

# théorie de décision $\rightarrow$ apprentissage

$$D_n = \{(X_i, Y_i)\}_{i=1}^n \sim P$$

$$l: A \times (X \times Y)$$

$$f: X \rightarrow Y$$

$$R_P(f) = \mathbb{E}_{(X,Y) \sim P} [l(f(X), Y)]$$

$$A(D_n) = \hat{f}_n$$

## tasks

classification

régression

estimation de densité

pertes

0-1

quadratique

$-\log P(y)$

alg

régression linéaire

approche prob.

MV

régression logistique etc.

## • surapprentissage $\Leftarrow$

complexité

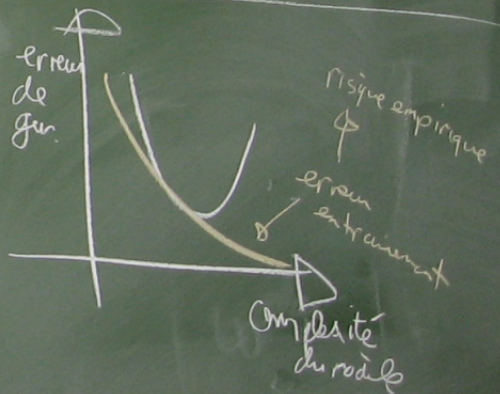
besoin de régularisation

- $\rightarrow$  réduire
- implicite
- explicite

$\rightarrow$  ML

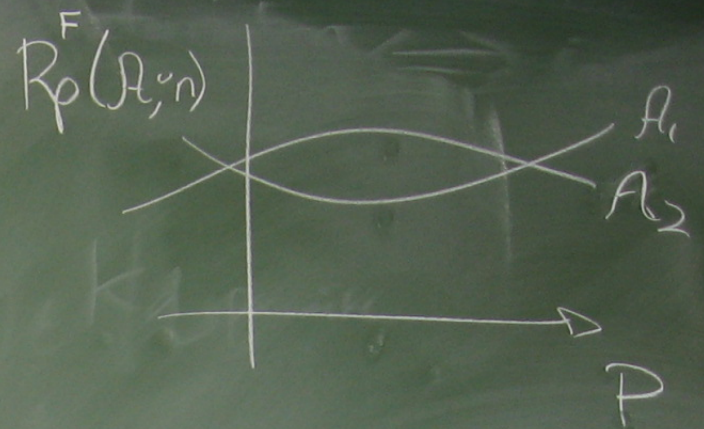
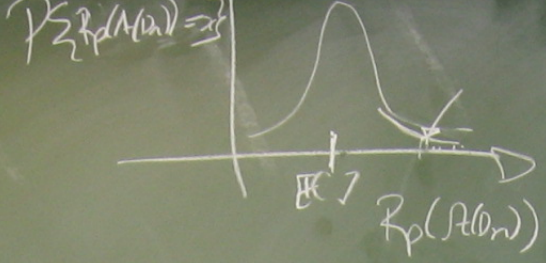
$\rightarrow$  borne d'Occam

$\rightarrow$  CL de Mallows





# THÉORIE



risque fréquentiste  $\mathbb{E} [R_p(A(D_n))]$   
 $D_n \sim P^{\otimes n}$

Sample complexity:  $n$  t.q.  $R_p^F(A, n) - R_p(f_p^*) \leq \epsilon$

Uniforme:  $\sup_{P \in \mathcal{P}} [R_p^F(A; n) - R_p(f_p^*)]$   
 espérance de l'excès de risque

PAC  $P\{R_p(A(D_n)) \leq \text{borne}(\_)\} \geq 1 - \delta$

Consistence  
 $\lim_{n \rightarrow \infty} [ \_ ] \rightarrow 0$

Sample complexity:  $n$  t.q.  $R_p(A(D_n)) - R_p(f_p^*) \leq \epsilon$   
 avec prob.  $\geq 1 - \delta$

Outils:

- loi des grands nombres
  - thm. centrale limite
  - inégalités de concentration
- asymptotique
- Chernoff  
 Hoeffding

excès de risque  
 Consistence uniforme universelle  
 → pas possible  $\nexists$  unifiée (no free lunch)  
 → possible  $\exists$  finie



$R_p^F$

→ perte quadratique

décomposition  
biais-variance

$$R_p^F(A; n) - R_p(f_p^*) = \underbrace{\mathbb{E}_X \left[ \left( \mathbb{E}_{D_n} [\hat{f}_n(X)] - f_p^*(X) \right)^2 \right]}_{\text{biais}} + \underbrace{\mathbb{E}_{D_n, X} \left[ \left( \hat{f}_n(X) - \mathbb{E}_{D_n} [\hat{f}_n(X) | X] \right)^2 \right]}_{\text{variance}}$$

- James-Stein
- $C_L$  Mallows  $\left. \begin{matrix} \text{design} \\ X \text{ fixe} \end{matrix} \right\}$  régression
- Consistance alg. partition (plug-in)  $\left. \begin{matrix} h_n \rightarrow 0 \\ nh_n^d \rightarrow +\infty \end{matrix} \right\}$  biais  $\rightarrow 0$   
variance  $\rightarrow 0$

borne deccam

$$R_p(f) \leq \hat{R}_n(f) + \frac{1}{\sqrt{2n}} \sqrt{\ln 2 |f|_{\infty} + \ln \frac{1}{\delta}} \approx \frac{1}{\sqrt{2n}} \sqrt{\ln 2 |f|_{\infty} + \ln \frac{1}{\delta}} \approx \frac{1}{\sqrt{2n}} \sqrt{\ln 2 |f|_{\infty} + \ln \frac{1}{\delta}}$$

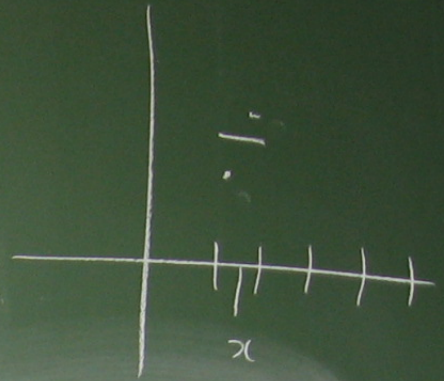
$\sqrt{\ln 2 |f|_{\infty} + \ln \frac{1}{\delta}}$  : complexité (f)  
 $\frac{1}{\sqrt{2n}}$  : à priori  
 $\approx \frac{1}{\sqrt{2n}}$  : [fléau de dimension]  
 $\forall P$  avec prob.  $\geq 1-\delta$



# Approches

• moyennage local  
(regression)

$$\hat{f}(x) = \sum_{i=1}^n W_i(x) Y_i$$



- histogramme  
(partition)

$$W_i(x) \propto \mathbb{I}\{X_i \in A(x)\}$$

$$\sum_i W_i(x) = 1$$

- K p.p.v.

$$\mathbb{I}\{X_i \in V_K(x)\}$$

↳ plug-in

$$Y_i \in \{0, 1\}$$

$$\hat{f}(x) = \mathbb{I}\{\hat{\pi}(x) \geq \frac{1}{2}\}$$

$$f^*(x) = \mathbb{E}[Y | X=x]$$

↳ pour perte quadratique

$$\mathbb{E} \left[ R(\hat{f}) \right] - R(f^*) \leq 2 \sqrt{\mathbb{E} [R(\hat{f})] - R(f^*)}$$



Min. risque empirique régularisé

$$\min_f \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) + \lambda \Omega(f)$$

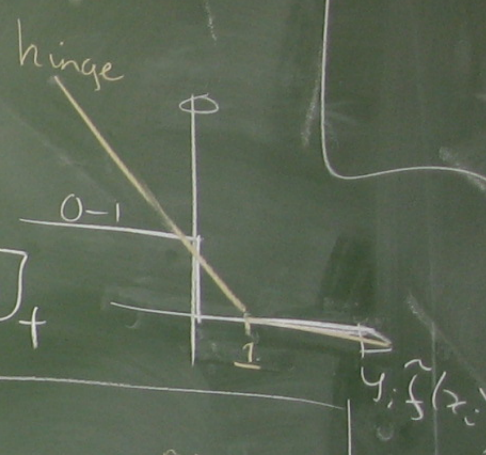
l-quadratique → régression ridge  $f: X \rightarrow \mathbb{R}$  → Problématique en classification NP-dur

plig-in  $\begin{cases} f: X \rightarrow \{-1, 1\} \\ \tilde{f}: X \rightarrow \mathbb{R}^q \end{cases}$

$\tilde{\ell}(\tilde{f}(x_i), y_i)$

exemple: pour SVM  $\tilde{\ell} \rightarrow \text{hinge} [1 - y_i \tilde{f}(x_i)]_+$

$f(x) \equiv \text{sgn}\{\tilde{f}(x) \cdot \mathbf{1}\}$



Maximum de vraisemblance → perte log

$f_w(x) \rightarrow P(y=1|x)$  → Régression logistique

$\ell(f_w(x), y) = -\log P(y|x)$

Surapprentissage

Complexité

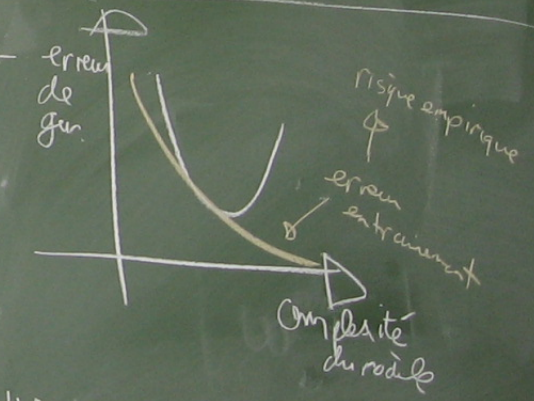
→ besoin de régularisation

→ réduire

- implicite
- explicite

→ borne d'Occam

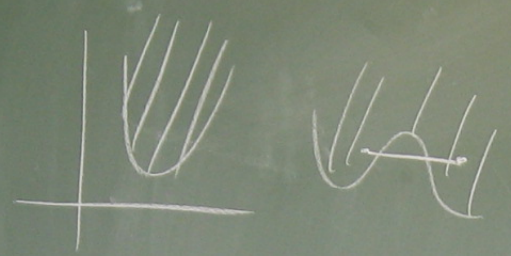
→ CL de Mallows





Optimisation & analyse Convexe  $\rightarrow$  classe de fct. "tractable"

- ensembles convexes
- fonctions convexes  $\left\{ \begin{array}{l} \text{épigraphe} \\ \text{épigraphes} \end{array} \right.$
- Convexité forte
- Jensen



Algorithmes

- descente de gradient
- ellipsoïde
- Newton  $\left( \begin{array}{l} \text{régression logistique} \\ \rightarrow \text{IRLS} \end{array} \right)$

Optimisation sous contrainte problème dual

$\min_x f(x)$   
 t.q.  $h_i(x) = 0$   
 $g_j(x) \leq 0$

$$\sup_{\lambda, \mu} \left[ \inf_x \underbrace{f(x) + \lambda^T h(x) + \mu^T g(x)}_{L(x, \lambda, \mu)} \right]$$

fct. dual  $q(\lambda, \mu)$

$$f(x) \geq f(x^*) \geq q(\lambda^*, \mu^*)$$

doublet

Slater

$\Rightarrow$  conditions KKT pour  $(x^*, \lambda^*, \mu^*)$

$x^*$  optimise  $L(x, \lambda^*, \mu^*)$   
 $x^*, \lambda^*, \mu^*$  faisible  
 $\mu_j^* g_j(x^*) = 0$  "Comp. Stricte"



\* Sélection de modèle / hyper-paramètre (A)

• validation croisée

•  $C_L$  Mallows

\* générative vs. conditionnelle  
 $p(x, y)$        $p(y|x)$

LDA, QDA

↓  
 regression logistique

\* astuce du noyau

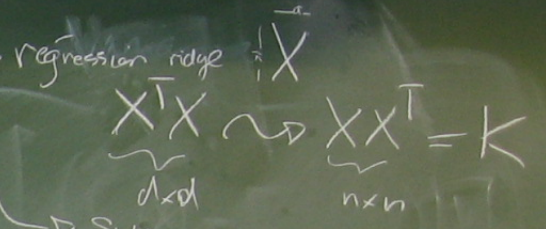
$$\langle \phi(x), \phi(x') \rangle = k(x, x')$$

- thm. du représentant
- dualité

(RBF)

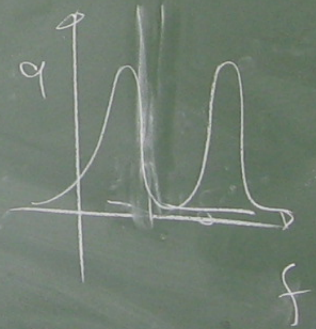
$$k(x, x') = \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right)$$

$$= \sum_{p=0}^{\infty} \frac{(-\|x-x'\|^2)^p}{2^p p!}$$



optionnelle  
 fréquentiste  $\hat{f}$

Bayésienne  $\hat{q}(f)$



$\max_f \hat{q}(f)$   
 $f \rightarrow \text{MAP}$

$$\pi(f) \quad p(D_n | f) \quad \hat{q}(f) \propto p(D_n | f) \pi(f)$$

$y = f(x)$

$$\int_{-f}^f f(x) \hat{q}(f) df \quad \text{posterior mean estimate}$$

$\max_f p(D_n | f)$



$$X X^T = K$$

$n \times n$

$$\frac{1}{2} \| \cdot \|^2$$

$$\frac{1}{2} \lambda \|\cdot\|^2$$

$$\max_f \hat{q}(f)$$

$f \rightarrow \text{MAP}$

$$\max_f P(\theta | f) \rightarrow \text{ML}$$

max. vrais.

