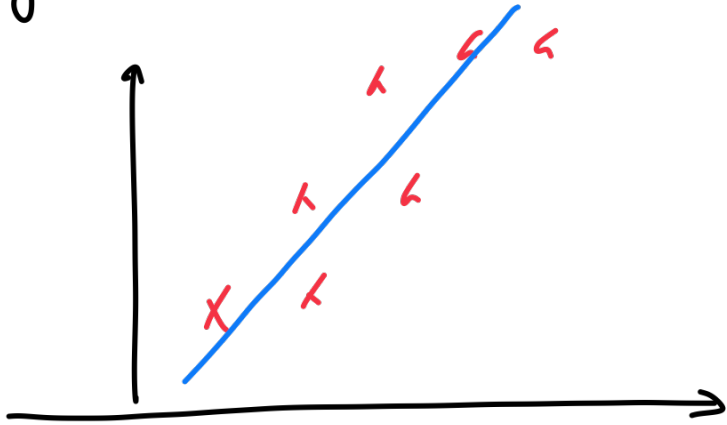


Introduction à l'apprentissage

→ Régression linéaire



$$\{x^{(i)}, t^{(i)}\} \quad x^{(i)} \in \mathbb{R}^D$$
$$t^{(i)} \in \mathbb{R}$$

→ Modèle linéaire

$$h_{\beta}(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_D x_D$$

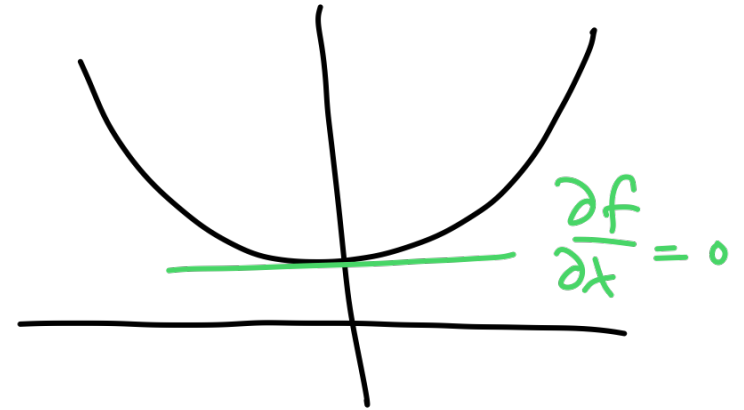
$$= \beta^T \tilde{x} = [\beta_0 \dots \beta_D]$$

$$\begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_D \end{bmatrix} \leftarrow x$$

→ Fonction de coût (loss)

→ Mesure de l'erreur (sum of squares)

$$L_{\beta}(x) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - h_{\beta}(x^{(i)}))^2$$



→ Fonction de coût Convexe

↳ 2 approches : → descente de gradient

→ Résolution des équations
normales

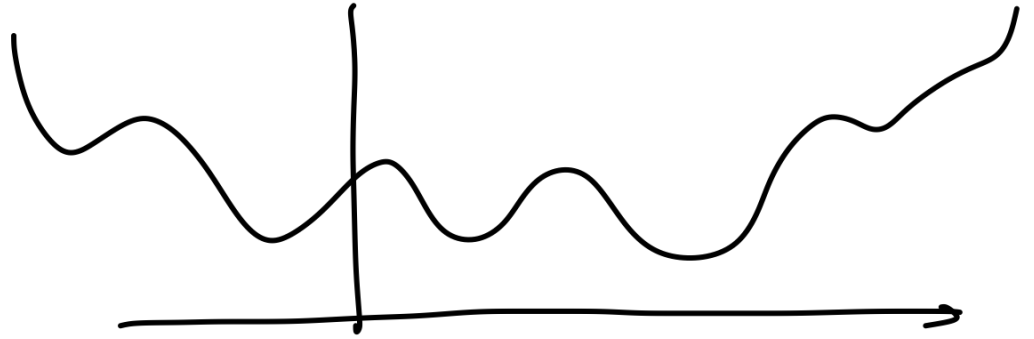
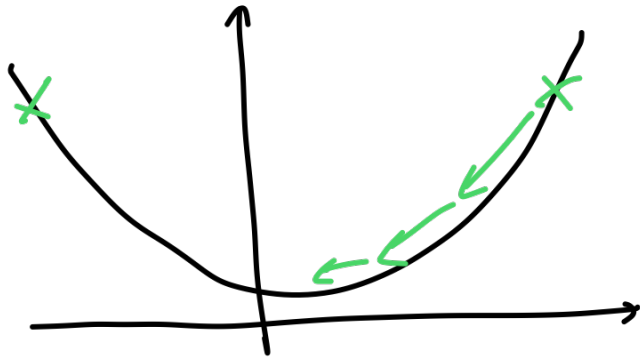
$$l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2 \leftarrow$$

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

→ Itérations de (descente de) gradient

$\beta^{(0)}$

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - \eta \cdot \text{grad}_{\beta} l(\beta)$$



→ Résolution de équations

$$t^{(i)} = \begin{bmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{bmatrix} \in \mathbb{R}^N \quad \vec{\beta} = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_D \end{bmatrix} \in \mathbb{R}^{D+1}$$

$$X = \begin{bmatrix} 1 & \vec{x}^{(1)} \\ \vdots & \vdots \\ 1 & \vec{x}^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times (D+1)}$$

$$\vec{\varepsilon} = \left(\vec{t} - \underset{=}{\tilde{X}} \underset{=}{\vec{\beta}} \right) \quad \begin{bmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{bmatrix} \approx \begin{bmatrix} 1 - \vec{x}^{(1)} - \\ \vdots \\ 1 - \vec{x}^{(N)} - \end{bmatrix} \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_D \end{bmatrix}$$

$$\rightarrow \mathcal{L}(\beta) = \frac{1}{N} \sum_{i=1}^N \varepsilon_i^2 = \|\vec{\varepsilon}\|^2 = \langle \vec{\varepsilon}, \vec{\varepsilon} \rangle = \vec{\varepsilon}^T \vec{\varepsilon}$$

Dérivée de $\mathcal{L}(\beta)$ par rapport à $\vec{\beta}$

$$\mathcal{L}(\beta) = \frac{1}{N} \left(\vec{t} - \underset{=}{\tilde{X}} \underset{=}{\vec{\beta}} \right)^T \left(\vec{t} - \underset{=}{\tilde{X}} \underset{=}{\vec{\beta}} \right)$$

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

transposé du produit $(A B)^T = B^T A^T$

$$-2 \vec{t}^T \underline{\tilde{X}}$$

derivé = 0

$$2 \underline{\tilde{X}}^T \underline{\tilde{X}} \vec{\beta}$$

$$\vec{t}^T \underline{\tilde{X}} = [t^{(2)} \dots t^{(N)}]$$

$$\in \mathbb{R}^{1 \times N}$$

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

$\mathbb{R}^{N \times (D+1)}$

$$= \vec{v}^T \vec{\beta} \quad \text{ou} \quad v = \vec{t}^T \underline{\tilde{X}}$$

$$\frac{\partial (\vec{v}^T \vec{\beta})}{\partial \vec{\beta}} = \frac{\partial}{\partial \vec{\beta}} (\underbrace{v_0}_{\beta_0} + \underbrace{v_1}_{\beta_1} + \dots + \underbrace{v_D}_{\beta_D}) = [\underbrace{v_0}_{\beta_0} \quad \underbrace{v_1}_{\beta_1} \quad \dots \quad \underbrace{v_D}_{\beta_D}]$$

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

$$\Rightarrow \frac{\partial}{\partial \beta} (\underline{\underline{E^T \tilde{X} \vec{\beta}}}) = \underline{\underline{2E^T \tilde{X}}}$$

$$[\beta_1 A_{11} + \beta_2 A_{12} \quad \beta_1 A_{12} + \beta_2 A_{22}]$$

pour le terme d'ordre 2

$$\frac{\partial}{\partial \beta} (\beta^T A \beta) \rightarrow [\beta_1 \quad \beta_2] \begin{bmatrix} A_{11} & A_{12} \\ A_{12} & A_{22} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$$

$$= (A_{11} \beta_1^2 + A_{22} \beta_2^2 + 2A_{12} \beta_1 \beta_2)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \beta_1} &= 2 A_{11} \beta_1 + 2 A_{12} \beta_2 \\ \frac{\partial \mathcal{L}}{\partial \beta_2} &= 2 A_{22} \beta_2 + 2 A_{12} \beta_1 \end{aligned} \quad = \frac{\partial (\beta^T A \beta)}{\partial \beta} = 2 \begin{bmatrix} A_{11} & A_{12} \\ A_{12} & A_{22} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = 2 A \beta$$

→ Regroupant les dérivés on obtient

$$\frac{\partial \mathcal{L}}{\partial \beta} = 2 \tilde{X}^T \tilde{X} \vec{\beta} - 2 \tilde{X}^T \vec{t} = 0$$

pour trouver
le minimum

$$A x - b = 0$$

$$\rightarrow \begin{cases} \bar{X}^T \tilde{X} \bar{\beta} = \tilde{X}^T \tilde{E} \\ \underline{\underline{=}} \underline{\underline{=}} \underline{\underline{\beta}} = \underline{\underline{\tilde{X}^T \tilde{E}}} \end{cases}$$

Equations Normales

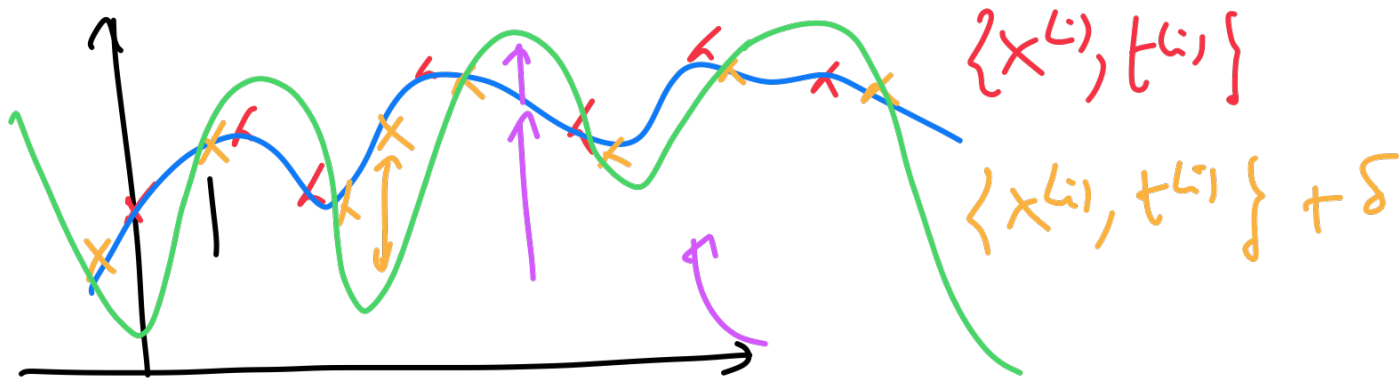
$$\underline{\underline{\beta}} = \underbrace{(\underline{\underline{\tilde{X}^T \tilde{X}}})^{-1}}_{\text{E} + \delta} \underline{\underline{\tilde{X}^T \tilde{E}}}$$

→ solution du problème de minimisation de la somme des carrés

→ OK si la matrice est inversible.

→ ? si la matrice est mal conditionnée

$$\begin{bmatrix} \lambda_1 & & \\ & \dots & \\ & & \lambda_{D+1} \end{bmatrix}^{-1} \rightarrow \begin{bmatrix} \lambda_1^{-2} & & \\ & \dots & \\ & & \lambda_{D+1}^{-2} \end{bmatrix}$$



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

→ si colonnes sont linéairement
dépendantes?

→ x_j sont corrélés

→ première approche : extraire un ensemble de colonnes orthogonales

Soit $\underline{z}_0 = \vec{c}_0$ la première colonne de $\underline{\tilde{X}} = \begin{bmatrix} | & | & \dots & | \\ \vec{c}_0 & \vec{c}_1 & \dots & \vec{c}_{D+1} \\ | & | & \dots & | \end{bmatrix}$

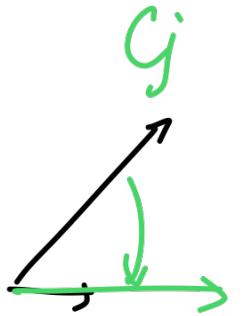
for $j = 1, \dots, D+1$

Calculer les projections $\hat{\gamma}_{lj} = \frac{\langle \vec{z}_l, \vec{c}_j \rangle}{\langle \vec{z}_l, \vec{z}_l \rangle}$

pour $l = 0, \dots, j-1$

Définir la nouvelle colonne \vec{z}_j via

$$\vec{z}_j = \vec{c}_j - \sum_{k=0}^{j-1} \hat{\gamma}_{kj} \vec{z}_k$$



Minimiser la fonction de coût $l(\beta) = \frac{1}{N} (\vec{E} - \underline{\tilde{X}} \beta)^T (\underline{\tilde{X}} \beta - \vec{E})$

$$\beta_0 \vec{c}_0 + \beta_1 \vec{c}_1 + \dots$$

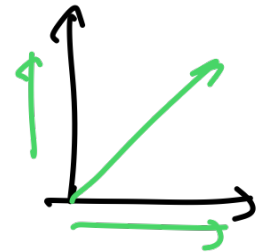
$\vec{\beta}$ → meilleure combinaison linéaire de colonnes de $\underline{\tilde{X}}$
au KLS des moindres carrés telle que $\underline{\tilde{E}} \approx \underline{\tilde{X}} \vec{\beta}$

$$\underline{\beta} = \left(\underline{\tilde{X}}^T \underline{\tilde{X}} \right)^{-1} \underline{\tilde{X}}^T \underline{\tilde{E}} \rightarrow \text{projection orthogonale sur}$$

l'espace des colonnes de $\underline{\tilde{X}}$

projection sur les z → On peut la calculer via

$$\alpha_l = \frac{\langle \vec{E}, z_l \rangle}{\langle z_l, z_l \rangle}$$



Plusieurs solutions pour améliorer le conditionnement

Best Subset Selection

→ Nombre de sous-ensembles de taille k d'un ensemble de caractéristiques de taille m : $\binom{D+1}{m}$

→ taille 2 taille 3 → nombre total de modèles à tester:

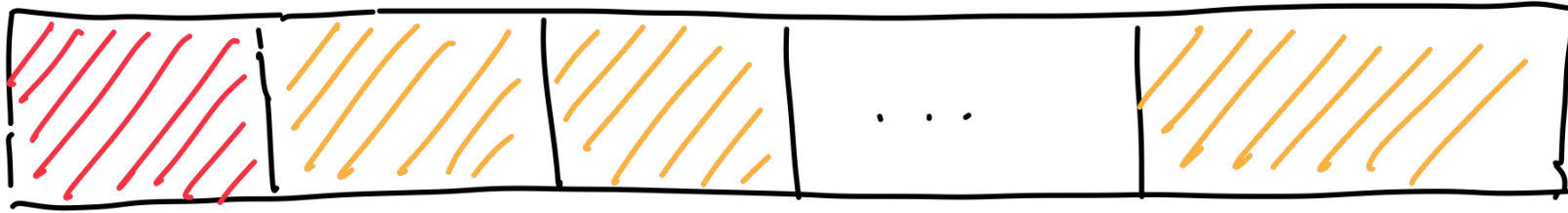
$\beta_0 \beta_1$	$\beta_0 \beta_1 \beta_2$	$\sum_{m=2}^{D+1} \binom{D+1}{m}$
$\beta_0 \beta_2$	$\beta_0 \beta_1 \beta_3$	
$\beta_0 \beta_3$	\vdots	
\vdots	$\beta_0 \beta_2 \beta_4$	
\vdots		

→ Validation croisée

→ ensemble d'entraînement
ensemble de test

→ Si bcp de données → il suffit de séparer l'ensemble
complet de données en un
ensemble test et un
ensemble d'entraînement

→ Si peu de données → k fold cross validation



test

K sous ensemble

↳ training

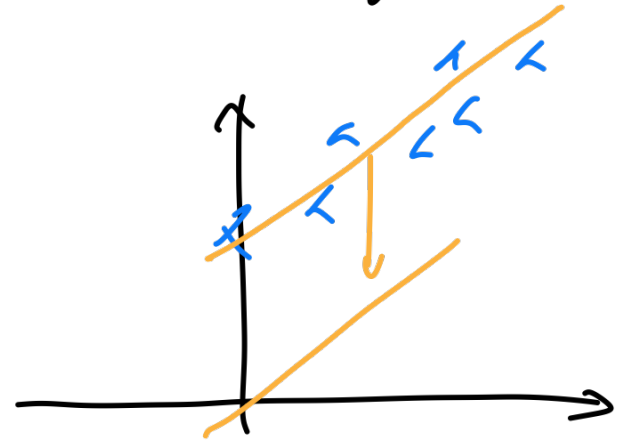
$$CV_K(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - h_{\beta}(x^{(i)}))^2$$

n'importe quel sous ensemble des coefficients de regression

donnée du terrain
↓

Modèle appris sur les données "hors terrain"

approche #2 → On étend la lecture de Coût en ajoutant
une pénalité sur les β_j .



→ Regression Ridge

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

→ Comment calculer Ridge ?

fit aux données

Model
Complexity

