{R}_P(\hat{f}(D_n)) \right] = \mathbb{E} \left[ c(\hat{f}(D_n; X), Y) \right]$$

- Exemple : règle par partition

# Exemple : plus proche voisin



$x \in \mathcal{X}$   
→ étiquette  $Y_i$   
du **plus proche**  
**voisin**  $X_i$  de  $x$   
parmi  $X_1, \dots, X_n$

Exemple :  $k$  plus proches voisins ( $k = 15$ )

$x \in \mathcal{X}$   
 $\rightarrow$  **vote majoritaire/moyenne**  
 parmi les  
 étiquettes des  
 **$k$  plus proches**  
**voisins** de  $x$  parmi  
 $X_1, \dots, X_n$

Figure : [Hastie et al., 2009]

# Consistance

- faible (pour  $P$ ) :

$$\mathbb{E} \left[ \mathcal{R}_P(\hat{f}(D_n)) \right] \xrightarrow{n \rightarrow \infty} \mathcal{R}^*$$

# Consistance

- **faible** (pour  $P$ ) :

$$\mathbb{E} \left[ \mathcal{R}_P(\hat{f}(D_n)) \right] \xrightarrow{n \rightarrow \infty} \mathcal{R}^*$$

- **forte** (pour  $P$ ) :

$$\mathcal{R}_P(\hat{f}(D_n)) \xrightarrow[n \rightarrow \infty]{p.s.} \mathcal{R}^*$$

# Consistance

- faible (pour  $P$ ) :

$$\mathbb{E} \left[ \mathcal{R}_P(\hat{f}(D_n)) \right] \xrightarrow[n \rightarrow \infty]{} \mathcal{R}^*$$

- forte (pour  $P$ ) :

$$\mathcal{R}_P(\hat{f}(D_n)) \xrightarrow[n \rightarrow \infty]{p.s.} \mathcal{R}^*$$

- **consistance universelle** (faible ou forte) : pour tout  $P$ .

# Consistance

- faible (pour  $P$ ) :

$$\mathbb{E} \left[ \mathcal{R}_P(\hat{f}(D_n)) \right] \xrightarrow[n \rightarrow \infty]{} \mathcal{R}^*$$

- forte (pour  $P$ ) :

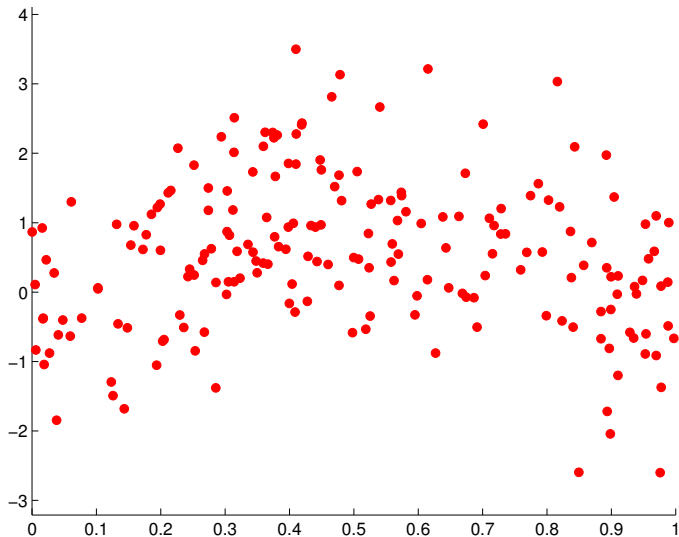
$$\mathcal{R}_P(\hat{f}(D_n)) \xrightarrow[n \rightarrow \infty]{p.s.} \mathcal{R}^*$$

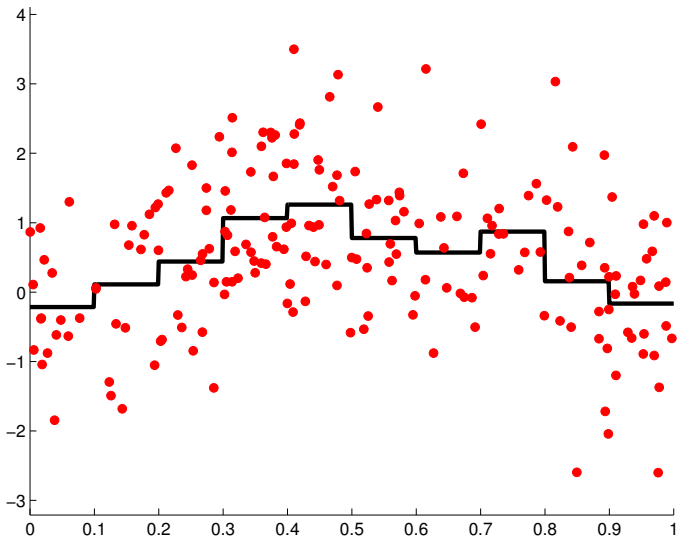
- **consistance universelle** (faible ou forte) : pour tout  $P$ .
- Attention ! Ne donne pas de **vitesse** de convergence !

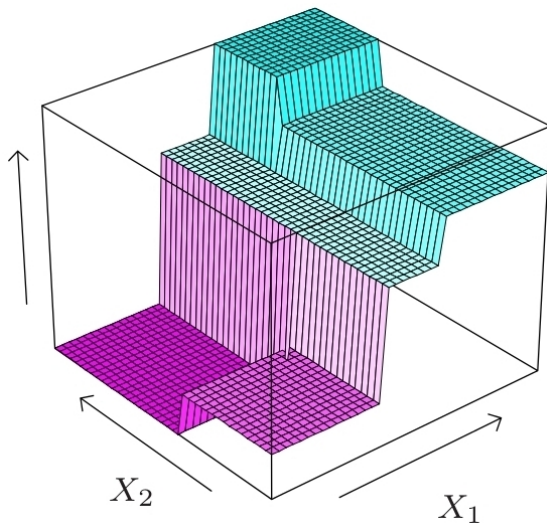


# Plan

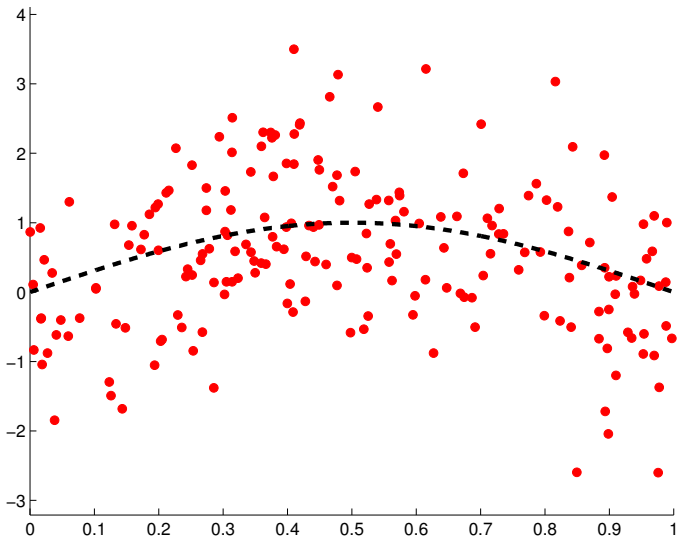
- 1 Prévision
- 2 Régression et classification
- 3 Minimisation du risque empirique
- 4 Moyennes locales
- 5 On n'a rien sans rien
- 6 Conclusion : enjeux de l'apprentissage

Régression :  $Y_i \in \mathcal{Y} = \mathbb{R}$ 

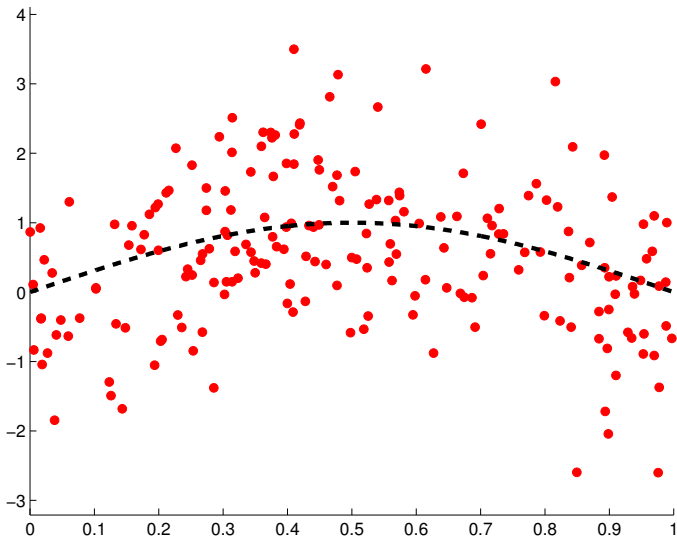
Règle par partition (régression),  $\mathcal{X} = \mathbb{R}$ 

Arbre de décision ( $\Rightarrow$  partition), régression,  $\mathcal{X} = \mathbb{R}^2$ 

# Fonction de régression : $\eta(X) := \mathbb{E}[Y | X]$



Régression :  $Y = \eta(X) + \varepsilon$  avec  $\mathbb{E}[\varepsilon | X] = 0$



# Règle par partition (régression)

- $\mathcal{A}$  : partition de  $\mathcal{X}$ , finie ou dénombrable.
- $\forall x \in \mathcal{X}$ ,  $\mathcal{A}(x)$  = l'unique élément de  $\mathcal{A}$  qui contient  $x$ .
- Règle de régression par partition associée à  $\mathcal{A}$  :

$\forall n \geq 1, x \in \mathcal{X}, (x_i, y_i)_{1 \leq i \leq n} \in (\mathcal{X} \times \mathcal{Y})^n$ ,

$$\hat{f}_{\mathcal{A}}^{\text{part-reg}}((x_i, y_i)_{1 \leq i \leq n}; x) := \frac{\sum_{i=1}^n y_i \mathbb{1}_{x_i \in \mathcal{A}(x)}}{\sum_{i=1}^n \mathbb{1}_{x_i \in \mathcal{A}(x)}}$$

- Convention :  $\frac{0}{0} = 0$

# Règle par partition (régression)

- $\mathcal{A}$  : partition de  $\mathcal{X}$ , finie ou dénombrable.
- $\forall x \in \mathcal{X}$ ,  $\mathcal{A}(x)$  = l'unique élément de  $\mathcal{A}$  qui contient  $x$ .
- Règle de régression par partition associée à  $\mathcal{A}$  :

$\forall n \geq 1$ ,  $x \in \mathcal{X}$ ,  $(x_i, y_i)_{1 \leq i \leq n} \in (\mathcal{X} \times \mathcal{Y})^n$ ,

$$\hat{f}_{\mathcal{A}}^{\text{part-reg}}((x_i, y_i)_{1 \leq i \leq n}; x) := \frac{\sum_{i=1}^n y_i \mathbb{1}_{x_i \in \mathcal{A}(x)}}{\sum_{i=1}^n \mathbb{1}_{x_i \in \mathcal{A}(x)}}$$

- Convention :  $\frac{0}{0} = 0$
- Partition cubique de pas  $h > 0$  sur  $\mathcal{X} = \mathbb{R}^p$  :

$$\mathcal{A}^{\text{cub}}(h) := \left( \prod_{i=1}^p [hk_i, h(k_i + 1)) \right)_{(k_1, \dots, k_p) \in \mathbb{Z}^p}$$



# Coût quadratique, régression idéale (1)

- **Coût quadratique** (des moindres carrés) :  $c(y, y') = (y - y')^2$
- Hypothèse :  $\mathbb{E}[Y^2] < +\infty$
- **Risque quadratique** de  $f \in \mathcal{F}$  :

$$\mathcal{R}_P(f) = \mathbb{E}[(f(X) - Y)^2]$$

# Coût quadratique, régression idéale (1)

- **Coût quadratique** (des moindres carrés) :  $c(y, y') = (y - y')^2$
- Hypothèse :  $\mathbb{E}[Y^2] < +\infty$
- **Risque quadratique** de  $f \in \mathcal{F}$  :

$$\begin{aligned}\mathcal{R}_P(f) &= \mathbb{E}\left[(f(X) - Y)^2\right] \\ &= \mathbb{E}\left[(f(X) - \eta(X) - \varepsilon)^2\right]\end{aligned}$$

# Coût quadratique, régression idéale (1)

- **Coût quadratique** (des moindres carrés) :  $c(y, y') = (y - y')^2$
- Hypothèse :  $\mathbb{E}[Y^2] < +\infty$
- **Risque quadratique** de  $f \in \mathcal{F}$  :

$$\begin{aligned}\mathcal{R}_P(f) &= \mathbb{E}\left[(f(X) - Y)^2\right] \\ &= \mathbb{E}\left[(f(X) - \eta(X) - \varepsilon)^2\right] \\ &= \mathbb{E}\left[(f(X) - \eta(X))^2\right] + \mathbb{E}[\varepsilon^2] + 2\mathbb{E}\left[(f(X) - \eta(X))\varepsilon\right]\end{aligned}$$

## Coût quadratique, régression idéale (1)

- **Coût quadratique** (des moindres carrés) :  $c(y, y') = (y - y')^2$
- Hypothèse :  $\mathbb{E}[Y^2] < +\infty$
- **Risque quadratique** de  $f \in \mathcal{F}$  :

$$\begin{aligned}
 \mathcal{R}_P(f) &= \mathbb{E}\left[(f(X) - Y)^2\right] \\
 &= \mathbb{E}\left[(f(X) - \eta(X) - \varepsilon)^2\right] \\
 &= \mathbb{E}\left[(f(X) - \eta(X))^2\right] + \mathbb{E}[\varepsilon^2] + \underbrace{2\mathbb{E}\left[(f(X) - \eta(X))\varepsilon\right]}_{=0}
 \end{aligned}$$

# Coût quadratique, régression idéale (1)

- **Coût quadratique** (des moindres carrés) :  $c(y, y') = (y - y')^2$
- Hypothèse :  $\mathbb{E}[Y^2] < +\infty$
- **Risque quadratique** de  $f \in \mathcal{F}$  :

$$\begin{aligned}\mathcal{R}_P(f) &= \mathbb{E}\left[(f(X) - Y)^2\right] \\ &= \mathbb{E}\left[(f(X) - \eta(X) - \varepsilon)^2\right] \\ &= \mathbb{E}\left[(f(X) - \eta(X))^2\right] + \mathbb{E}[\varepsilon^2]\end{aligned}$$

⇒ **minimal si et seulement si  $f(X) = \eta(X)$  p.s.**

# Coût quadratique, régression idéale (1)

- **Coût quadratique** (des moindres carrés) :  $c(y, y') = (y - y')^2$
- Hypothèse :  $\mathbb{E}[Y^2] < +\infty$
- **Risque quadratique** de  $f \in \mathcal{F}$  :

$$\begin{aligned} \mathcal{R}_P(f) &= \mathbb{E}[(f(X) - Y)^2] \\ &= \mathbb{E}[(f(X) - \eta(X) - \varepsilon)^2] \\ &= \mathbb{E}[(f(X) - \eta(X))^2] + \mathbb{E}[\varepsilon^2] \end{aligned}$$

⇒ minimal si et seulement si  $f(X) = \eta(X)$  p.s.

⇒ **excès de risque**  $\mathbb{E}[(f(X) - \eta(X))^2]$

## Coût quadratique, régression idéale (2)

## Proposition (2.1)

- (i)  $\eta$  est un prédicteur de Bayes
- (ii)  $f \in \mathcal{F}$  est un prédicteur de Bayes  $\Leftrightarrow f(X) = \eta(X)$  p.s.
- (iii) Risque de Bayes :

$$\mathcal{R}_P^* = \mathbb{E}[(Y - \eta(X))^2] = \text{var}(Y | X) = \mathbb{E}[\varepsilon^2]$$

- (iv) Excès de risque de  $f \in \mathcal{F}$  :

$$\ell(f^*, f) = \mathbb{E}[(f(X) - \eta(X))^2] = \|f - \eta\|_{L^2(P_X)}^2$$

## Coût quadratique, régression idéale (2)

## Proposition (2.1)

- (i)  $\eta$  est un prédicteur de Bayes
- (ii)  $f \in \mathcal{F}$  est un prédicteur de Bayes  $\Leftrightarrow f(X) = \eta(X)$  p.s.
- (iii) Risque de Bayes :

$$\mathcal{R}_P^* = \mathbb{E}[(Y - \eta(X))^2] = \text{var}(Y | X) = \mathbb{E}[\varepsilon^2]$$

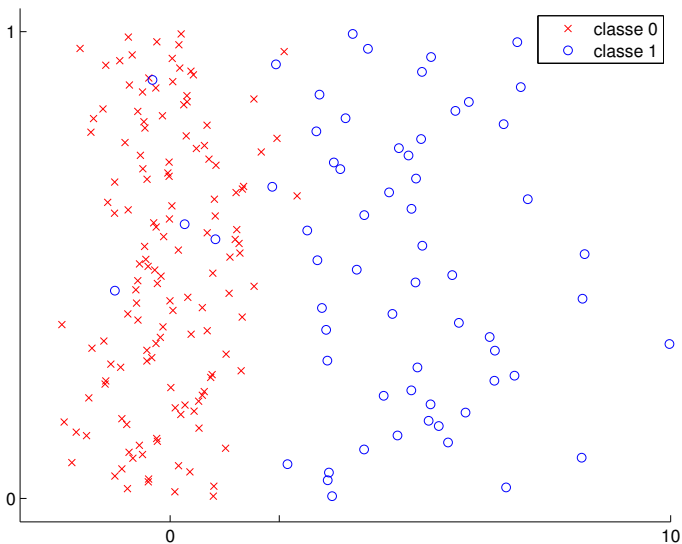
- (iv) Excès de risque de  $f \in \mathcal{F}$  :

$$\ell(f^*, f) = \mathbb{E}[(f(X) - \eta(X))^2] = \|f - \eta\|_{L^2(P_X)}^2$$

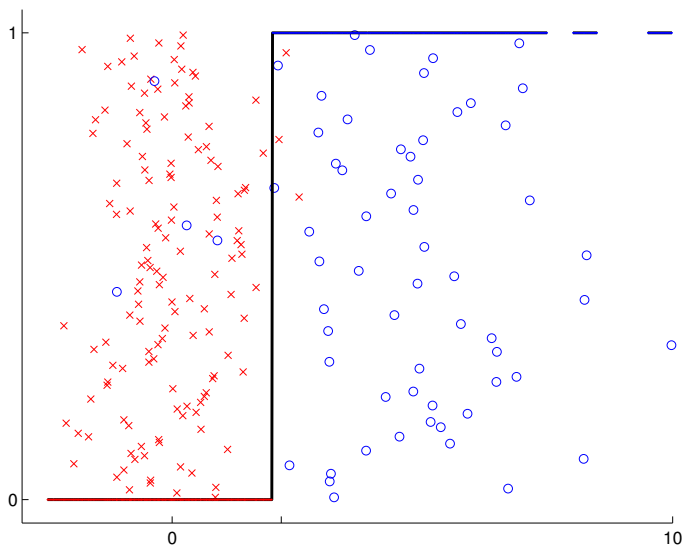
Exercice (2.1–2.2) : avec le coût valeur absolue  $c(y, y') = |y - y'|$  ?



# Classification (binaire supervisée) : $Y_i \in \mathcal{Y} = \{0, 1\}$



# Règle par partition (classification)



# Règle par partition (classification)

- $\mathcal{A}$  : partition de  $\mathcal{X}$ , finie ou dénombrable.
- $\forall x \in \mathcal{X}$ ,  $\mathcal{A}(x)$  = l'unique élément de  $\mathcal{A}$  qui contient  $x$ .
- Règle de classification par partition associée à  $\mathcal{A}$  :

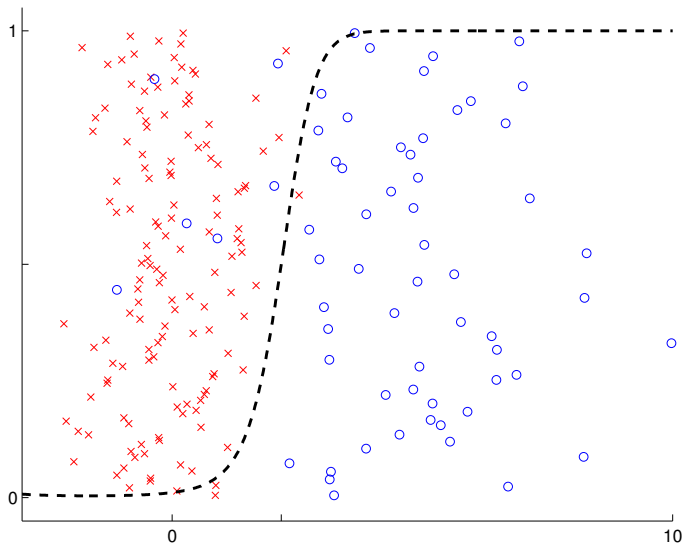
$\forall n \geq 1$ ,  $x \in \mathcal{X}$ ,  $(x_i, y_i)_{1 \leq i \leq n} \in (\mathcal{X} \times \mathcal{Y})^n$ ,

$$\hat{f}_{\mathcal{A}}^{\text{part-class}}((x_i, y_i)_{1 \leq i \leq n}; x) := \begin{cases} 1 & \text{si } \text{Card}\{i / y_i = 1 \text{ et } x_i \in \mathcal{A}(x)\} \\ & > \text{Card}\{i / y_i = 0 \text{ et } x_i \in \mathcal{A}(x)\} \\ 0 & \text{sinon} \end{cases}$$

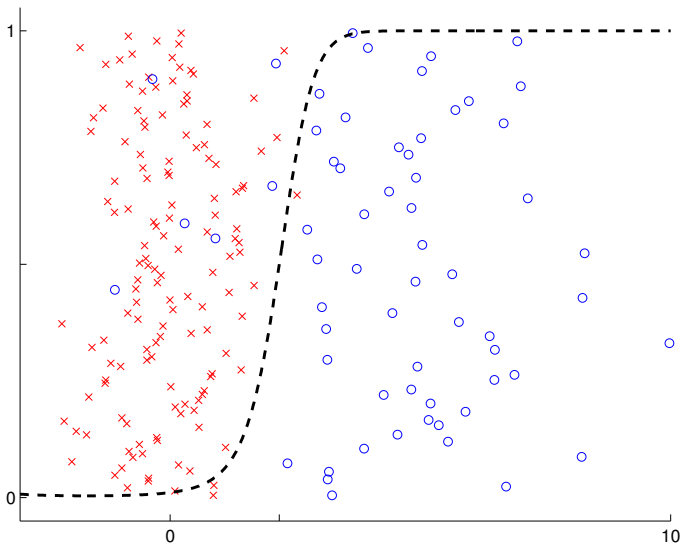
$\Leftrightarrow$  vote majoritaire dans chaque case, avec 0 par défaut

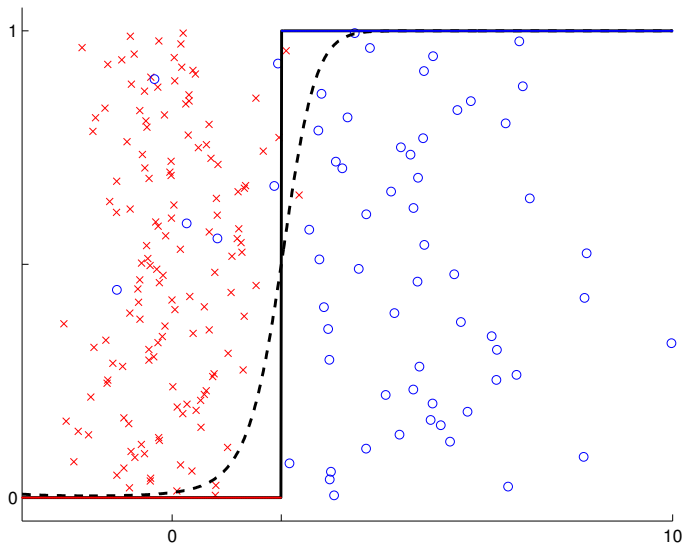
- Partition cubique de pas  $h > 0$  sur  $\mathcal{X} = \mathbb{R}^p$  :

$$\mathcal{A}^{\text{cub}}(h) := \left( \prod_{i=1}^p [hk_i, h(k_i + 1)) \right)_{(k_1, \dots, k_p) \in \mathbb{Z}^p}$$

Fonction de régression :  $\eta(X) := \mathbb{E}[Y | X] = \mathbb{P}(Y = 1 | X)$ 

# Coût 0–1 $c(y, y') = \mathbb{1}_{y \neq y'}$ : classification idéale ?



Coût 0-1  $c(y, y') = \mathbb{1}_{y \neq y'}$  : classification idéale

# Coût 0–1 : classification idéale

Risque 0–1 :  $\mathcal{R}_P : f \in \mathcal{F} \mapsto \mathbb{E}[\mathbb{1}_{f(X) \neq Y}] = \mathbb{P}(f(X) \neq Y)$

## Proposition (2.2)

- (i)  $f^* : x \mapsto \mathbb{1}_{\eta(x) > 1/2}$  est un *classifieur de Bayes*
- (ii)  $f \in \mathcal{F}$  est un *classifieur de Bayes*  $\Leftrightarrow f(X) = \mathbb{1}_{\eta(X) > 1/2}$  p.s.,  
sauf éventuellement sur  $\{\eta(X) = 1/2\}$
- (iii) *Risque de Bayes* :

$$\mathcal{R}_P^* = \mathbb{E}[\min\{\eta(X), 1 - \eta(X)\}]$$

- (iv) *Excès de risque* de  $f \in \mathcal{F}$  :

$$\ell(f^*, f) = \mathbb{E}[\lvert 2\eta(X) - 1 \rvert \mathbb{1}_{f^*(X) \neq f(X)}]$$

# Coût asymétrique

- **Coût asymétrique** :  $c_w : (y, y') \in \{0, 1\} \mapsto w_{y'} \mathbb{1}_{y \neq y'}$  avec  $w = (w_0, w_1) \in [0, +\infty[^2$  et  $w_0 + w_1 > 0$
- **Motivations** : spams, diagnostic médical, etc.
- Risque associé :

$$\begin{aligned} \mathcal{R}_P^w(f) &= \mathbb{E}[w_Y \mathbb{1}_{f(X) \neq Y}] \\ &= w_1 \mathbb{P}(f(X) = 0 \text{ et } Y = 1) + w_0 \mathbb{P}(f(X) = 1 \text{ et } Y = 0). \end{aligned}$$

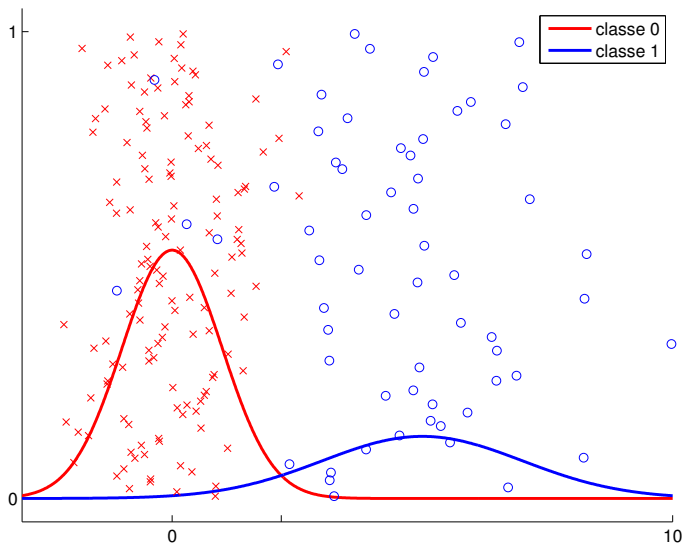
- **Classifieur de Bayes** (Proposition 2.3) :

$$f_w^* : x \mapsto \mathbb{1}_{\eta(x) > \frac{w_0}{w_0 + w_1}}$$

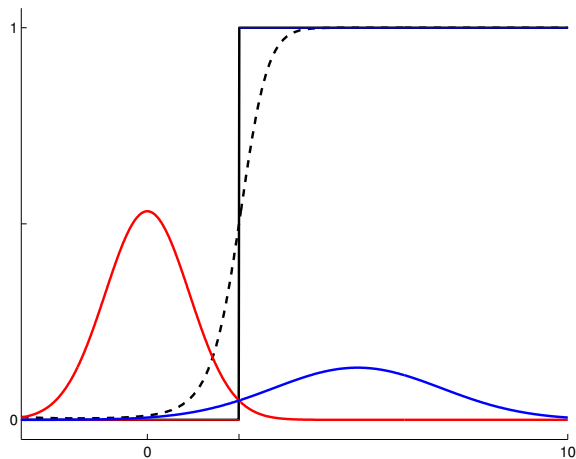
- ⇔ **test de rapport de vraisemblance** de  $H_0 : \ll x \sim \mathcal{L}(X | Y = 0) \gg$  contre  $H_1 : \ll x \sim \mathcal{L}(X | Y = 1) \gg$ .



## Classification binaire et test d'hypothèse (Rq. 2.21)



## Classification binaire et test d'hypothèse (Rq. 2.21)



Classifieur de  
Bayes  $f_w^*$

$\Leftrightarrow$

test de rapport de  
vraisemblance

$$\mathbb{1} \left\{ \frac{g_1(x)}{g_0(x)} \geq \frac{\mathbb{P}(Y=0)w_0}{\mathbb{P}(Y=1)w_1} \right\}$$

[Figure :  $\mathbb{P}(Y=0) = \mathbb{P}(Y=1) = 1/2$ ,  $w_0 = w_1 = 1$ ]

# Classifieur « plug-in »

- Idée :

$$s^*(x) = \mathbb{1}_{\eta(x) > \frac{1}{2}}$$

⇒ si  $\hat{\eta}(D_n)$  estime  $\eta$  (régression),

$$\hat{f}_{\hat{\eta}}(D_n; x) = \mathbb{1}_{\hat{\eta}(D_n; x) > \frac{1}{2}}$$

est la règle de classification par plug-in associée à  $\hat{\eta}$

# Classifieur « plug-in »

- Idée :

$$s^*(x) = \mathbb{1}_{\eta(x) > \frac{1}{2}}$$

⇒ si  $\hat{\eta}(D_n)$  estime  $\eta$  (régression),

$$\hat{f}_{\hat{\eta}}(D_n; x) = \mathbb{1}_{\hat{\eta}(D_n; x) > \frac{1}{2}}$$

est la **règle de classification par plug-in associée à  $\hat{\eta}$**

- **Exemples :**
  - règles par partitions :

$$\hat{f}_{\mathcal{A}}^{\text{part-class}}(D_n; x) = \mathbb{1}_{\hat{f}_{\mathcal{A}}^{\text{part-reg}}(D_n; x) > 1/2}$$

- $k$  plus proches voisins
- moyennes locales (Section 2.4)

# Borne de risque pour un classifieur « plug-in »

## Proposition (2.4)

*Pour le coût 0–1 en classification :*

$$\ell(f^*, \hat{f}_{\hat{\eta}}) \leq 2\mathbb{E}|\hat{\eta}(X) - \eta(X)| \leq 2\sqrt{\mathbb{E}[(\hat{\eta}(X) - \eta(X))^2]}$$

## Borne de risque pour un classifieur « plug-in »

## Proposition (2.4)

Pour le coût 0–1 en classification :

$$\ell(f^*, \hat{f}_{\hat{\eta}}) \leq 2\mathbb{E}|\hat{\eta}(X) - \eta(X)| \leq 2\sqrt{\mathbb{E}[(\hat{\eta}(X) - \eta(X))^2]}$$

En particulier :

$\hat{\eta}$  faiblement consistante (pour  $P$ , en régression avec le coût quadratique)

$\Rightarrow \hat{f}_{\hat{\eta}}$  faiblement consistante (pour  $P$ , en classification avec le coût 0–1).

## Coûts convexes : pseudo-classifieurs et $\Phi$ -risque

- Pourquoi un coût convexe ?

# Coûts convexes : pseudo-classifieurs et $\Phi$ -risque

- Pourquoi un coût convexe ?
- Convention (ici seulement) :  $Y_i \in \mathcal{Y} = \{-1, 1\}$



# Coûts convexes : pseudo-classifieurs et $\Phi$ -risque

- Pourquoi un coût convexe ?
  - Convention (ici seulement) :  $Y_i \in \mathcal{Y} = \{-1, 1\}$
  - « **Pseudo-classifieur** »  $g \in \overline{\mathcal{F}}$  : fonction  $\mathcal{X} \rightarrow \overline{\mathbb{R}}$
- $\Rightarrow$  classifieur associé :

$$\text{signe}(g) : x \in \mathcal{X} \mapsto \begin{cases} 1 & \text{si } g(x) > 0 \\ -1 & \text{sinon} \end{cases}$$

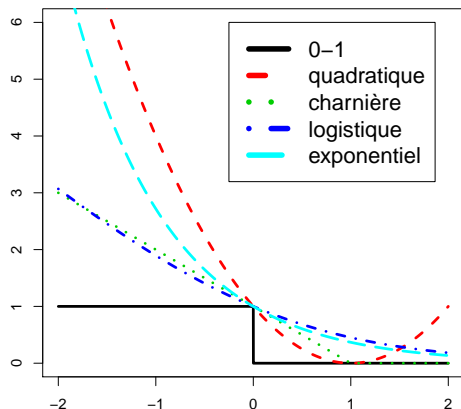
Coûts convexes : pseudo-classifieurs et  $\Phi$ -risque

- Pourquoi un coût convexe ?
  - Convention (ici seulement) :  $Y_i \in \mathcal{Y} = \{-1, 1\}$
  - « **Pseudo-classifieur** »  $g \in \overline{\mathcal{F}}$  : fonction  $\mathcal{X} \rightarrow \overline{\mathbb{R}}$
- $\Rightarrow$  classifieur associé :

$$\text{signe}(g) : x \in \mathcal{X} \mapsto \begin{cases} 1 & \text{si } g(x) > 0 \\ -1 & \text{sinon} \end{cases}$$

- $\Phi : \overline{\mathbb{R}} \rightarrow [0, +\infty] \Rightarrow$   **$\Phi$ -risque** de  $g \in \overline{\mathcal{F}}$  :

$$\mathcal{R}_P^\Phi(g) := \mathbb{E}[\Phi(Yg(X))]$$

Coûts convexes : Exemples de fonctions  $\Phi$ 

- **0-1** :  $\Phi_{0-1}(u) = \mathbb{1}_{u \leq 0}$
- **quadratique** :  
 $\Phi(u) = (1 - u)^2$
- **charnière (hinge)** :  
 $\Phi(u) = \max\{1 - u, 0\}$   
 $\Rightarrow$  SVM
- **logistique** :  
 $\Phi(u) = \ln_2(1 + e^{-u})$   
 $\Rightarrow$  régression logistique
- **exponentiel** :  $\Phi(u) = e^{-u}$   
 $\Rightarrow$  boosting, AdaBoost

$$\mathcal{R}_P^\Phi(g) := \mathbb{E}[\Phi(Yg(X))]$$

# Lien entre $\Phi$ -risque et risque 0–1

Définition (Définition 2.3 & Proposition 2.6)

$\Phi$  est **calibrée pour la classification** si pour toute loi  $P$ ,  
 $g_{\Phi}^* \in \operatorname{argmin}_{g \in \overline{\mathcal{F}}} \mathcal{R}_P^{\Phi}(g) \Rightarrow \operatorname{sign}(g_{\Phi}^*)$  classifieur de Bayes pour le  
risque 0–1.

Lien entre  $\Phi$ -risque et risque 0–1

## Définition (Définition 2.3 &amp; Proposition 2.6)

$\Phi$  est **calibrée pour la classification** si pour toute loi  $P$ ,  
 $g_\Phi^* \in \operatorname{argmin}_{g \in \mathcal{F}} \mathcal{R}_P^\Phi(g) \Rightarrow \operatorname{sign}(g_\Phi^*)$  classifieur de Bayes pour le  
 risque 0–1.

## Proposition (2.7)

*On suppose  $\Phi$  convexe. Alors :*

*$\Phi$  calibrée pour la classification  $\Leftrightarrow \Phi$  dérivable en 0 et  $\Phi'(0) < 0$ .*

Exemples :

- quadratique
- charnière
- logistique
- exponentiel

Lien entre  $\Phi$ -risque et risque 0-1

## Définition (Définition 2.3 &amp; Proposition 2.6)

$\Phi$  est **calibrée pour la classification** si pour toute loi  $P$ ,  
 $g_{\Phi}^* \in \operatorname{argmin}_{g \in \overline{\mathcal{F}}} \mathcal{R}_P^{\Phi}(g) \Rightarrow \operatorname{signe}(g_{\Phi}^*)$  classifieur de Bayes pour le  
 risque 0-1.

## Proposition (2.7)

*On suppose  $\Phi$  convexe. Alors :*

$\Phi$  calibrée pour la classification  $\Leftrightarrow \Phi$  dérivable en 0 et  $\Phi'(0) < 0$ .

## Théorème (2.1)

*Si  $\Phi$  est convexe et calibrée pour la classification, il existe une  
 fonction  $\Psi : [-1, 1] \rightarrow [0, +\infty[$  telle que :*

$$\forall P, \forall g \in \overline{\mathcal{F}}, \quad \Psi\left(\mathcal{R}_P^{0-1}(\operatorname{signe}(g)) - \mathcal{R}_P^{0-1*}\right) \leq \mathcal{R}_P^{\Phi}(g) - \mathcal{R}_P^{\Phi*}.$$

Lien entre  $\Phi$ -risque et risque 0–1

## Théorème (2.1)

Si  $\Phi$  est convexe et calibrée pour la classification, il existe une fonction  $\Psi : [-1, 1] \rightarrow [0, +\infty[$  telle que :

$$\forall P, \forall g \in \overline{\mathcal{F}}, \quad \Psi\left(\mathcal{R}_P^{0-1}(\text{signe}(g)) - \mathcal{R}_P^{0-1*}\right) \leq \mathcal{R}_P^\Phi(g) - \mathcal{R}_P^{\Phi*}.$$

Exemples :

- quadratique :  $\Psi(\theta) = \theta^2$
- charnière :  $\Psi(\theta) = |\theta|$
- logistique :  $\Psi(\theta) \geq \theta^2 / (2 \ln 2)$
- exponentiel :  $\Psi(\theta) = 1 - \sqrt{1 - \theta^2} \geq \theta^2 / 2$

# Plan

- 1 Prévision
- 2 Régression et classification
- 3 Minimisation du risque empirique
- 4 Moyennes locales
- 5 On n'a rien sans rien
- 6 Conclusion : enjeux de l'apprentissage



# Minimisation du risque empirique

- **Risque empirique** de  $f \in \mathcal{F}$  :

$$\widehat{\mathcal{R}}_n(f) = \frac{1}{n} \sum_{i=1}^n c(f(X_i), Y_i)$$

- Pour tout  $f \in \mathcal{F}$  :

$$\mathbb{E}[\widehat{\mathcal{R}}_n(f)] = \mathbb{E}[c(f(X), Y)] = \mathcal{R}_P(f)$$

# Minimisation du risque empirique

- **Risque empirique** de  $f \in \mathcal{F}$  :

$$\widehat{\mathcal{R}}_n(f) = \frac{1}{n} \sum_{i=1}^n c(f(X_i), Y_i)$$

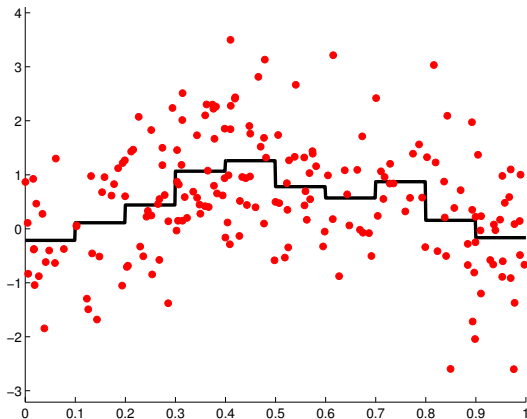
- Pour tout  $f \in \mathcal{F}$  :

$$\mathbb{E}[\widehat{\mathcal{R}}_n(f)] = \mathbb{E}[c(f(X), Y)] = \mathcal{R}_P(f)$$

- **Minimiseur du risque empirique** sur le **modèle**  $S \subset \mathcal{F}$  :

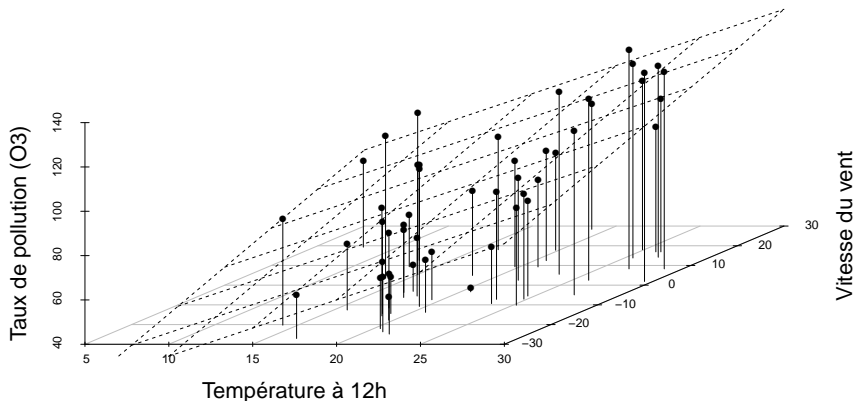
$$\widehat{f}_S \in \operatorname{argmin}_{f \in S} \{\widehat{\mathcal{R}}_n(f)\}$$

## Exemple : partition en régression



$$S_{\text{reg}}^{\text{part}}(\mathcal{A}) = \overline{\text{vect}\{\mathbb{1}_A, A \in \mathcal{A}\}}$$

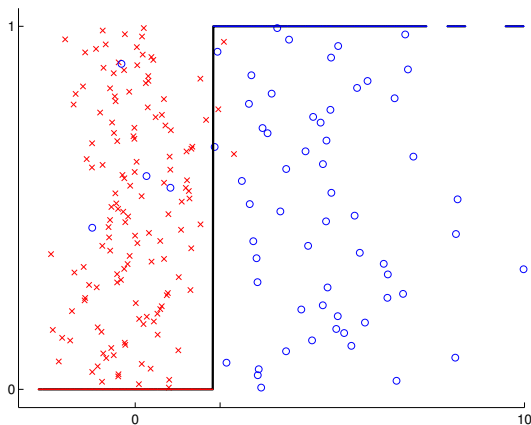
## Exemple : régression linéaire



$$S^{\text{lin}} = \left\{ x \in \mathcal{X} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle = \sum_{j=1}^p w_j x_j / \mathbf{w} \in \mathbb{R}^p \right\}$$

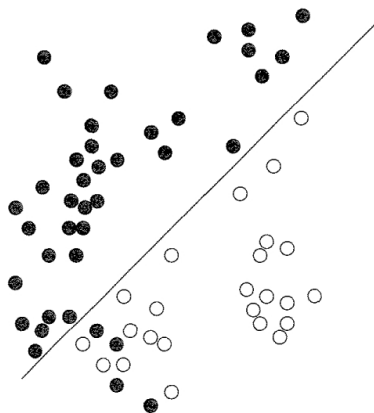
Figure : [Cornillon and Matzner-Løber, 2011]

## Exemple : partition en classification



$$S_{\text{class}}^{\text{part}}(\mathcal{A}) = S_{\text{reg}}^{\text{part}}(\mathcal{A}) \cap \mathcal{F}(\mathcal{X}, \{0, 1\})$$

## Exemple : discrimination linéaire



$$S_{\text{class}}^{\text{lin}} = \{x \in \mathcal{X} \mapsto \mathbb{1}_{\langle \mathbf{w}, x \rangle \geq 0} / \mathbf{w} \in \mathbb{R}^P\}$$

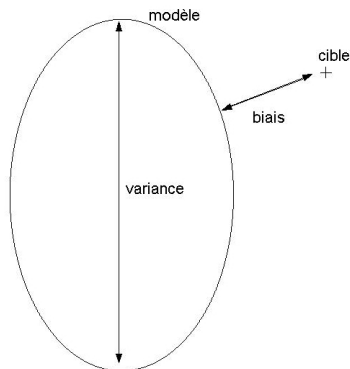
## Erreurs d'approximation et d'estimation

**Erreur d'approximation** (ou biais)

$$l(f^*, S) := \inf_{f \in S} l(f^*, f)$$

**Erreur d'estimation** (ou variance)

$$\begin{aligned} & \mathcal{R}_P(\hat{f}_S) - \inf_{f \in S} \mathcal{R}_P(f) \\ &= l(f^*, \hat{f}_S) - l(f^*, S) \geq 0 \end{aligned}$$



$$l(f^*, \hat{f}_S) = \underbrace{l(f^*, S)}_{\text{erreur d'approximation}} + \underbrace{l(f^*, \hat{f}_S) - l(f^*, S)}_{\text{erreur d'estimation}}$$

# Majoration générale de l'erreur d'estimation

Pour tout  $g \in \mathcal{S}$ , on a :

$$\begin{aligned}
 & \mathcal{R}_P(\hat{f}_S) - \mathcal{R}_P(g) \\
 = & \underbrace{\mathcal{R}_P(\hat{f}_S) - \hat{\mathcal{R}}_n(\hat{f}_S)}_{\leq \sup_{f \in \mathcal{S}} |\mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f)|} + \underbrace{\hat{\mathcal{R}}_n(\hat{f}_S) - \hat{\mathcal{R}}_n(g)}_{\leq 0} + \underbrace{\hat{\mathcal{R}}_n(g) - \mathcal{R}_P(g)}_{\leq \sup_{f \in \mathcal{S}} |\mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f)|} \\
 \leq & 2 \sup_{f \in \mathcal{S}} |\mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f)|.
 \end{aligned}$$



# Majoration générale de l'erreur d'estimation

Pour tout  $g \in S$ , on a :

$$\begin{aligned}
 & \mathcal{R}_P(\hat{f}_S) - \mathcal{R}_P(g) \\
 = & \underbrace{\mathcal{R}_P(\hat{f}_S) - \hat{\mathcal{R}}_n(\hat{f}_S)}_{\leq \sup_{f \in S} |\mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f)|} + \underbrace{\hat{\mathcal{R}}_n(\hat{f}_S) - \hat{\mathcal{R}}_n(g)}_{\leq 0} + \underbrace{\hat{\mathcal{R}}_n(g) - \mathcal{R}_P(g)}_{\leq \sup_{f \in S} |\mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f)|} \\
 \leq & 2 \sup_{f \in S} |\mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f)|.
 \end{aligned}$$

En prenant le sup sur  $g \in S$ , on obtient :

## Proposition (2.8)

Si  $\hat{f}_S$  minimise le risque empirique sur  $S$  :

$$\ell(f^*, \hat{f}_S) - \ell(f^*, S) \leq 2 \sup_{f \in S} |\mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f)|.$$

# Majoration générale de l'erreur d'estimation

## Proposition (2.8)

Si  $\hat{f}_S$  minimise le risque empirique sur  $S$  :

$$\ell(f^*, \hat{f}_S) - \ell(f^*, S) \leq 2 \sup_{f \in S} |\mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f)|.$$

## Proposition (2.9)

Si  $\hat{f}_S$  minimise le risque empirique sur  $S$  :

$$\mathbb{E}[\ell(f^*, \hat{f}_S) - \ell(f^*, S)] \leq \mathbb{E} \left[ \sup_{f \in S} \{\mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f)\} \right].$$

# Erreur d'estimation et **minimisation approchée**

## Proposition (2.8)

Si  $\hat{f}_{S,\rho}$  minimise le risque empirique sur  $S$  à  $\rho$  près :

$$\ell(f^*, \hat{f}_{S,\rho}) - \ell(f^*, S) \leq 2 \sup_{f \in S} |\mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f)| + \rho.$$

## Proposition (2.9)

Si  $\hat{f}_{S,\rho}$  minimise le risque empirique sur  $S$  à  $\rho$  près :

$$\mathbb{E}[\ell(f^*, \hat{f}_{S,\rho}) - \ell(f^*, S)] \leq \mathbb{E} \left[ \sup_{f \in S} \{ \mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f) \} \right] + \rho.$$

# Coût borné, modèle fini

## Proposition (2.11)

Si  $\hat{f}_S$  minimise le risque empirique sur  $S$  et si

$$c(f(X), Y) \in [a, b] \quad \text{presque sûrement,}$$

alors,

$$\mathbb{E}[\ell(f^*, \hat{f}_S)] \leq \ell(f^*, S) + (b - a) \sqrt{\frac{\ln(\text{Card } S)}{2n}}.$$

## Coût borné, modèle fini

## Proposition (2.11)

Si  $\hat{f}_S$  minimise le risque empirique sur  $S$  et si

$$c(f(X), Y) \in [a, b] \quad \text{presque sûrement,}$$

alors,

$$\mathbb{E}[\ell(f^*, \hat{f}_S)] \leq \ell(f^*, S) + (b - a) \sqrt{\frac{\ln(\text{Card } S)}{2n}}.$$

Application : règle par **partition finie** en classification 0-1 :

$$\mathbb{E}[\ell(f^*, \hat{f}_A^{\text{part-class}})] \leq \ell(f^*, S_{\text{class}}^{\text{part}}(\mathcal{A})) + \sqrt{\frac{\ln(2)}{2}} \sqrt{\frac{\text{Card}(\mathcal{A})}{n}}$$

# Classification 0–1, modèle quelconque

- **Tout se passe comme si  $S$  était fini** de cardinal :

$$\text{Card}\left\{ (f(X_i))_{1 \leq i \leq n} / f \in S \right\} =: \exp(H_S(X_1, \dots, X_n))$$

# Classification 0–1, modèle quelconque

- Tout se passe comme si  $S$  était fini de cardinal :

$$\text{Card}\left\{ (f(X_i))_{1 \leq i \leq n} / f \in S \right\} =: \exp(H_S(X_1, \dots, X_n))$$

- On a toujours :

$$\sup_{x_1, \dots, x_n \in \mathcal{X}} \text{Card}\left\{ (f(x_i))_{1 \leq i \leq n} / f \in S \right\} \leq 2^n$$

# Classification 0–1, modèle quelconque

- Tout se passe comme si  $S$  était fini de cardinal :

$$\text{Card}\left\{ (f(X_i))_{1 \leq i \leq n} / f \in S \right\} =: \exp(H_S(X_1, \dots, X_n))$$

- On a toujours :

$$\sup_{x_1, \dots, x_n \in \mathcal{X}} \text{Card}\left\{ (f(x_i))_{1 \leq i \leq n} / f \in S \right\} \leq 2^n$$

- **Classe de Vapnik-Chervonenkis de dimension  $V(S)$  :**
  - égalité pour  $n = 1, \dots, V(S)$  ( $\exists x_1, \dots, x_n$  « pulvérisés » par  $S$ )
  - **inégalité stricte pour  $n > V(S)$**
- **Exemple : discrimination linéaire dans  $\mathbb{R}^d$  :  $V(S) = d$**



# Classification 0–1, modèle quelconque

- Tout se passe comme si  $S$  était fini de cardinal :

$$\text{Card}\left\{ (f(X_i))_{1 \leq i \leq n} / f \in S \right\} =: \exp(H_S(X_1, \dots, X_n))$$

- On a toujours :

$$\sup_{x_1, \dots, x_n \in \mathcal{X}} \text{Card}\left\{ (f(x_i))_{1 \leq i \leq n} / f \in S \right\} \leq 2^n$$

- **Classe de Vapnik-Chervonenkis de dimension  $V(S)$**  :
  - égalité pour  $n = 1, \dots, V(S)$  ( $\exists x_1, \dots, x_n$  « pulvérisés » par  $S$ )
  - inégalité stricte pour  $n > V(S)$
- Exemple : discrimination linéaire dans  $\mathbb{R}^d$  :  $V(S) = d$
- **Lemme de Sauer** :  
 $n > 2V(S) \Rightarrow H_S(x_1, \dots, x_n) \leq V(S) \ln(en/V(S))$

## Classification 0-1, modèle quelconque : récapitulatif

$$\begin{aligned}
& \mathbb{E}[\ell(f^*, \hat{f}_S)] \\
& \leq \ell(f^*, S) + \mathbb{E} \left[ \sup_{f \in S} \{ \mathcal{R}_P(f) - \hat{\mathcal{R}}_n(f) \} \right] \quad (\text{Prop. 2.9}) \\
& \leq \ell(f^*, S) + 2 \mathbb{E} \left[ \sup_{f \in S} \left\{ \frac{1}{n} \sum_{i=1}^n \varepsilon_i c(f(X_i), Y_i) \right\} \right] \quad (\text{symétrisation}) \\
& \leq \ell(f^*, S) + \frac{2\sqrt{2}}{\sqrt{n}} \mathbb{E} \left[ \sqrt{H_S(X_1, \dots, X_n)} \right] \quad (\text{entropie combinatoire}) \\
& \leq \ell(f^*, S) + \frac{2\sqrt{2}}{\sqrt{n}} \sqrt{\sup_{x_1, \dots, x_n \in \mathcal{X}} H_S(x_1, \dots, x_n)} \\
& \leq \ell(f^*, S) + 2\sqrt{2} \sqrt{\frac{V(S)}{n} \ln \left( \frac{en}{V(S)} \right)}. \quad (\text{Vapnik-Chervonenkis})
\end{aligned}$$

# Classification zéro-erreur

## Proposition (2.12)

*Hypothèses :*

- *classification 0-1*
- $\eta(X) \in \{0, 1\}$  *p.s. (classification zéro-erreur),*
- $\hat{f}_S$  *minimise le risque empirique sur S fini avec  $f^* \in S$*

*Alors :*

$$\mathbb{E} \left[ \ell(f^*, \hat{f}_S(D_n)) \right] \leq \frac{1 + \ln(\text{Card}(S))}{n}.$$

# Classification zéro-erreur

## Proposition (2.12)

Hypothèses :

- classification 0-1
- $\eta(X) \in \{0, 1\}$  p.s. (classification zéro-erreur),
- $\hat{f}_S$  minimise le risque empirique sur  $S$  fini avec  $f^* \in S$

Alors :

$$\mathbb{E} \left[ \ell(f^*, \hat{f}_S(D_n)) \right] \leq \frac{1 + \ln(\text{Card}(S))}{n}.$$

Application : règle par **partition finie** en classification 0-1 :

$$\text{si } f^* \in S_{\text{class}}^{\text{part}}(\mathcal{A}), \quad \mathbb{E} \left[ \ell(f^*, \hat{f}_{\mathcal{A}}^{\text{part-class}}) \right] \leq \frac{1 + \ln(2) \text{Card}(\mathcal{A})}{n}.$$

# Classification zéro-erreur

## Proposition (2.12)

Hypothèses :

- classification 0-1
- $\eta(X) \in \{0, 1\}$  p.s. (classification zéro-erreur),
- $\hat{f}_S$  minimise le risque empirique sur  $S$  fini avec  $f^* \in S$

Alors :

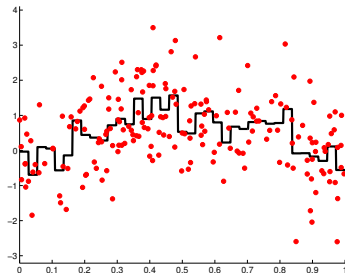
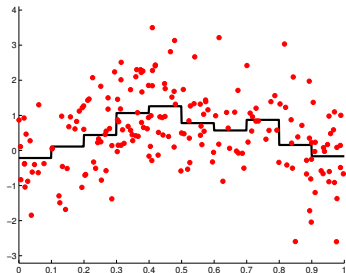
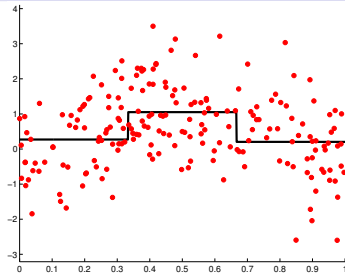
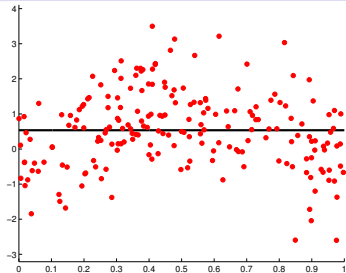
$$\mathbb{E} \left[ \ell(f^*, \hat{f}_S(D_n)) \right] \leq \frac{1 + \ln(\text{Card}(S))}{n}.$$

Application : règle par **partition finie** en classification 0-1 :

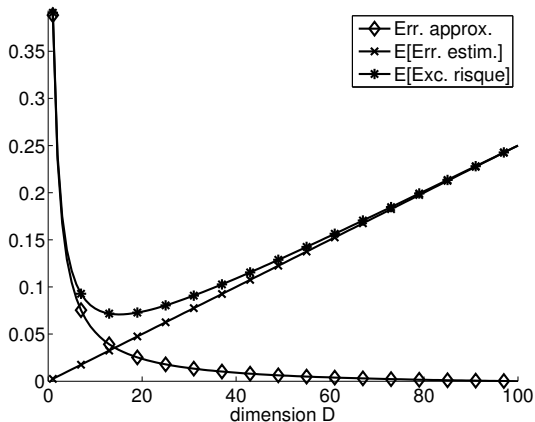
$$\text{si } f^* \in S_{\text{class}}^{\text{part}}(\mathcal{A}), \quad \mathbb{E} \left[ \ell(f^*, \hat{f}_{\mathcal{A}}^{\text{part-class}}) \right] \leq \frac{1 + \ln(2) \text{Card}(\mathcal{A})}{n}.$$

Se généralise à  $S$  quelconque (entropie combinatoire, dimension de Vapnik-Chervonenkis, ...).

# Sélection de modèles : partitions cubiques



# Compromis biais-variance



Sous-apprentissage

Sur-apprentissage

# Sélection de modèles

- famille de modèles  $(S_m)_{m \in \mathcal{M}} \Rightarrow$  famille des règles  $\hat{f}_m = \hat{f}_{S_m}$



# Sélection de modèles

- famille de modèles  $(S_m)_{m \in \mathcal{M}} \Rightarrow$  famille des règles  $\hat{f}_m = \hat{f}_{S_m}$
- $\Rightarrow$  Objectif de prévision : choisir  $\hat{m} = \hat{m}(D_n)$  tel que  $\mathcal{R}(\hat{f}_{\hat{m}(D_n)}(D_n))$  est minimal ?

# Sélection de modèles

- famille de modèles  $(S_m)_{m \in \mathcal{M}} \Rightarrow$  famille des règles  $\hat{f}_m = \hat{f}_{S_m}$
- $\Rightarrow$  Objectif de prévision : choisir  $\hat{m} = \hat{m}(D_n)$  tel que  $\mathcal{R}(\hat{f}_{\hat{m}(D_n)}(D_n))$  est minimal ?
- **Optimalité asymptotique** :

$$\frac{\ell(f^*, \hat{f}_{\hat{m}(D_n)}(D_n))}{\inf_{m \in \mathcal{M}} \left\{ \ell(f^*, \hat{f}_m(D_n)) \right\}} \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} 1.$$

- **Inégalité oracle** (en espérance / avec grande probabilité) :

$$\ell(f^*, \hat{f}_{\hat{m}(D_n)}(D_n)) \leq C_n \inf_{m \in \mathcal{M}} \left\{ \ell(f^*, \hat{f}_m(D_n)) \right\} + R_n$$

# Sélection de modèles

- famille de modèles  $(S_m)_{m \in \mathcal{M}} \Rightarrow$  famille des règles  $\hat{f}_m = \hat{f}_{S_m}$
- $\Rightarrow$  Objectif de prévision : choisir  $\hat{m} = \hat{m}(D_n)$  tel que  $\mathcal{R}(\hat{f}_{\hat{m}(D_n)}(D_n))$  est minimal ?
- **Optimalité asymptotique** :

$$\frac{\ell(f^*, \hat{f}_{\hat{m}(D_n)}(D_n))}{\inf_{m \in \mathcal{M}} \left\{ \ell(f^*, \hat{f}_m(D_n)) \right\}} \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} 1.$$

- **Inégalité oracle** (en espérance / avec grande probabilité) :

$$\ell(f^*, \hat{f}_{\hat{m}(D_n)}(D_n)) \leq C_n \inf_{m \in \mathcal{M}} \left\{ \ell(f^*, \hat{f}_m(D_n)) \right\} + R_n$$

- **Attention à l'interprétation de  $\hat{m}$**  : le meilleur modèle peut être faux / le vrai modèle peut être moins bon qu'un modèle plus simple.

# Méthodes de sélection de modèle

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \{\operatorname{crit}(m; D_n)\}$$

## Méthodes de sélection de modèle

$$\hat{m} \in \underset{m \in \mathcal{M}}{\operatorname{argmin}} \{ \operatorname{crit}(m; D_n) \}$$

- **Estimation sans biais du risque** (AIC,  $C_p$ , validation croisée...):

$$\mathbb{E}[\operatorname{crit}(m; D_n)] \approx \mathbb{E}[\mathcal{R}_P(\hat{f}_m(D_n))] \quad (\text{à translation près})$$

- ⇒  $\operatorname{crit}(m; D_n) \approx \mathcal{R}(\hat{f}_{\hat{m}(D_n)}(D_n))$  simultanément pour tous les  $m \in \mathcal{M}$  (si  $\mathcal{M}$  « petite »)
- ⇒ inégalité oracle « optimale »

## Méthodes de sélection de modèle

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \{\operatorname{crit}(m; D_n)\}$$

- Estimation sans biais du risque (AIC,  $C_p$ , validation croisée...):

$$\mathbb{E}[\operatorname{crit}(m; D_n)] \approx \mathbb{E}[\mathcal{R}_P(\hat{f}_m(D_n))] \quad (\text{à translation près})$$

- ⇒  $\operatorname{crit}(m; D_n) \approx \mathcal{R}(\hat{f}_{\hat{m}(D_n)}(D_n))$  simultanément pour tous les  $m \in \mathcal{M}$  (si  $\mathcal{M}$  « petite »)
- ⇒ inégalité oracle « optimale »

- Avec une **majoration du risque** :

$$\operatorname{crit}(m; D_n) \geq \mathcal{R}_P(\hat{f}_m(D_n)) \quad (\text{à translation près})$$

- ⇒ inégalité oracle **si la majoration n'est pas trop large**

# Plan

- 1 Prévision
- 2 Régression et classification
- 3 Minimisation du risque empirique
- 4 Moyennes locales
- 5 On n'a rien sans rien
- 6 Conclusion : enjeux de l'apprentissage

# Moyennes locales

## Définition (2.7)

$$\forall x, x_1, \dots, x_n \in \mathcal{X}, \forall i \in \{1, \dots, n\}, \quad W_i(x_{1\dots n}; x) \geq 0$$

- Régression :

$$\hat{\eta}((x_i, y_i)_{1 \leq i \leq n}; x) := \sum_{i=1}^n W_i(x_{1\dots n}; x) y_i$$



## Moyennes locales

## Définition (2.7)

$$\forall x, x_1, \dots, x_n \in \mathcal{X}, \forall i \in \{1, \dots, n\}, \quad W_i(x_{1\dots n}; x) \geq 0$$

- Régression :

$$\hat{\eta}((x_i, y_i)_{1 \leq i \leq n}; x) := \sum_{i=1}^n W_i(x_{1\dots n}; x) y_i$$

- Classification (plug-in) :

$$\hat{f}((x_i, y_i)_{1 \leq i \leq n}; x) := \mathbb{1}_{\hat{\eta}((x_i, y_i)_{1 \leq i \leq n}; x) > 1/2}$$

## Moyennes locales

## Définition (2.7)

$$\forall \mathbf{x}, x_1, \dots, x_n \in \mathcal{X}, \forall i \in \{1, \dots, n\}, \quad W_i(x_{1\dots n}; \mathbf{x}) \geq 0$$

- Régression :

$$\hat{\eta}((x_i, y_i)_{1 \leq i \leq n}; \mathbf{x}) := \sum_{i=1}^n W_i(x_{1\dots n}; \mathbf{x}) y_i$$

- Classification (plug-in) :

$$\hat{f}((x_i, y_i)_{1 \leq i \leq n}; \mathbf{x}) := \mathbb{1}_{\hat{\eta}((x_i, y_i)_{1 \leq i \leq n}; \mathbf{x}) > 1/2}$$

- Hypothèse générale :

$$\forall \mathbf{x}, x_1, \dots, x_n \in \mathcal{X}, \quad \sum_{i=1}^n W_i(x_{1\dots n}; \mathbf{x}) \approx 1$$

## Moyennes locales : exemples

- **Partitions :**

$$W_i^{\mathcal{A}}(x_{1\dots n}; \mathbf{x}) := \begin{cases} \frac{\mathbb{1}_{x_i \in \mathcal{A}(x)}}{\sum_{j=1}^n \mathbb{1}_{x_j \in \mathcal{A}(x)}} & \text{si } \sum_{j=1}^n \mathbb{1}_{x_j \in \mathcal{A}(x)} > 0 \\ 0 & \text{sinon.} \end{cases}$$

$$\Rightarrow \sum_{i=1}^n W_i^{\mathcal{A}}(x_{1\dots n}; \mathbf{x}) = \mathbb{1}_{\sum_{j=1}^n \mathbb{1}_{x_j \in \mathcal{A}(x)} > 0} \leq 1$$

# Moyennes locales : exemples

- **Partitions :**

$$W_i^{\mathcal{A}}(x_{1\dots n}; \mathbf{x}) := \begin{cases} \frac{\mathbb{1}_{x_i \in \mathcal{A}(x)}}{\sum_{j=1}^n \mathbb{1}_{x_j \in \mathcal{A}(x)}} & \text{si } \sum_{j=1}^n \mathbb{1}_{x_j \in \mathcal{A}(x)} > 0 \\ 0 & \text{sinon.} \end{cases}$$

$$\Rightarrow \sum_{i=1}^n W_i^{\mathcal{A}}(x_{1\dots n}; \mathbf{x}) = \mathbb{1}_{\sum_{j=1}^n \mathbb{1}_{x_j \in \mathcal{A}(x)} > 0} \leq 1$$

- **$k$  plus proches voisins :**

$$W_i^{k\text{-ppv}}(x_{1\dots n}; \mathbf{x}) := \frac{1}{k} \times \mathbb{1}_{\{x_i \in \{k\text{-ppv de } x \text{ parmi } x_{1\dots n}\}\}}$$

## Moyennes locales : exemples

- **Partitions :**

$$W_i^{\mathcal{A}}(x_{1\dots n}; x) := \begin{cases} \frac{\mathbb{1}_{x_i \in \mathcal{A}(x)}}{\sum_{j=1}^n \mathbb{1}_{x_j \in \mathcal{A}(x)}} & \text{si } \sum_{j=1}^n \mathbb{1}_{x_j \in \mathcal{A}(x)} > 0 \\ 0 & \text{sinon.} \end{cases}$$

$$\Rightarrow \sum_{i=1}^n W_i^{\mathcal{A}}(x_{1\dots n}; x) = \mathbb{1}_{\sum_{j=1}^n \mathbb{1}_{x_j \in \mathcal{A}(x)} > 0} \leq 1$$

- **k plus proches voisins :**

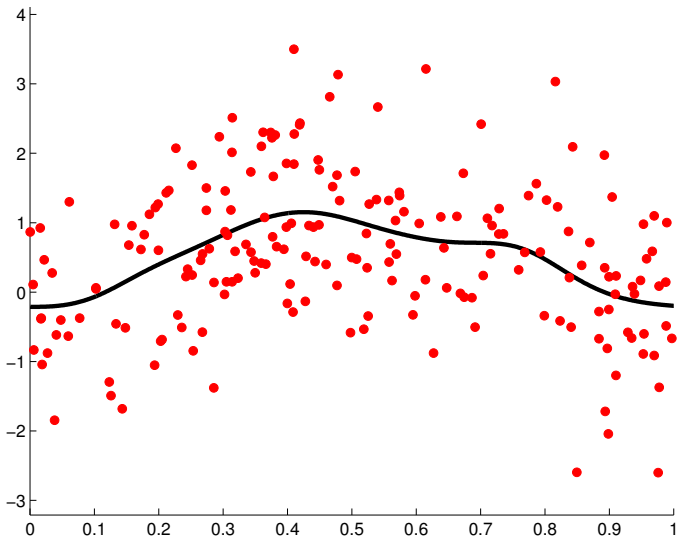
$$W_i^{k\text{-ppv}}(x_{1\dots n}; x) := \frac{1}{k} \times \mathbb{1}_{\{x_i \in \{k\text{-ppv de } x \text{ parmi } x_{1\dots n}\}\}}$$

- **Noyaux** ( $\mathcal{X} = \mathbb{R}^p$ ) :

$$W_i^{K,h}(x_{1\dots n}; x) = \frac{K\left(\frac{x_i - x}{h}\right)}{\sum_{j=1}^n K\left(\frac{x_j - x}{h}\right)}$$

Convention :  $\frac{0}{0} = 0$ .

# Exemple : noyaux (Nadaraya-Watson)



# Consistance : Théorème de Stone

## Théorème (2.2)

Hypothèses :

- $\mathcal{X} = \mathbb{R}^p$ , classification, coût 0-1

$$(a) \sum_{i=1}^n W_{i,n}(X_{1\dots n}; X) \xrightarrow[n \rightarrow \infty]{(L^1)} 1 \quad \text{et} \quad \sum_{i=1}^n W_{i,n}(X_{1\dots n}; X) \leq c_a \text{ p.s.}$$

$$(b) \text{ « pas de sous-apprentissage » : } \forall a > 0, \\ \lim_{n \rightarrow \infty} \mathbb{E} \left[ \sum_{i=1}^n W_{i,n}(X_{1\dots n}; X) \mathbb{1}_{\|X_i - X\| \geq a} \right] = 0$$

$$(c) \text{ « pas de sur-apprentissage » : } \\ \lim_{n \rightarrow \infty} \mathbb{E} \left[ \max_{1 \leq i \leq n} W_{i,n}(X_{1\dots n}; X) \right] = 0$$

$$(d) \exists c_d > 0, \forall f \in L^1(P_X), \mathbb{E} \left[ \sum_{i=1}^n W_{i,n}(X_{1\dots n}; X) f(X_i) \right] \leq c_d \mathbb{E} [f(X)]$$

Alors, la règle par moyennes locales  $\hat{f}$  associée aux poids  $(W_{i,n})_{1 \leq i \leq n}$  est **faiblement consistante pour  $P$** .

## Partitions : consistance universelle

## Corollaire (2.1)

Hypothèses :

- $\mathcal{X} = \mathbb{R}^p$ , classification, coût 0-1

(b') « pas de sous-apprentissage » :

$$\lim_{n \rightarrow \infty} \max_{A \in \mathcal{A}_n} \{\text{diam}(A)\} = 0$$

(c') « pas de sur-apprentissage » :  $\forall r > 0$ ,

$$\lim_{n \rightarrow \infty} \frac{\text{Card}\{A \in \mathcal{A}_n / A \cap \mathcal{B}(0, r) \neq \emptyset\}}{n} = 0$$

Alors,  $\hat{f}_{(\mathcal{A}_n)}^{\text{part-class}}$  est *faiblement universellement consistante*.



## Partitions : consistance universelle

## Corollaire (2.1)

Hypothèses :

- $\mathcal{X} = \mathbb{R}^p$ , classification, coût 0-1

(b') « pas de sous-apprentissage » :

$$\lim_{n \rightarrow \infty} \max_{A \in \mathcal{A}_n} \{\text{diam}(A)\} = 0$$

(c') « pas de sur-apprentissage » :  $\forall r > 0$ ,

$$\lim_{n \rightarrow \infty} \frac{\text{Card}\{A \in \mathcal{A}_n / A \cap \mathcal{B}(0, r) \neq \emptyset\}}{n} = 0$$

Alors,  $\hat{f}_{(\mathcal{A}_n)}^{\text{part-class}}$  est *faiblement universellement consistante*.

Exemple : partitions cubiques avec

(b'')  $h_n \rightarrow 0$

(c'')  $nh_n^p \rightarrow 0$

# Plan

- 1 Prévision
- 2 Régression et classification
- 3 Minimisation du risque empirique
- 4 Moyennes locales
- 5 **On n'a rien sans rien**
- 6 Conclusion : enjeux de l'apprentissage

# Consistance universelle uniforme ?

- consistance faible universelle :

$$\sup_{P \text{ loi sur } \mathcal{X} \times \mathcal{Y}} \lim_{n \rightarrow +\infty} \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] = 0$$

- consistance faible universelle **uniforme** :

$$\lim_{n \rightarrow +\infty} \sup_{P \text{ loi sur } \mathcal{X} \times \mathcal{Y}} \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] = 0$$

c'est-à-dire, vitesse d'apprentissage valable quelle que soit  $P$  ?

# Consistance universelle uniforme ?

- consistance faible universelle :

$$\sup_{P \text{ loi sur } \mathcal{X} \times \mathcal{Y}} \lim_{n \rightarrow +\infty} \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] = 0$$

- consistance faible universelle **uniforme** :

$$\lim_{n \rightarrow +\infty} \sup_{P \text{ loi sur } \mathcal{X} \times \mathcal{Y}} \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] = 0$$

c'est-à-dire, vitesse d'apprentissage valable quelle que soit  $P$  ?

- Oui **si  $\mathcal{X}$  est fini**.
- Non sinon.

# On n'a rien sans rien

## Théorème (2.3)

En classification, avec le coût 0-1, si  $\mathcal{X}$  est *infini*,  
pour toute règle de classification  $\hat{f}$  et tout  $n \geq 1$  :

$$\sup_{P \in \mathcal{M}_1(\mathcal{X} \times \{0,1\})} \left\{ \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] \right\} \geq \frac{1}{2}.$$

$\mathcal{M}_1(\mathcal{X} \times \{0,1\})$  : ensemble des mesures de probabilité sur  $\mathcal{X} \times \{0,1\}$

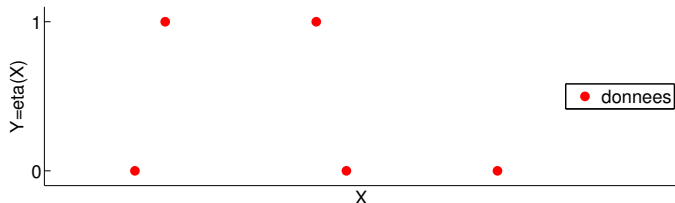
## On n'a rien sans rien

## Théorème (2.3)

En classification, avec le coût 0-1, si  $\mathcal{X}$  est *infini*, pour toute règle de classification  $\hat{f}$  et tout  $n \geq 1$  :

$$\sup_{P \in \mathcal{M}_1(\mathcal{X} \times \{0,1\})} \left\{ \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] \right\} \geq \frac{1}{2}.$$

$\mathcal{M}_1(\mathcal{X} \times \{0,1\})$  : ensemble des mesures de probabilité sur  $\mathcal{X} \times \{0,1\}$



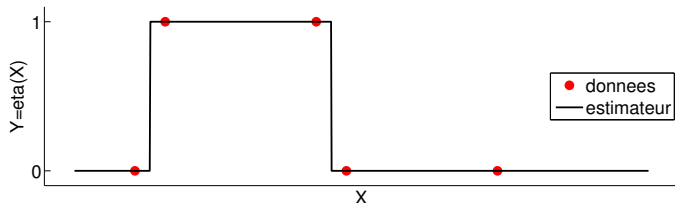
## On n'a rien sans rien

## Théorème (2.3)

En classification, avec le coût 0-1, si  $\mathcal{X}$  est *infini*, pour toute règle de classification  $\hat{f}$  et tout  $n \geq 1$  :

$$\sup_{P \in \mathcal{M}_1(\mathcal{X} \times \{0,1\})} \left\{ \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] \right\} \geq \frac{1}{2}.$$

$\mathcal{M}_1(\mathcal{X} \times \{0,1\})$  : ensemble des mesures de probabilité sur  $\mathcal{X} \times \{0,1\}$



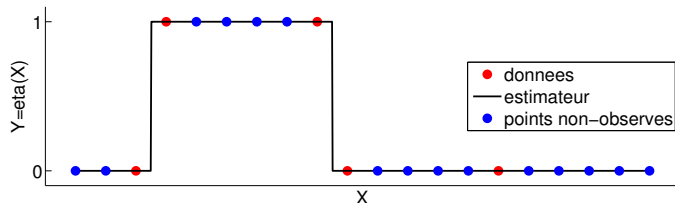
## On n'a rien sans rien

## Théorème (2.3)

En classification, avec le coût 0-1, si  $\mathcal{X}$  est *infini*, pour toute règle de classification  $\hat{f}$  et tout  $n \geq 1$  :

$$\sup_{P \in \mathcal{M}_1(\mathcal{X} \times \{0,1\})} \left\{ \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] \right\} \geq \frac{1}{2}.$$

$\mathcal{M}_1(\mathcal{X} \times \{0,1\})$  : ensemble des mesures de probabilité sur  $\mathcal{X} \times \{0,1\}$





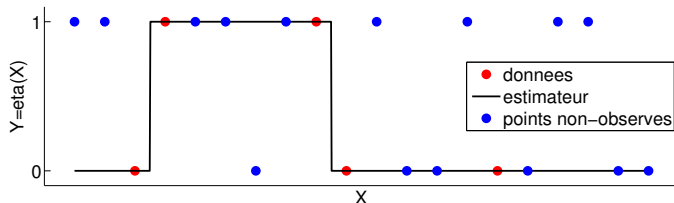
## On n'a rien sans rien

## Théorème (2.3)

En classification, avec le coût 0-1, si  $\mathcal{X}$  est *infini*, pour toute règle de classification  $\hat{f}$  et tout  $n \geq 1$  :

$$\sup_{P \in \mathcal{M}_1(\mathcal{X} \times \{0,1\})} \left\{ \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] \right\} \geq \frac{1}{2}.$$

$\mathcal{M}_1(\mathcal{X} \times \{0,1\})$  : ensemble des mesures de probabilité sur  $\mathcal{X} \times \{0,1\}$



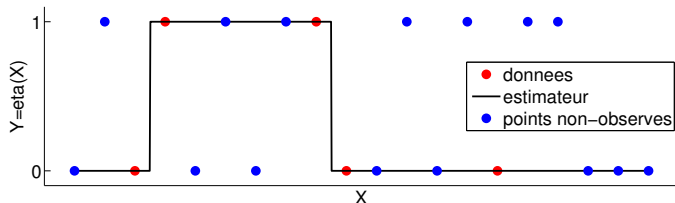
## On n'a rien sans rien

## Théorème (2.3)

En classification, avec le coût 0-1, si  $\mathcal{X}$  est *infini*,  
pour toute règle de classification  $\hat{f}$  et tout  $n \geq 1$  :

$$\sup_{P \in \mathcal{M}_1(\mathcal{X} \times \{0,1\})} \left\{ \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] \right\} \geq \frac{1}{2}.$$

$\mathcal{M}_1(\mathcal{X} \times \{0,1\})$  : ensemble des mesures de probabilité sur  $\mathcal{X} \times \{0,1\}$



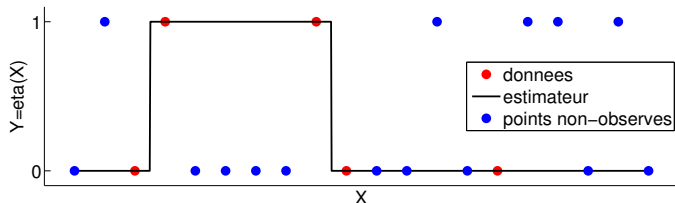
## On n'a rien sans rien

## Théorème (2.3)

En classification, avec le coût 0-1, si  $\mathcal{X}$  est *infini*,  
pour toute règle de classification  $\hat{f}$  et tout  $n \geq 1$  :

$$\sup_{P \in \mathcal{M}_1(\mathcal{X} \times \{0,1\})} \left\{ \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] \right\} \geq \frac{1}{2}.$$

$\mathcal{M}_1(\mathcal{X} \times \{0,1\})$  : ensemble des mesures de probabilité sur  $\mathcal{X} \times \{0,1\}$



Classification sur  $\mathcal{X}$  fini

Règle de majorité  $\hat{f}^{\text{maj}}$  :

$\hat{f}^{\text{maj}}((x_i, y_i)_{1 \leq i \leq n}; x)$  réalise un vote majoritaire parmi  $\{y_i / x_i = x\}$

## Proposition (2.14)

En classification, avec le coût 0–1, si  $\mathcal{X}$  est fini, pour tout  $n \geq 1$  :

$$\sup_{P \in \mathcal{M}_1(\mathcal{X} \times \{0,1\})} \left\{ \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}^{\text{maj}}(D_n)) \right] \right\} \leq \sqrt{\frac{\ln(2) \text{Card}(\mathcal{X})}{2n}}.$$

Classification sur  $\mathcal{X}$  fini

Règle de majorité  $\hat{f}^{\text{maj}}$  :

$\hat{f}^{\text{maj}}((x_i, y_i)_{1 \leq i \leq n}; x)$  réalise un vote majoritaire parmi  $\{y_i / x_i = x\}$

## Proposition (2.14)

En classification, avec le coût 0–1, si  $\mathcal{X}$  est fini, pour tout  $n \geq 1$  :

$$\sup_{P \in \mathcal{M}_1(\mathcal{X} \times \{0,1\})} \left\{ \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}^{\text{maj}}(D_n)) \right] \right\} \leq \sqrt{\frac{\ln(2) \text{Card}(\mathcal{X})}{2n}}.$$

$\Rightarrow \hat{f}^{\text{maj}}$  est uniformément universellement consistante.

# Classification sur $\mathcal{X}$ fini

Règle de majorité  $\hat{f}^{\text{maj}}$  :

$\hat{f}^{\text{maj}}((x_i, y_i)_{1 \leq i \leq n}; x)$  réalise un vote majoritaire parmi  $\{y_i / x_i = x\}$

## Proposition (2.14)

En classification, avec le coût 0–1, si  $\mathcal{X}$  est *fini*, pour tout  $n \geq 1$  :

$$\sup_{P \in \mathcal{M}_1(\mathcal{X} \times \{0,1\})} \left\{ \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}^{\text{maj}}(D_n)) \right] \right\} \leq \sqrt{\frac{\ln(2) \text{Card}(\mathcal{X})}{2n}}.$$

**Attention aux constantes** : on peut avoir  $\text{Card}(\mathcal{X}) \geq n$

$\Rightarrow$  attention aux résultats asymptotiques et aux  $\mathcal{O}(\cdot)$  qui peuvent cacher de telles constantes dans des termes de 1<sup>er</sup> ou 2<sup>e</sup> ordre.

# On n'a rien sans rien ?

Vitesse universelle à constante près :

$$\exists c(P) \in \mathbb{R}, \quad \forall n \geq 1, \quad \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] \leq c(P) u_n \xrightarrow[n \rightarrow +\infty]{} 0 ?$$

# On n'a **vraiment** rien sans rien

Vitesse universelle à constante près : **Impossible !**

$$\exists c(P) \in \mathbb{R}, \quad \forall n \geq 1, \quad \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] \leq c(P) u_n \xrightarrow{n \rightarrow +\infty} 0 ?$$

## Théorème (2.4)

*En classification, avec le coût 0-1, si  $\mathcal{X}$  est **infini**, pour toute règle de classification  $\hat{f}$  et toute suite  $(a_n)_{n \geq 1}$  décroissant vers zéro, avec  $a_1 \leq 1/16$  :*

$$\exists P \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y}), \forall n \geq 1, \quad \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \ell_P(f_P^*, \hat{f}(D_n)) \right] \geq a_n.$$

$\Rightarrow$  **impossible d'avoir  $\frac{c(P)}{\log \log n}$  comme borne de risque universelle !**



# Plan

- 1 Prévision
- 2 Régression et classification
- 3 Minimisation du risque empirique
- 4 Moyennes locales
- 5 On n'a rien sans rien
- 6 Conclusion : enjeux de l'apprentissage

# Risque minimax

## Définition (2.8)

$\mathcal{P}$  : ensemble de lois de probabilité sur  $\mathcal{X} \times \mathcal{Y}$ .

**Risque minimax** sur  $\mathcal{P}$  avec  $n$  observations :

$$\mathcal{R}_{\text{minimax}}(\mathcal{P}, n) := \inf_{\hat{f}} \sup_{P \in \mathcal{P}} \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \mathcal{R}_P(\hat{f}(D_n)) - \mathcal{R}_P^* \right]$$

# Risque minimax

## Définition (2.8)

$\mathcal{P}$  : ensemble de lois de probabilité sur  $\mathcal{X} \times \mathcal{Y}$ .

**Risque minimax** sur  $\mathcal{P}$  avec  $n$  observations :

$$\mathcal{R}_{\text{minimax}}(\mathcal{P}, n) := \inf_{\hat{f}} \sup_{P \in \mathcal{P}} \mathbb{E}_{D_n \sim P^{\otimes n}} \left[ \mathcal{R}_P(\hat{f}(D_n)) - \mathcal{R}_P^* \right]$$

**On n'a rien sans rien** (classification, coût 0-1) :

$$\forall n \geq 1, \quad \mathcal{R}_{\text{minimax}}\left(\mathcal{M}_1(\mathcal{X} \times \{0, 1\})\right) = \frac{1}{2}$$

$\Rightarrow$  « **pile ou face** » **optimal** (en pire cas) !

# Risque minimax et adaptation

$S$  classe de Vapnik-Chervonenkis,  $h \in [0, 1]$ .

- Proposition 2.15 :

$$\mathcal{R}_{\text{minimax}}\left(\{P / f^* \in S\}, n\right) \propto \sqrt{\frac{V(S)}{n}}$$

(atteint en minimisant le risque empirique sur  $S$ )

# Risque minimax et adaptation

$S$  classe de Vapnik-Chervonenkis,  $h \in [0, 1]$ .

- Proposition 2.15 :

$$\mathcal{R}_{\text{minimax}}\left(\{P / f^* \in S\}, n\right) \propto \sqrt{\frac{V(S)}{n}}$$

(atteint en minimisant le risque empirique sur  $S$ )

- Proposition 2.16 : à un facteur  $\ln(n)$  près,

$$\mathcal{R}_{\text{minimax}}\left(\left\{P / f^* \in S \text{ et } |2\eta_P(X) - 1| \geq h \text{ p.s.}\right\}, n\right) \propto \sqrt{\frac{V(S)}{n}} \wedge \frac{V(S)}{nh}$$

(atteint en minimisant le risque empirique sur  $S$ )

# Risque minimax et adaptation

$S$  classe de Vapnik-Chervonenkis,  $h \in [0, 1]$ .

- Proposition 2.15 :

$$\mathcal{R}_{\text{minimax}}\left(\{P / f^* \in S\}, n\right) \propto \sqrt{\frac{V(S)}{n}}$$

(atteint en minimisant le risque empirique sur  $S$ )

- Proposition 2.16 : à un facteur  $\ln(n)$  près,

$$\mathcal{R}_{\text{minimax}}\left(\left\{P / f^* \in S \text{ et } |2\eta_P(X) - 1| \geq h \text{ p.s.}\right\}, n\right) \propto \sqrt{\frac{V(S)}{n}} \wedge \frac{V(S)}{nh}$$

(atteint en minimisant le risque empirique sur  $S$ )

**Adaptation** : même règle (quasi) optimale  $\forall S ? \forall h ? \forall (h, S) ?$

## Conclusion : enjeux de l'apprentissage

- On ne peut pas apprendre « n'importe quoi »  
⇒ Quelles hypothèses « intéressantes » sur  $P$  ?

## Conclusion : enjeux de l'apprentissage

- On ne peut pas apprendre « n'importe quoi »  
⇒ Quelles hypothèses « intéressantes » sur  $P$  ?
- Hypothèse ⇒ vitesse optimale ? pour quelle(s) règle(s) ?






# Conclusion : enjeux de l'apprentissage

- On ne peut pas apprendre « n'importe quoi »  
⇒ Quelles **hypothèses « intéressantes »** sur  $P$  ?
- **Hypothèse** ⇒ **vitesse** optimale ? pour quelle(s) règle(s) ?
- Règle d'apprentissage ⇒ **hypothèses implicites** ?  
( $\Leftrightarrow$  identifier les lois  $P$  où elle fonctionne / ne fonctionne pas)
  - minimisation du risque empirique :  $f^*$  « proche » de  $S$
  - moyennes locales ( $k$ -ppv) :  $f^*$  « régulière »
  - réseaux de neurones (profonds), forêts aléatoires, SVM : ???

# Conclusion : enjeux de l'apprentissage

- On ne peut pas apprendre « n'importe quoi »  
⇒ Quelles **hypothèses « intéressantes »** sur  $P$  ?
- **Hypothèse** ⇒ **vitesse** optimale ? pour quelle(s) règle(s) ?
- Règle d'apprentissage ⇒ **hypothèses implicites** ?  
( $\Leftrightarrow$  identifier les lois  $P$  où elle fonctionne / ne fonctionne pas)
  - minimisation du risque empirique :  $f^*$  « proche » de  $S$
  - moyennes locales ( $k$ -ppv) :  $f^*$  « régulière »
  - réseaux de neurones (profonds), forêts aléatoires, SVM : ???
- Problèmes clés : **choix des paramètres** & **coût de calcul**  
**Données massives** ⇒ nouveaux défis

# Questions ?

-  Cornillon, P.-A. and Matzner-Løber, E. (2011).  
*Régression avec R.*  
Pratique R. Springer.
-  Devroye, L. P., Györfi, L., and Lugosi, G. (1996).  
*A Probabilistic Theory of Pattern Recognition*, volume 31 of  
*Applications of Mathematics (New York)*.  
Springer-Verlag, New York.
-  Hastie, T., Tibshirani, R., and Friedman, J. (2009).  
*The Elements of Statistical Learning.*  
Springer Series in Statistics. Springer, New York, second  
edition.  
Data Mining, Inference, and Prediction.