

# Notes de cours 3

Slide 5. On fixe la partition  $\mathcal{A} = \{A_1, \dots, A_M\}$ . on définit

$$\hat{c}_A = \operatorname{argmin}_{c \in \mathcal{F}_A} \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i c(x_i) < 0)$$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i c(x_i) < 0) &= \frac{1}{n} \sum_{m=1}^M \sum_{i: x_i \in A_m} \mathbb{1}(y_i c(x_i) < 0) \\ &= \frac{1}{n} \sum_{m=1}^M \left\{ \sum_{\substack{i: x_i \in A_m \\ y_i = 1}} \mathbb{1}(y_i c(x_i) < 0) + \sum_{\substack{i: x_i \in A_m \\ y_i = -1}} \mathbb{1}(y_i c(x_i) < 0) \right\} \end{aligned}$$

On note  $c_m$  la valeur que prend  $c$  sur  $A_m$ .

$$\begin{aligned} \text{Si } c_m = 1 \quad \text{alors } \left\{ \right\} &= \# \{i: x_i \in A_m, y_i = -1\} \\ c_m = -1 \quad \text{--- } \left\{ \right\} &= \# \{i: x_i \in A_m, y_i = 1\} \end{aligned}$$

On veut minimiser le nombre d'erreurs. on doit poser

$$\begin{aligned} c_m = 1 \quad \text{si } \# \{i: x_i \in A_m, y_i = 1\} \\ > \# \{i: x_i \in A_m, y_i = -1\} \end{aligned}$$

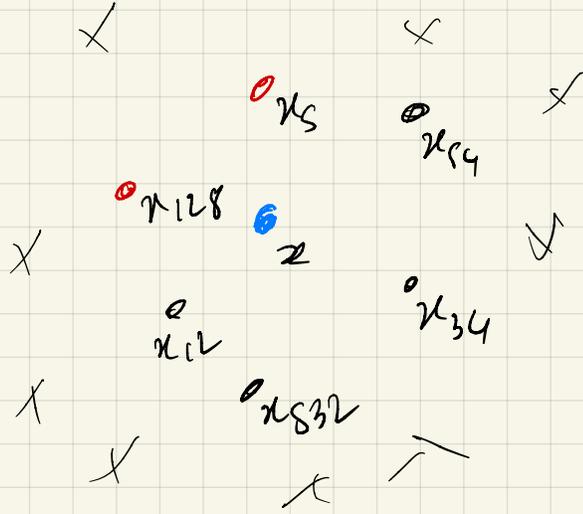
$$c_m = -1 \quad \text{sinon}$$

→ ça correspond donc à un vote à la majorité. → il reste à construire la partition!

# Slide 9

$k = 6$ .

L'algo  $k$ -nn  
colore le point  
bleu en noir.



On note  $I_x = \{ \text{des numéros des } k + \text{ proches voisins de } x \}$   
 $= \{ 5, 12, 32, 34, 54, 128 \}$   
 $= \{ i : \|x - x_i\| \text{ sont les } k = 6 + \text{ petites distances} \}$ .

On vote à la majorité:

$$\hat{c}(x) = \underset{l \in \{1, \dots, 1\}}{\operatorname{argmax}} \{ i \in I_x, y_i = l \}$$

On remarque  $I_x$  est un sous-ensemble à  $k$  éléments de  $\llbracket 1, n \rrbracket$ , il appartient à  $\{ \phi^1, \dots, \phi^M \}$  avec  $M = \binom{n}{k}$ .

On note

$$A_m = \{ x \in X, I_x = \phi^m \}$$

→ on en déduit que les  $k$ -nn sont bien dans la classe des classifieurs constants sur une partition.

Slide 29  $\Delta$   $\hat{C}_b^*$  sont dépendants  
(corrélés)

v.a.r.

$z_1, \dots, z_B$  de  $\hat{m}$  les tq  $\text{cov}(z_i, z_j) = \rho$

$$\text{On a } W\left(\frac{1}{B} \sum_{b=1}^B z_b\right) = \frac{1}{B^2} \sum_{b=1}^B \sum_{b'=1}^B \text{cov}(z_b, z_{b'})$$

$$= \frac{1}{B^2} \left\{ \sum_{b=1}^B W(z_b) + \sum_{b=1}^B \sum_{b' \neq b} \text{cov}(z_b, z_{b'}) \right\}$$

$$= \frac{1}{B^2} \left\{ B W(z_b) + \sum_{b=1}^B \sum_{b' \neq b} W(z_b) \rho \right\} \quad \text{ici} \quad \text{cov}(z_b, z_{b'}) = \frac{\text{cov}(z_b, z_{b'})}{\sqrt{W(z_b)W(z_{b'})}}$$

$$= \frac{1}{B^2} \left\{ B W(z_b) + B(B-1) \rho W(z_b) \right\} = \frac{\text{cov}(z_b, z_{b'})}{W(z_b)}$$

$$= \frac{1}{B} W(z_b) + \frac{(B-1)}{B} \rho W(z_b)$$

$$= \rho W(z_b) + \frac{1}{B} (1-\rho) W(z_b)$$

il peut arriver que  $W\left(\frac{1}{B} \sum z_b\right) > W(z_b)$  

Slide 39. Principe du gradient boosting.

On veut construire  $B$

$$g^{(B)}(x) = \sum_{b=1}^B \eta^{(b)} h^{(b)}(x) \quad \text{avec} \quad \eta^{(b)} > 0.$$

et on veut minimiser la perte empirique

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, g^{(B)}(x_i))$$

"brutalement" on voudrait définir

$$\hat{g}^{(B)}(x) = \sum_{b=1}^B \hat{\eta}^{(b)} \hat{h}^{(b)}(x)$$

$$\left( \underbrace{\hat{\eta}^{(1)}, \dots, \hat{\eta}^{(B)}}_{\in \mathbb{R}_+}, \underbrace{\hat{h}^{(1)}, \dots, \hat{h}^{(B)}}_{\in \mathcal{H}} \right)$$

$$= \operatorname{argmin}_{\eta^{(1)}, \dots, \eta^{(B)}, h^{(1)}, \dots, h^{(B)}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, \sum_{b=1}^B \eta^{(b)} h^{(b)}(x_i))$$

le problème est qu'il faudrait chercher pour les  $h$  dans un ensemble de cardinal  $\#(\mathcal{H})^B \rightarrow$  trop long... NP dur  
On ne fait pas, comme ça!

On procède pas-à-pas via une sorte de descente de gradient.

À l'itération  $b$ : on a  $\hat{g}^{(b)}$

Pour  $b$ : on cherche  $\hat{g}^{(b+1)}$  de la forme  $\hat{g}^{(b)} + \eta^{(b+1)} \hat{h}^{(b+1)}$

il reste à définir  $\hat{\eta}^{(b+1)}$  et  $\hat{h}^{(b+1)}$ .

on pourrait faire:

$$\hat{\eta}^{(b+1)}, \hat{h}^{(b+1)} = \operatorname{argmin}_{\eta, h} \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{g}^{(b)}(x_i) + \eta h(x_i))$$

C'est encore trop long...

On utilise la descente de gradient sur un pas seulement.

On cherche à faire un pas de descente de gradient pour

$$u \mapsto \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{g}^{(b)}(x_i) + \eta u(x_i))$$

en enlevant, pour l'instant la contrainte  $u \in \mathcal{H}$ .

On remarque que  $u: \mathbb{R}^d \rightarrow \mathbb{R}$  et qu'on s'intéresse seulement à ses valeurs  $\begin{pmatrix} u(x_1) \\ \vdots \\ u(x_n) \end{pmatrix}$  on va identifier la fonction  $u$  au vecteur  $\begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} \in \mathbb{R}^n$ .

On cherche à minimiser sur  $\underline{\mathbb{R}^n}$

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{g}^{(b)}(x_i) + \eta u_i)$$

$$\left| \hat{g}^{(b)} = \hat{g}^{(b)} + 0 \right.$$

On fait un pas de descente de gradient.

$$\nabla_u \left( \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{g}^{(b)}(x_i) + \eta u_i) \right)$$

$$= \frac{\eta}{n} \cdot \begin{pmatrix} \nabla_{y'} \ell(y_1, \hat{g}^{(b)}(x_1) + \eta u_1) \\ \vdots \\ \nabla_{y'} \ell(y_n, \hat{g}^{(b)}(x_n) + \eta u_n) \end{pmatrix} \in \mathbb{R}^n. \quad (y, y') \rightarrow \ell(y, \underline{y'})$$

Le pas de gradient à considérer est dans la

$$\text{direction } \underline{g}^{(b+1)} = \frac{\eta}{n} \begin{pmatrix} \nabla_{y'} \ell(y_1, \hat{g}^{(b)}(x_1) + \eta u_1) \\ \vdots \\ \nabla_{y'} \ell(y_n, \hat{g}^{(b)}(x_n) + \eta u_n) \end{pmatrix}_0$$

$$= \frac{\eta}{n} \begin{pmatrix} \nabla_{y'} \ell(y_1, \hat{g}^{(b)}(x_1)) \\ \vdots \\ \nabla_{y'} \ell(y_n, \hat{g}^{(b)}(x_n)) \end{pmatrix}.$$

Dans le modèle linéaire

$$l(y, y') = \frac{1}{2} (y - y')^2$$

$$\nabla_{y'} l(y, y') = -(y - y')$$

dans ce cas, le pas de gradient est donné par

$$\delta_{\text{lm}}^{(b+1)} = \frac{\eta}{n} \begin{pmatrix} -(y_1 - \hat{g}^{(b)}(x_1)) \\ \vdots \\ -(y_n - \hat{g}^{(b)}(x_n)) \end{pmatrix}$$

on remarque que ce sont des résidus de l'itération  $b$ .

le problème: en faisant comme si en vérité que le pas de gradient est dans  $\mathcal{H}$ !

$\hat{g}^{(b)} = \hat{g}^{(b)} + 0$  la valeur courante de " $\eta h$ " est  $\begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$

dans les descentes de gradient

$$\begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} - \text{gradient} = -\delta^{(b+1)} = - \begin{pmatrix} -\delta^{(b+1)}(1) \\ \vdots \\ -\delta^{(b+1)}(n) \end{pmatrix}$$

On voudrait pour " $\eta h^{(b+1)}$ " proportionnel à  $-\delta^{(b+1)}$

si on fait ça, on n'a aucune assurance que  $\hat{h}^{(b+1)}$  est dans  $\mathcal{H}$ .

On cherche  $\hat{h}^{(b+1)}$  comme

$$\left( \hat{h}^{(b+1)}, \hat{\nu} \right) = \underset{h \in \mathcal{H}, \nu \in \mathbb{R}}{\text{argmin}} \frac{\eta}{n} \sum_{i=1}^n \left( \underbrace{\delta^{(b+1)}(i)}_{\text{descente idéale}} - \nu \underbrace{h(x_i)}_{\text{descente contrainte à rester dans } \mathcal{H}} \right)^2$$

Slide 47. En régression linéaire  $\delta^{(b+1)}(i) = -(y_i - \hat{g}^{(b)}(x_i))$

on veut résoudre  $\underset{h \in \mathcal{H}, \nu \in \mathbb{R}}{\text{argmin}} \frac{\eta}{n} \sum_{i=1}^n \left( -(y_i - \hat{g}^{(b)}(x_i)) - \nu h(x_i) \right)^2$

Cela revient à chercher la fonction  $h \in \mathcal{H}$  qui "fitte" le mieux les résidus de l'étape  $b$ .

Ada boost. Slide 52.

$$l(y, u) = \exp(-yu)$$

À l'itération  $b$ .  $\hat{g}^{(b)}$

À l'itération  $b+1$  on cherche  $\eta$  et  $h$  tels

$$\begin{aligned} \hat{\eta}^{(b+1)}, \hat{h}^{(b+1)} &= \underset{\eta, h}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \exp(-y_i (\hat{g}^{(b)}(x_i) + \eta h(x_i))) \\ &= \frac{1}{n} \sum_{i=1}^n \exp(-y_i \hat{g}^{(b)}(x_i)) \exp(-y_i \eta h(x_i)) \\ (*) &= \frac{1}{n} \sum_{i=1}^n w^{(b)}(i) \exp(-y_i \eta h(x_i)) \end{aligned}$$

risque empirique associé à la perte exponentielle pondéré par les  $w^{(b)}(i)$

$w^{(b)}(i) \propto \exp(-y_i \hat{g}^{(b)}(x_i))$   
est grand si on a mal classé à l'étape  $b$   
est petit si on a bien classé.

On a soit  $y_i h(x_i) = 1$   
soit  $y_i h(x_i) = -1$

On met la contrainte  $\frac{1}{n} \sum w^{(b)}(i) = 1$

$$\begin{aligned} (*) &= \frac{1}{n} \left\{ \sum_{i: y_i h(x_i) = 1} w^{(b)}(i) \exp(-\eta) + \sum_{i: y_i h(x_i) = -1} w^{(b)}(i) \exp(\eta) \right\} \\ &= \frac{1}{n} \left\{ \sum_{i=1}^n w^{(b)}(i) \exp(-\eta) - \sum_{i: y_i h(x_i) = -1} w^{(b)}(i) \exp(-\eta) + \sum_{i: y_i h(x_i) = -1} w^{(b)}(i) \exp(\eta) \right\} \\ &= \exp(-\eta) + (\exp(\eta) - \exp(-\eta)) \frac{1}{n} \sum_{i: y_i h(x_i) = -1} w^{(b)}(i) \end{aligned}$$

Pour la minimisation en  $h$  on va chercher  $\hat{h}^{(b+1)} = \underset{h}{\operatorname{argmin}} \frac{1}{n} \sum_{i: y_i h(x_i) = -1} w^{(b)}(i) = \underset{h}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n w^{(b)}(i) \mathbb{1}(y_i h(x_i) < 0)$

Donc  $\hat{h}^{(b+1)}$  est le meilleur classifieur qui minimise la perte ou la perte.

Si on optimise en  $\eta$ , on trouve :

$$\hat{\eta}^{(b+1)} = \frac{1}{2} \log \left( \frac{1 - E^{(b+1)}}{E^{(b+1)}} \right)$$

$$\text{où } E^{(b+1)} = \sum_{i=1}^n w^{(b)}(i) \frac{1}{y_i} \hat{h}^{(b+1)}(x_i) < 0.$$