

Apprentissage supervisé $\{x^{(i)} \in \mathbb{R}^D, t^{(i)}\}_{i=1}^N$

→ Objectif : apprendre un modèle h_β

→ 1^{re} solution $\underset{\text{OLS}}{l(\beta)} = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - h_\beta(x^{(i)}))^2$

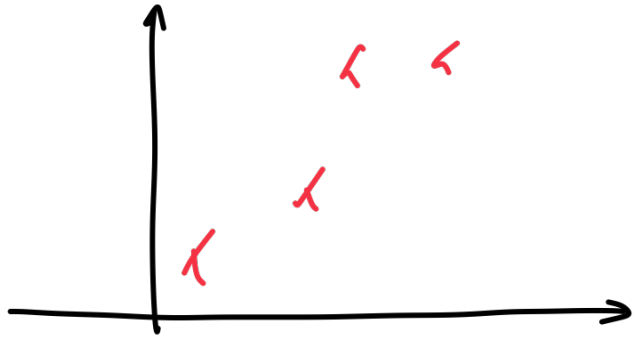
$$h_\beta = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_D x_D$$

↳ 2 approches : descente de gradient

→ Equations normales : $\beta = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T t$

→ Dans le cas de caractéristiques corrélées / matrice \tilde{X} singulière → on peut toujours coder un processus d'orthogonalisation

$\beta_D = \frac{\langle t_E, z_D \rangle}{\langle z_D, z_D \rangle}$ → Dans le cas de corrélation → amplitude de l'écart de β_j



→ Pour déterminer la complexité optimale du modèle

↳ On peut étudier l'erreur quadratique moyenne

$$MSE(x) = \text{bias}^2 + \text{variance}$$

Sélection automatique de modèles / caractéristiques

→ Best subset selection

→ Regression Ridge

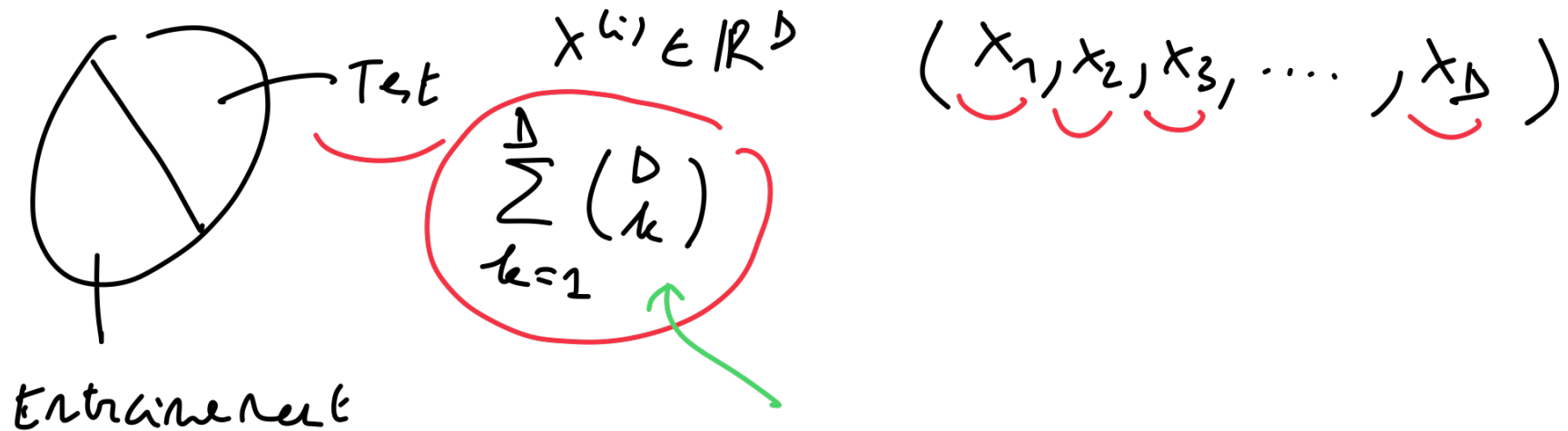
→ Regression LASSO

→ Classification → 2 classes via Perceptron

→ > 2 classes

→ regression logistique

Question: Comment sélectionner les caractéristiques les plus représentatives



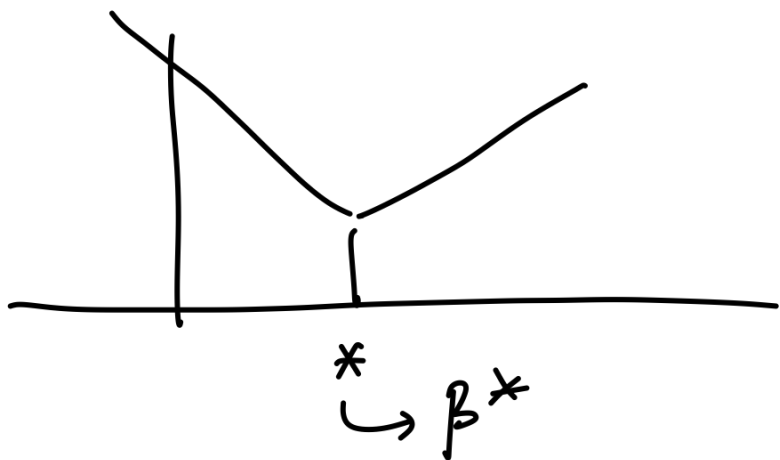
→ Si peu de données → validation croisée (cross validation)

K-fold cross validation



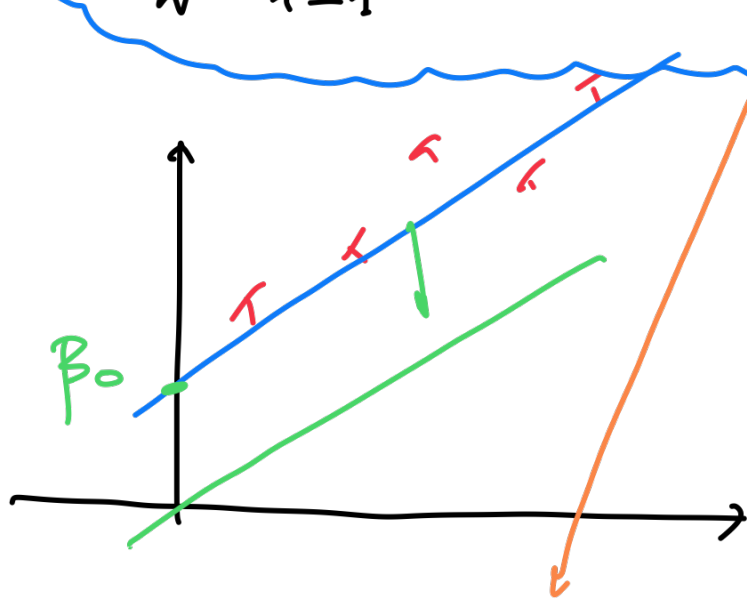
K

$$CV(\beta) = \frac{1}{N} \sum_{i=2}^N \left(t^{(i)} - h_{\beta}^{-K(i)}(x^{(i)}) \right)^2$$



$$l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - h_{\beta}(x^{(i)}))^2 + \lambda \sum_{j=2}^D |\beta_j|^2$$

$$= \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2 + \lambda \sum_{j=1}^D |\beta_j|^2$$



Fidélité
au données

RIDGE

penalty
complexité

2 possibilités

→ descente du gradient

→ Résolution des équations
normales

Equations normales

étape 1 calculer les vecteurs $x^{(i)} \in \mathbb{R}^D$

$$x^{(i)} \leftarrow x^{(i)} - \frac{1}{N} \sum_{i=1}^N x^{(i)}$$

$$\frac{1}{N} \sum x_j^{(i)} = 0$$

$$\frac{1}{N} \sum_{i=1}^N \vec{x}^{(i)} = 0$$

$$l_{\text{RIDGE}}(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2 + \lambda \sum_{j=1}^D |\beta_j|^2$$

$$\frac{\partial l}{\partial \beta_0} = \frac{2}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})) (-1)$$

$$= \frac{2}{N} \sum_{i=1}^N t^{(i)} + \frac{2}{N} \sum_{i=1}^N \beta_0 + \frac{2}{N} \sum_{i=1}^N (\cancel{\beta_1 x_1^{(i)}} + \dots + \cancel{\beta_D x_D^{(i)}})$$

$$2\beta_0 = \frac{2}{N} \sum t^{(i)} \rightarrow \beta_0 = \frac{1}{N} \sum_{i=1}^N t^{(i)}$$

Etape 2: calculer le $t_c^{(i)}$ $t_c^{(i)} = t^{(i)} - \frac{1}{N} \sum_{i=1}^N t^{(i)}$

$$\mathcal{L}(\beta) = \frac{1}{N} \sum_{i=1}^N (t_c^{(i)} - (\beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2 + \lambda \sum_{j=1}^D \beta_j^2$$

$$\frac{\partial \mathcal{L}}{\partial \beta_j} = \frac{2}{N} \sum_{i=1}^N (t_c^{(i)} - (\beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})) (-x_j^{(i)}) + 2\lambda \beta_j$$

$$= \frac{2}{N} \sum_{i=1}^N \underbrace{e_i}_{\textcircled{e_i}} \cdot \underbrace{x_j^{(i)}}_j + 2\lambda \beta_j \leftarrow$$

$$X = \begin{bmatrix} - & x^{(1)} & - \\ - & x^{(2)} & - \\ & \vdots & \\ - & x^{(N)} & - \end{bmatrix}$$

$$\text{grad}_{\beta} = \left[\frac{\partial \mathcal{L}}{\partial \beta_1}, \frac{\partial \mathcal{L}}{\partial \beta_2}, \dots, \frac{\partial \mathcal{L}}{\partial \beta_D} \right] =$$

$$= -\frac{2}{N} \sum_{i=1}^N \underbrace{e_i}_{\textcircled{e_i}} \cdot \underbrace{\bar{x}^{(i)}}_{\textcircled{\bar{x}^{(i)}}} + 2\lambda \beta$$

$$e_i = \bar{t}_i - \underline{(X \vec{\beta})}_i$$

$$= \left(-\frac{2}{N} (t - X\beta)^T X \right)^T + 2\lambda \beta \leftarrow$$

$$t = \begin{bmatrix} t^{(2)} \\ \vdots \\ t^{(N)} \end{bmatrix}$$

$$= -\frac{2}{N} X^T (t - X\beta) + 2\lambda \beta \leftarrow$$

$$= -\frac{2}{N} X^T t + 2 \underbrace{X^T X \beta} + \underbrace{2\lambda \beta}$$

$$\Rightarrow \left(2 \underbrace{X^T X} + 2\lambda I \right) \beta = \frac{2}{N} X^T t$$

$$\beta_{\text{RIDGE}} = \left(\frac{2}{N} X^T X + 2\lambda I \right)^{-1} \frac{2}{N} X^T t$$

$$\underbrace{\sum e_i \vec{x}^{(i)}}_{\sum_{i=1}^N e_i \left[\begin{array}{c} x^{(2)} \\ \vdots \\ x^{(N)} \end{array} \right]} + 2\lambda \beta$$

$$X^T X v = \alpha v$$

$$\underbrace{(X^T X + \lambda I)}_{\text{}} v = \alpha' v$$

$$\begin{aligned} X^T X v + \lambda v &= \alpha v + \lambda v \\ &= (\alpha + \lambda) v \end{aligned}$$

$$\mathcal{L}_{\text{LASSO}}(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2 + \lambda \sum_{j=2}^D |\beta_j|$$

formulations constraints

$$\min_{\beta} l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2$$

s.t

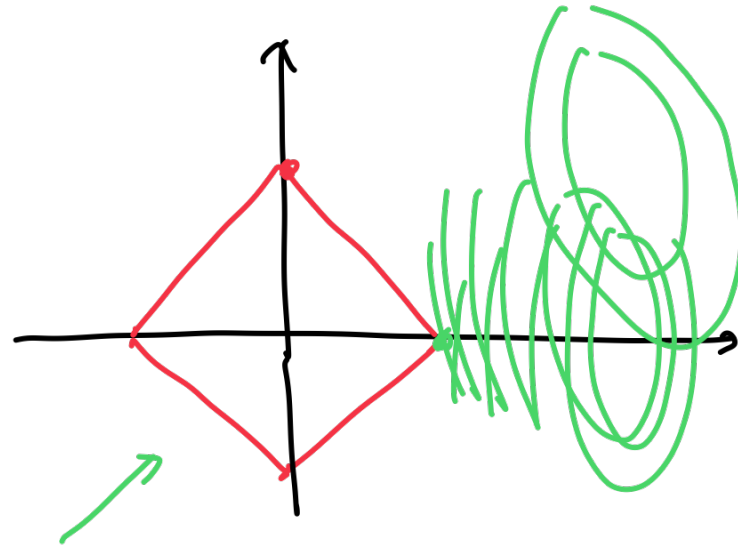
$$\left(\sum_{j=1}^D \beta_j^2 \leq t \right)$$

(Ridge)



$$\sum_{j=1}^D |\beta_j| \leq t$$

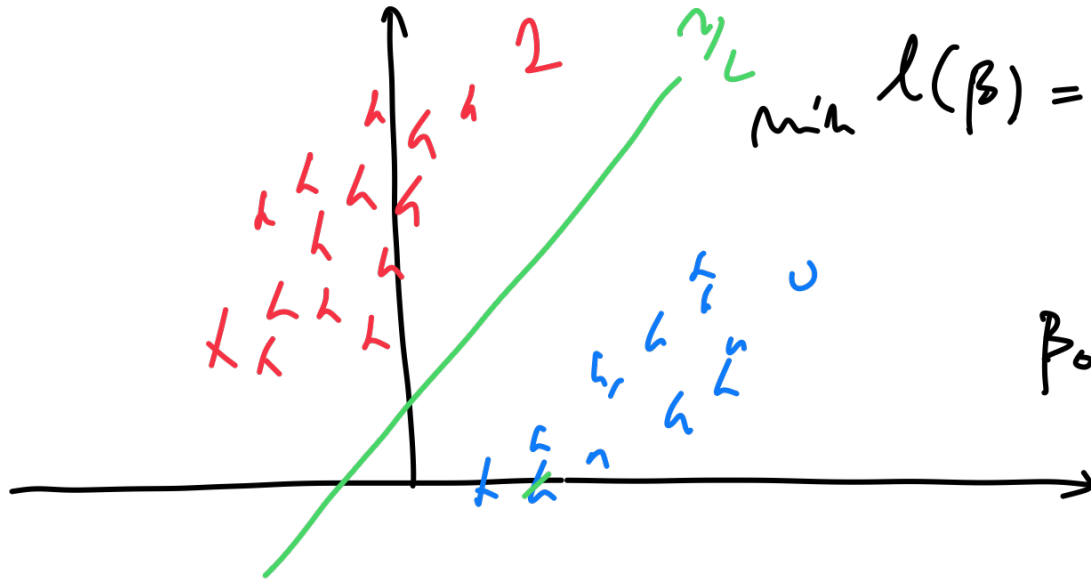
(Lasso)



$$\|v\|_2 = \sum_{i=1}^N |v_i|$$

Classification

$\{x^{(i)}, t^{(i)}\} \quad t^{(i)} \in \{1, \dots, K\}$



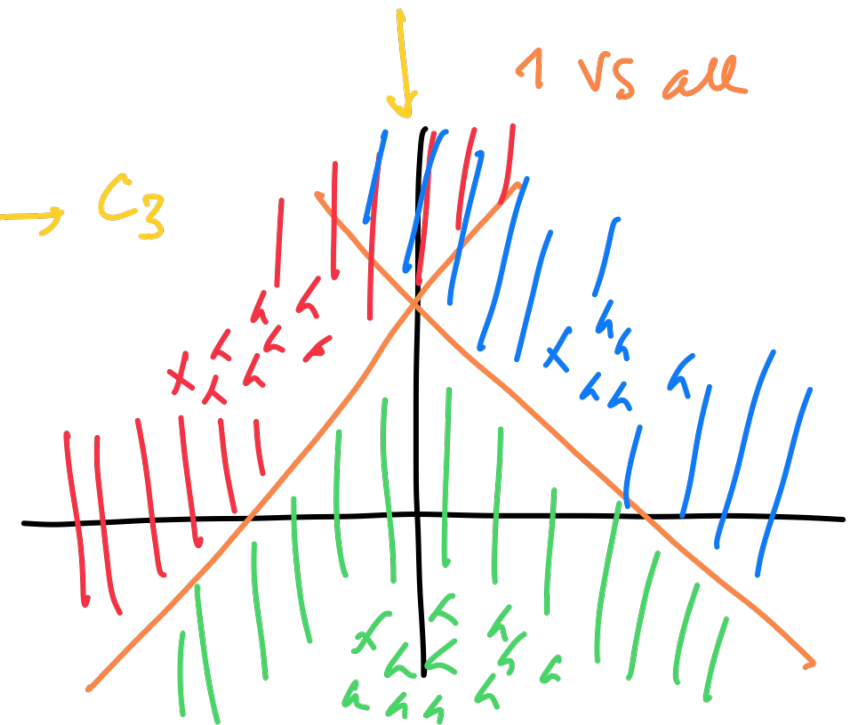
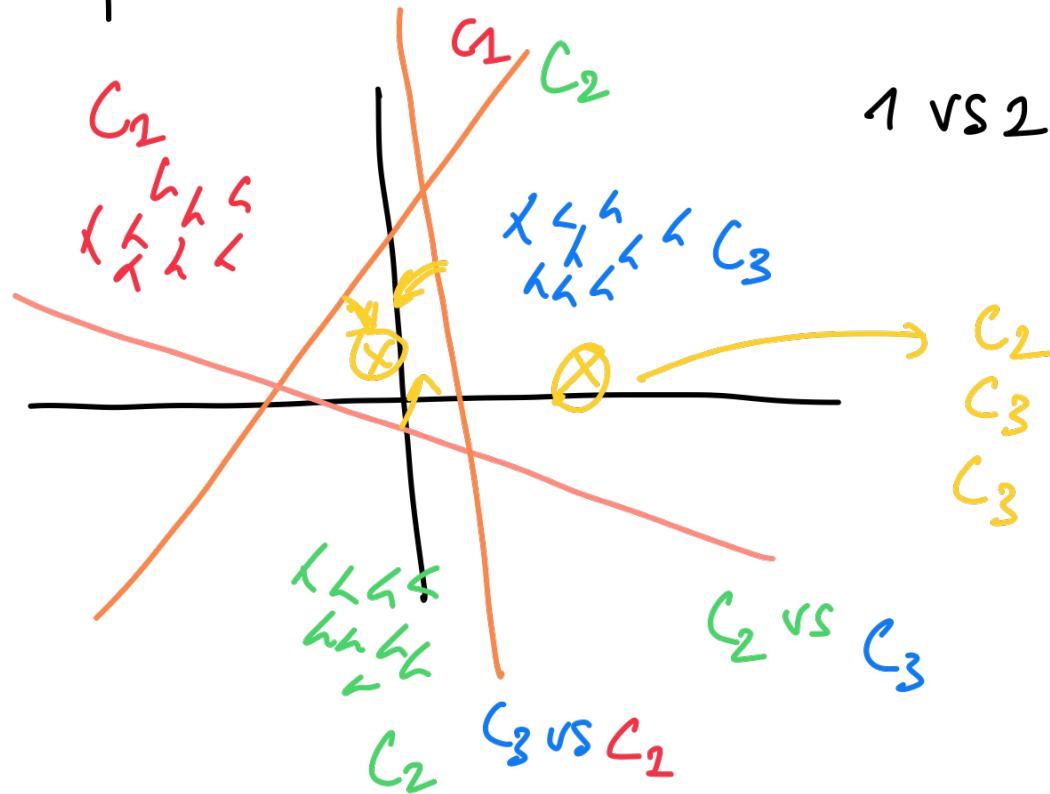
$$\min_{\beta} l(\beta) = \frac{1}{N} \sum_{i=1}^N \left(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}) \right)^2$$

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$$

$$x_2 = \frac{-\beta_1 x_1 - \beta_0}{\beta_2}$$

Q1? 2 classes \rightarrow class multiples

Q2? points outliers



$$T = \begin{bmatrix} \text{---} t^{(1)} \text{---} \\ \vdots \\ \text{---} t^{(N)} \text{---} \end{bmatrix}$$

$$t^{(i)} \in \{1, \dots, k\}$$

$$t^{(i)} = [0, \textcolor{green}{1}, 0, \dots, 0]$$

$$B = \begin{bmatrix} \text{---} \vec{\beta}_1^T \text{---} \\ \vdots \\ \text{---} \vec{\beta}_k^T \text{---} \end{bmatrix} \begin{matrix} \textcolor{green}{\left(\begin{array}{c} | \\ x^{(1)} \\ | \end{array} \right)} & \textcolor{orange}{\left(\begin{array}{c} | \\ x^{(2)} \\ | \end{array} \right)} & \dots & x^{(N)} \end{matrix} = \begin{bmatrix} \textcolor{green}{1} & 0 \\ 0 & \textcolor{orange}{1} \\ \vdots & | \\ 0 & 0 \end{bmatrix}$$

$$T \simeq \begin{matrix} X B^T \\ B X^T \end{matrix}$$

$$\min_B \frac{1}{NK} \sum_{i=1}^N \sum_{k=1}^k (t_k^{(i)} - (\vec{\beta}_k)^T x^{(i)})^2$$

$$B = (X^T X)^{-2} X^T T$$