

# Apprentissage CM3

→ Apprentissage supervisé

↳ Modèle de régression linéaire

$$h_{\beta}(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_D x_D \quad x \in \mathbb{R}^D$$
$$\beta \in \mathbb{R}^{D+1}$$

→ Fonction coût

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

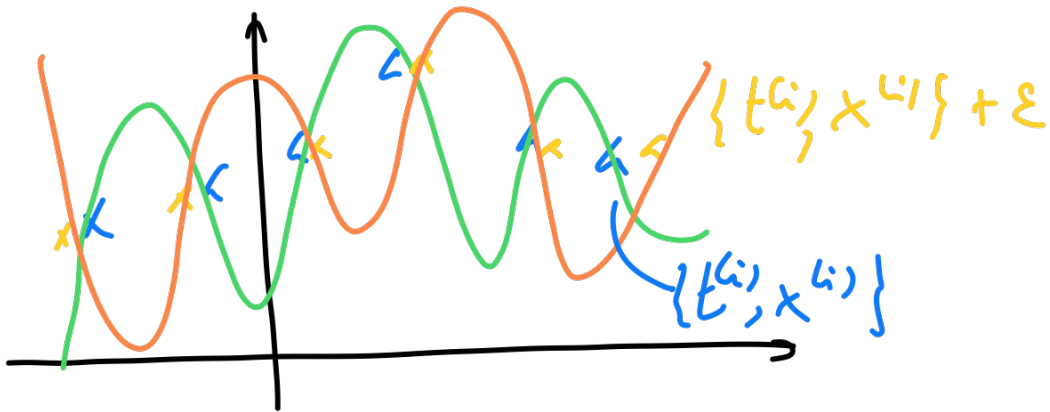
→ 2 approches pour trouver le  $\bar{\beta}$  optimal

→ Descente de gradient

→ Calcul et annulation des dérivées  
+ résolution des équations normales

$$\rightarrow \underline{\underline{\vec{\beta}}} = (\underline{\underline{X^T X}})^{-1} \underline{\underline{X^T E}}$$

?  $X^T X$  n'est pas inversible ou mal conditionné



$$\beta_{\text{RL}} = \frac{\langle t, t \rangle}{\langle x, x \rangle} \gg \gg$$

→ Régularisation

#1 Sélection du meilleur sous ensemble (Best Subset Selection)

→ for  $k = 2 : D+1$

Évaluer le modèle entraîné sur un sous ensemble de taille  $k$  des  $D+1$   $\beta_j$

→ sélectionner le sous ensemble qui donne la meilleure prédiction sur l'ensemble test  
(via cross validation)

$$\sum_{k=2}^{D+1} \binom{D+1}{k}$$

#2 Ajouter une penalité sur les coefficients  $\beta_j$  à la fonction coût  $l(\beta)$  → SOMME DES CARRÉS DES  $\beta_j$

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

→ Regression Ridge

→ fidélité aux données

penalité sur la complexité du Modèle

$$\vec{\varepsilon} = (\vec{t} - \tilde{X} \vec{\beta}) \quad l(\beta) = \frac{1}{N} \sum_{i=1}^N (\varepsilon^{(i)})^2 = \frac{1}{N} \|\vec{\varepsilon}\|^2 = \frac{1}{N} \langle \vec{\varepsilon}, \vec{\varepsilon} \rangle = \frac{1}{N} \vec{\varepsilon}^T \vec{\varepsilon}$$

$$l(\beta) = \frac{1}{N} (\vec{t} - \tilde{X} \vec{\beta})^T (\vec{t} - \tilde{X} \vec{\beta})$$

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

$$= \frac{2}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})) \cdot (-1) = 0$$

$$\rightarrow \sum_{i=1}^N t^{(i)} - \sum_{i=1}^N \beta_0 - \sum_{i=1}^N (\beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}) = 0$$

$$0 = \sum_{i=1}^N t^{(i)} - N \beta_0 - \sum_{i=1}^N \beta^T \vec{x}^{(i)} \quad \text{ou } \beta = (\beta_1, \dots, \beta_D)$$

*à partir de  $\beta_1$*

$$\beta_0 = \frac{1}{N} \sum_{i=1}^N t^{(i)} - \frac{1}{N} \sum_{i=1}^N \beta^T \vec{x}^{(i)}$$

$$\tilde{X} = [1, x_1, \dots, x_D]$$

$$* \tilde{\beta} = [\beta_0, \vec{\beta}]$$

Si on centre les données

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

$$t_c^{(i)} \leftarrow t^{(i)} - \frac{1}{N} \sum_{j=1}^N t^{(j)} = t^{(i)} - \bar{t}$$

Pour  $x_c^{(i)}$  et  $t_c^{(i)}$  si on calcule le  $\beta_0$  on retrouve

$$\begin{aligned} \beta_0 &= \frac{1}{N} \sum_{i=1}^N t_c^{(i)} - \frac{1}{N} \sum_{i=1}^N \beta^T x_c^{(i)} \quad \leftarrow ** = * (\text{données centrées}) \\ &= \frac{1}{N} \sum_{i=1}^N \left( t^{(i)} - \frac{1}{N} \sum_{j=1}^N t^{(j)} \right) - \frac{1}{N} \sum_{i=1}^N \beta^T \left( x^{(i)} - \frac{1}{N} \sum_{j=1}^N x^{(j)} \right) \\ &= \frac{1}{N} \sum_{i=1}^N t^{(i)} - \frac{1}{N} \sum_{i=1}^N \cdot \frac{1}{N} \sum_{j=1}^N t^{(j)} = \frac{1}{N} \sum_{i=1}^N t^{(i)} - \frac{1}{N} N \cdot \bar{t} \end{aligned}$$

$$= \frac{1}{N} \sum_{i=1}^N t^{(i)} - \bar{t}$$

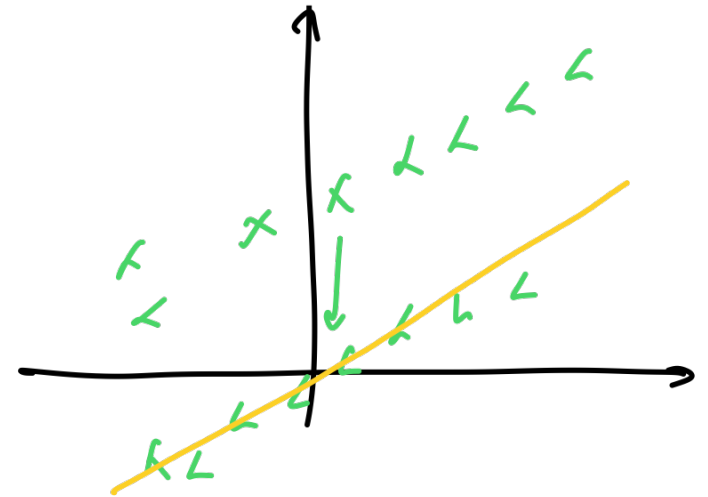
$$= \bar{t} - \bar{t} = 0$$

pour le second terme, on a de même équivalence

$$-\frac{1}{N} \sum_{i=1}^N \beta^T x^{(i)} + \frac{1}{N} \sum_{i=1}^N \beta^T \left( \frac{1}{N} \sum_{j=1}^N x^{(j)} \right)$$

$$-\frac{1}{N} \beta^T \sum_{i=1}^N x^{(i)} + \frac{1}{N} N \cdot \beta^T \bar{x}$$

$$-\beta^T \bar{x} + \beta^T \bar{x} = 0$$



Pour des données centrées on se ramène à une fonction de coût

du type 
$$l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_1 x_1^{(i)} + \beta_2 x_2^{(i)} + \dots + \beta_D x_D^{(i)}))^2 + \lambda \underbrace{\sum_{j=1}^D \beta_j^2}$$

à partir de maintenant on utilisera  $\hat{\beta}$  pour représenter le vecteur  $[\beta_0, \beta_1, \dots, \beta_D]$  et  $\vec{\beta} = [\beta_1, \beta_2, \dots, \beta_D]$

→ Pour la formulation de type Ridge, on a

$$l(\beta) = \frac{1}{N} \vec{\epsilon}^T \vec{\epsilon} = \frac{1}{N} (\vec{t} - \underline{X} \vec{\beta})^T (\vec{t} - \underline{X} \vec{\beta}) + \lambda \vec{\beta}^T \vec{\beta}$$



$$l(\beta) = \frac{1}{N} \left( \vec{E}^T \vec{E} - \vec{\beta}^T \underline{\underline{X}}^T \vec{E} - \vec{E}^T \underline{\underline{X}} \beta + \vec{\beta}^T \underline{\underline{X}}^T \underline{\underline{X}} \vec{\beta} \right) + \lambda \beta^T \beta \leftarrow$$

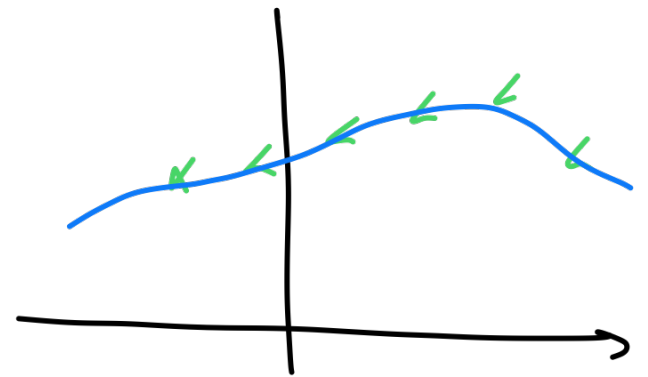
$$\frac{\partial}{\partial \beta} (v^T \beta) = v$$

$$\frac{\partial}{\partial \beta} = -2 \underline{\underline{X}}^T \vec{E} + 2 \underline{\underline{X}}^T \underline{\underline{X}} \beta \quad 2\lambda \beta$$

$$\frac{\partial}{\partial \beta} l(\beta) = - \underline{\underline{X}}^T \vec{E} + \underline{\underline{X}}^T \underline{\underline{X}} \beta + \lambda \beta = 0$$

$$\underline{\underline{X}}^T \vec{E} = \left( \underline{\underline{X}}^T \underline{\underline{X}} + \lambda \mathbf{I}_D \right) \vec{\beta}$$

$$\vec{\beta} = \left( \underline{\underline{X}}^T \underline{\underline{X}} + \lambda \mathbf{I}_D \right)^{-1} \underline{\underline{X}}^T \vec{E}$$



$$A \vec{x} = \gamma \vec{x}$$

$x, \gamma$  vecteur  
et valeur  
propre de  $A$

$$(A + \lambda I) \vec{x} = A \vec{x} + \lambda \vec{x}$$

$$= \gamma \vec{x} + \lambda \vec{x}$$

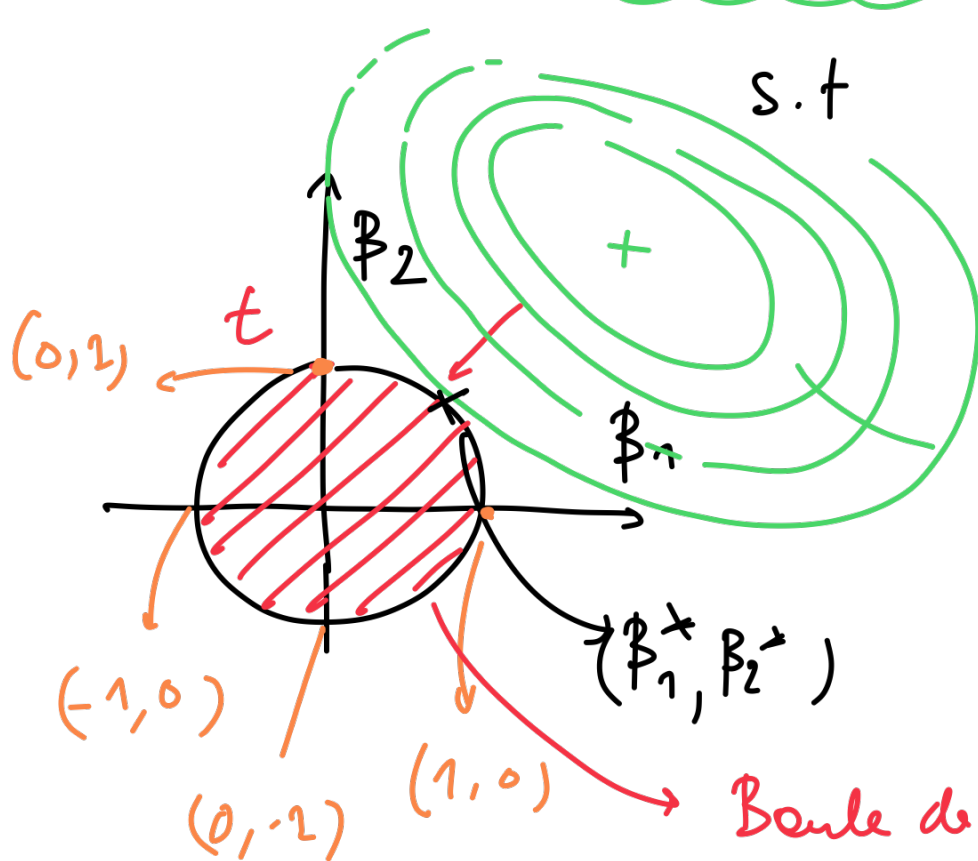
$$= (\gamma + \lambda) \vec{x}$$

→ la plus petite valeur propre de

$$\left( \underline{\underline{X^T X}} + \lambda I_b \right) \text{ est à présent défini par } \lambda + \lambda_{\min}(\underline{\underline{X^T X}})$$

La formulation Ridge peut aussi s'écrire de façon  
contrainte

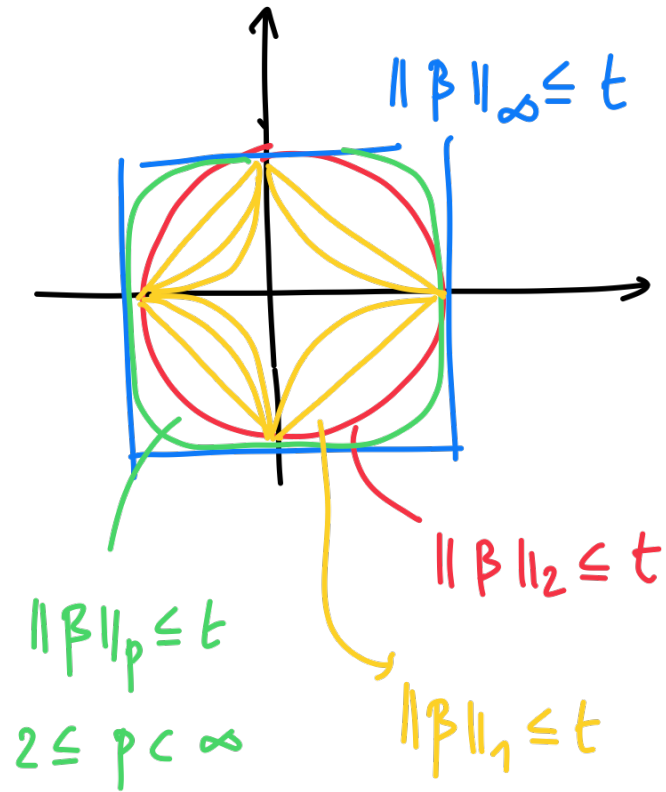
$$\min_{\beta} l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_1 x_1^{(i)} + \beta_2 x_2^{(i)} + \dots + \beta_D x_D^{(i)}))^2$$



$$\text{s.t.} \quad \sum_{j=1}^D \beta_j^2 \leq t \quad \leftarrow$$

$$\|\vec{\beta}\|_2^2 = \left( \sqrt{\beta_1^2 + \beta_2^2 + \dots + \beta_D^2} \right)^2$$

Boule de la norme  $l_2$



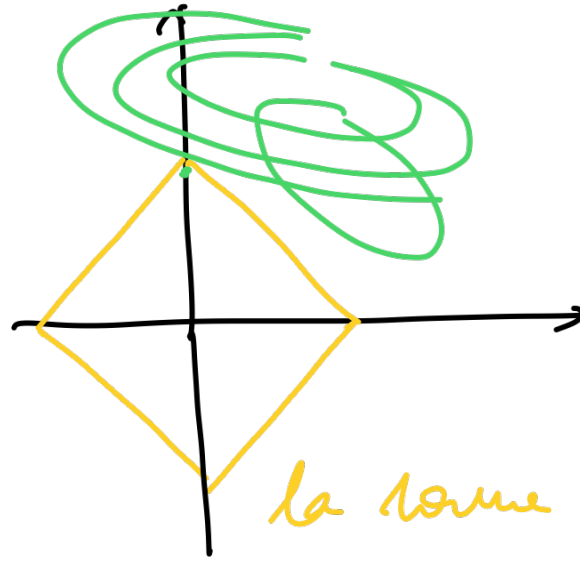
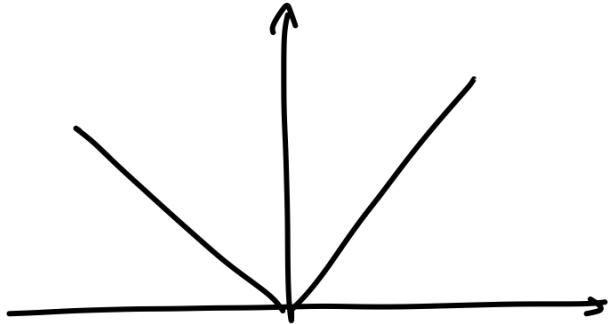
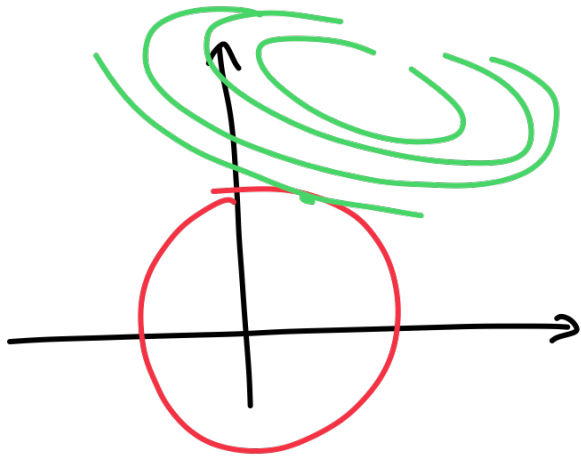
$$\|\vec{\beta}\|_p = \left( \sum_{j=1}^D |\beta_j|^p \right)^{1/p}$$

$$\|\vec{\beta}\|_\infty = \max_j |\beta_j|$$

$$\|\vec{\beta}\|_\infty \leq t$$

$p=2$  #3 Modèle de régression LASSO

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



la norme  $\|\beta\|_1$   
est plus efficace en  
terme de sélection de  
Caractéristiques

