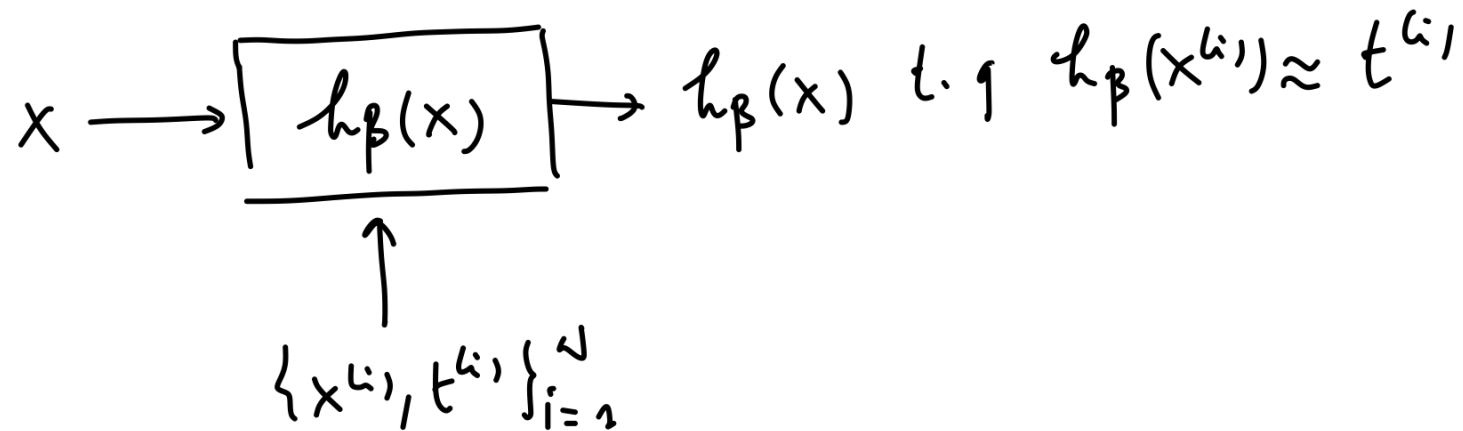


Apprentissage Supervisé $\{x^{(i)}, t^{(i)}\}_{i=1}^N$



→ Apprentissage : 2 approches : - Descente de gradient
- Résolution des équations normales

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

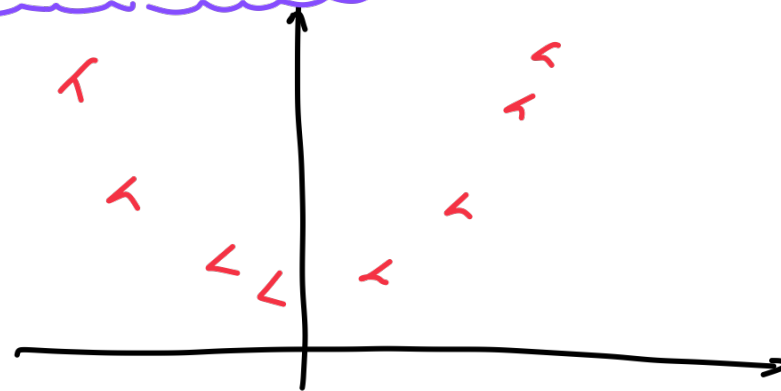
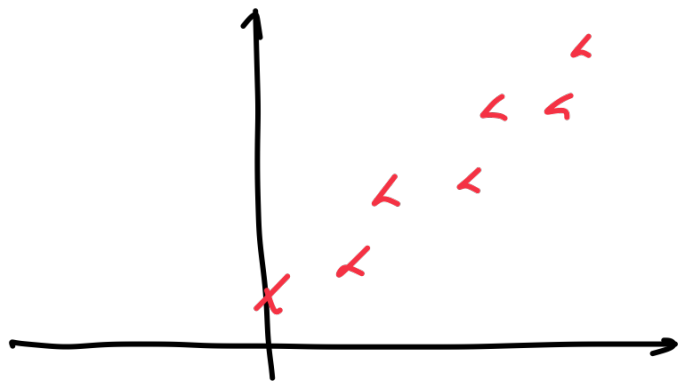
Descente de gradient : $\vec{\beta}^{(0)} \in \mathbb{R}^{D+1}$

for iter in MaxIter

$$\vec{\beta} \leftarrow \vec{\beta} - \eta \text{grad } \ell(\beta)$$

Equations normales

$$\vec{\beta}_{OLS} = (\underline{\tilde{X}}^T \underline{\tilde{X}})^{-1} \underline{\tilde{X}}^T \underline{t}$$



Pour D suffisamment grand, il existera toujours un modèle linéaire capable de "filtrer" les observations

→ ? $\tilde{X}^T \tilde{X}$ non inversible ? ex: quand certaines caractéristiques sont dépendantes.

$$\langle x, y \rangle = X^T y$$

→ Première solution: Gram-Schmidt

Démarrer avec $c_0 = z_0$ la première colonne de \tilde{X}

Fer $j = 1, \dots, D$

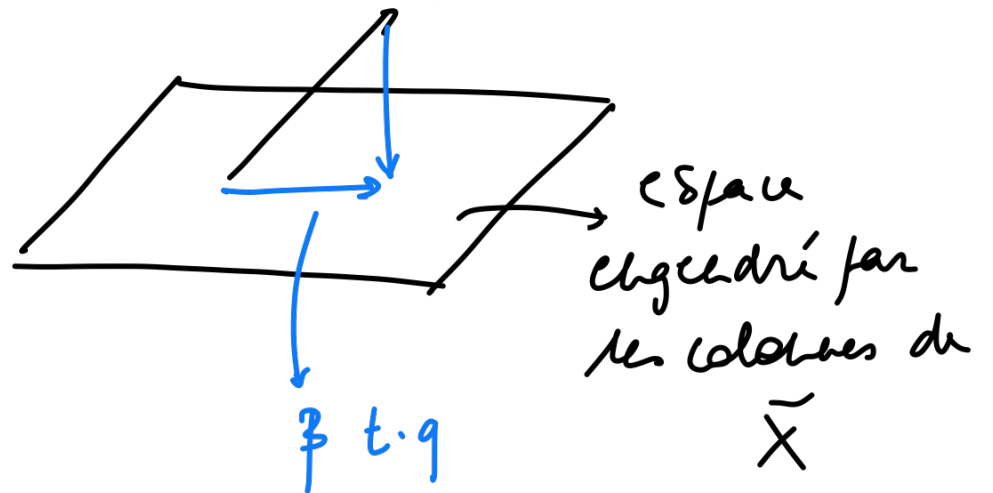
Calculer les coefficients $\hat{\gamma}_{lj} = \langle c_j, z_l \rangle / \langle z_l, z_l \rangle$

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

On définit la nouvelle colonne $\underbrace{z_j}_{\text{green}} = \underbrace{c_j}_{\text{green}} - \underbrace{\sum_{l=0}^{j-1} \hat{\gamma}_{lj} z_l}_{\text{green}}$

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

$$\min_{\beta} \frac{1}{N} \|\vec{t} - \tilde{X}\beta\|_2^2 \quad (*)$$



$\|\vec{t} - \tilde{X}\beta\|$ est la plus petite possible

la solution de (*) peut être facilement calculée dans la base des vecteurs z_j via

$$\frac{\langle \vec{t}, z_j \rangle}{\langle z_j, z_j \rangle} = \alpha_j$$

$$\begin{aligned} \text{or a } h_{\beta}(x) &= \sum_j \alpha_j z_j = \alpha_D z_D + \sum_{j=0}^{D-1} \alpha_j z_j \\ &= \alpha_D C_D + \dots + \sum_{j=0}^{D-1} \alpha_j z_j \end{aligned}$$

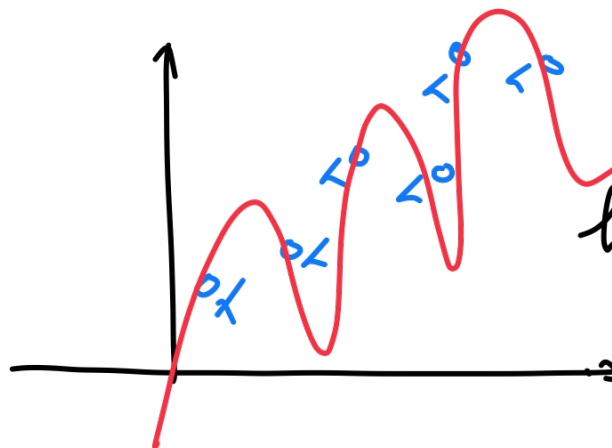
$$\beta_D = \alpha_D = \frac{\langle t, z_D \rangle}{\langle z_D, z_D \rangle}$$

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

$$= \beta^T \tilde{x}$$

$$\{x^{(i)}, t^{(i)}\}_{i=1}^N$$

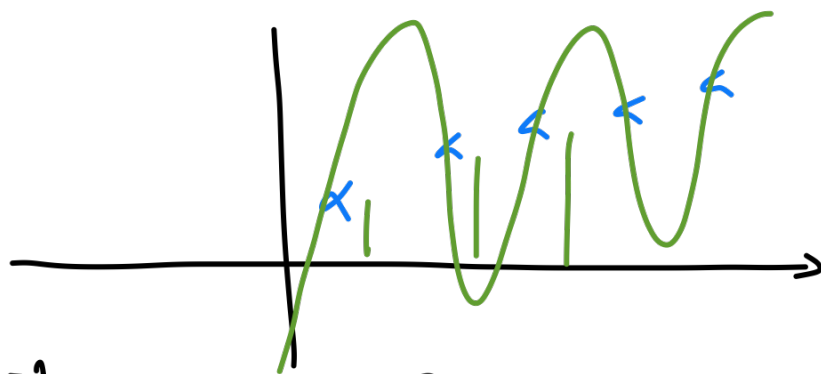
$$\{x^{(i)}, t^{(i)} + \Delta\}_{i=1}^N$$



$$h_\beta(\tilde{x} + \Delta) = \tilde{\beta}^T (\tilde{x} + \Delta)$$

$$h_\beta(\tilde{x}) = \tilde{\beta}^T \tilde{x}$$

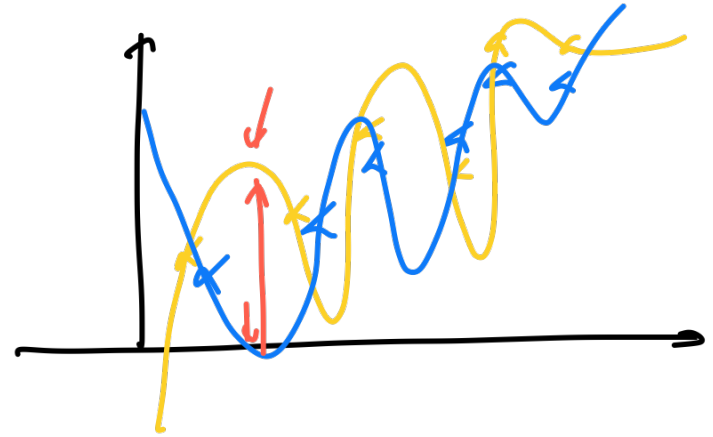
$$\beta_{OLS} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T t$$



$$\begin{bmatrix} \lambda_2 & & & \\ & \ddots & & \\ & & \lambda_{D+1} & \\ & & & \ddots \end{bmatrix}^{-1} = \begin{bmatrix} \lambda_2^{-1} & & & \\ & \ddots & & \\ & & \lambda_{D+1}^{-1} & \\ & & & \ddots \end{bmatrix}$$

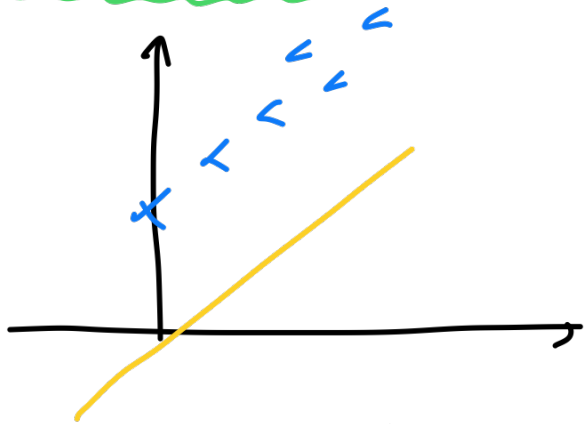
$$\beta_{OLS}^2 = (\tilde{X}^T \tilde{X})^{-2} \tilde{X}^T t$$

$$\beta_{OLS}^2 = (\tilde{X}^T \tilde{X})^{-2} \tilde{X}^T (t + \Delta)$$

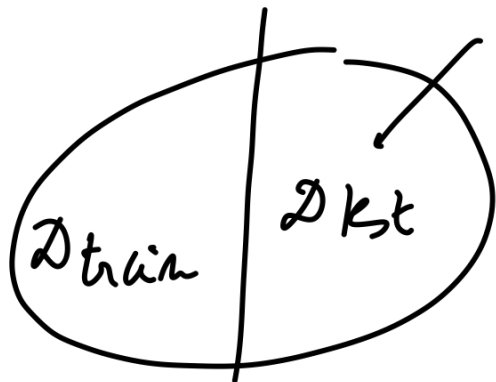


(Sélection de caractéristiques)

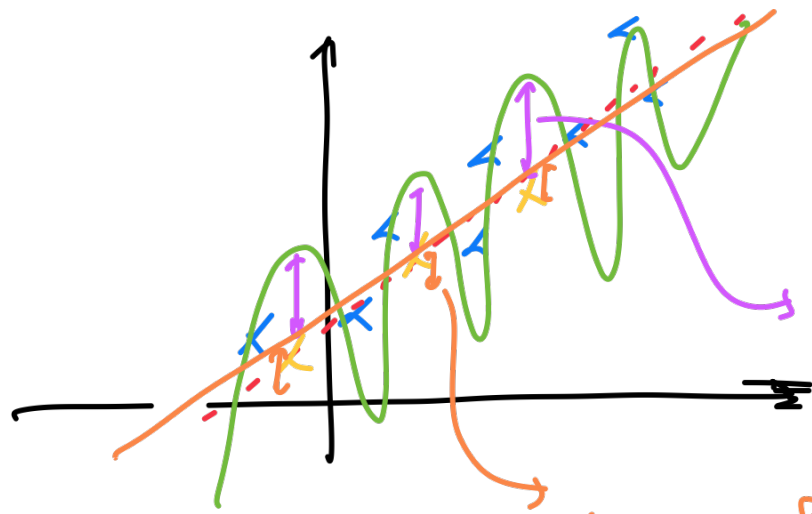
Approche 2: Sélection du meilleur sous ensemble



$$\sum_{k=1}^D \binom{D}{k} \text{ sous ensemble}$$



→ 2 étapes: Pour chaque sous ensemble, on entraîne le modèle sur D_{train} et on calcule l'erreur de prédiction sur D_{test}



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

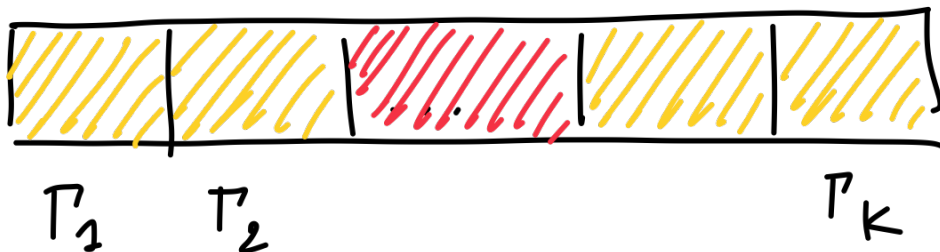
erreur sur D_{test}

$D \gg \gg 1$

erreur sur D_{test} pour un modèle

plus simple

Dans le cas où l'ensemble des observations est limité on se tourne vers la validation croisée (à K compartiments)



1 Compartiment de test et $k-1$ compartiments d'entraînement

$$\text{erreur}_{CV} = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - h_{\beta}^{-k(i)}(x^{(i)}))^2$$

$h_{\beta}^{-k(i)}$ = le modèle entraîné sur tous les compartiments
sauf celui qui contient l'observation $\{x^{(i)}, t^{(i)}\}$

Approche #3 :

terme de fidélité aux données

penalti sur la
complexité du
modèle

$$\mathcal{L}(\beta) = \underbrace{\frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}_{\text{terme de fidélité aux données}} + \underbrace{\lambda \sum_{j=1}^D \beta_j^2}_{\text{penalti sur la complexité du modèle}}$$

→ Modèle de Régression Ridge (Régularisation)

β_1 β_2 pénalité Ridge : $\beta_1^2 + \beta_2^2$

$$\min_{\beta} \ell(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)}))^2 + \lambda \sum_{j=1}^2 \beta_j^2$$

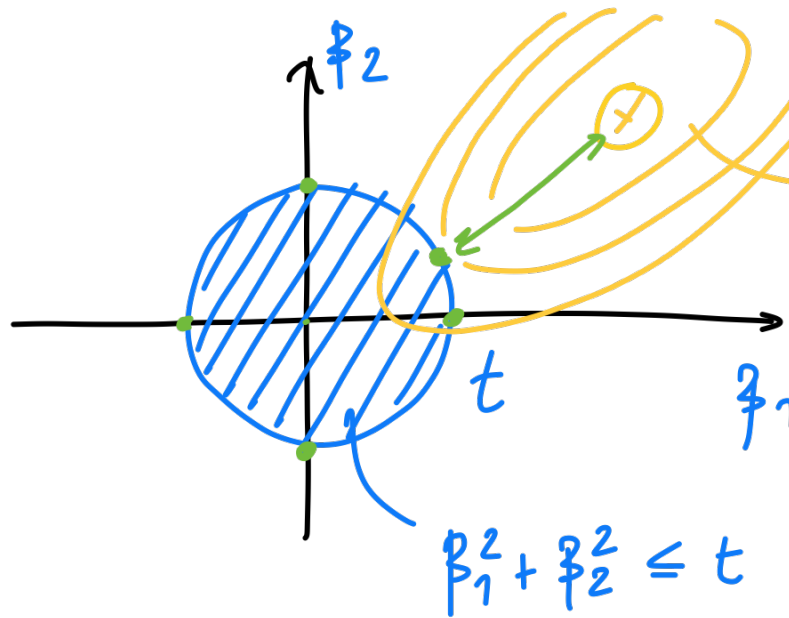
(formulation non contrainte)

$$\min_{\beta} \ell(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)}))^2$$

tel que $\sum_{j=1}^2 \beta_j^2 \leq t^2$

(Formulation contrainte)

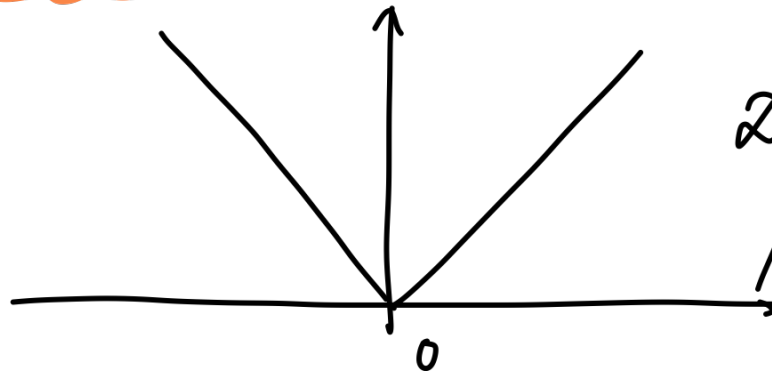
$$\beta_1^2 + \beta_2^2 \leq t^2$$



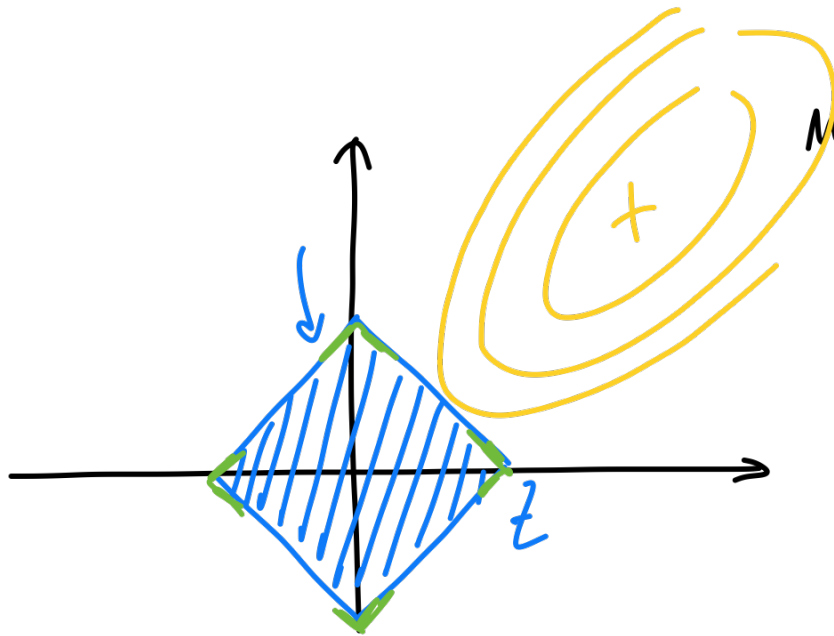
$$\frac{1}{N} \sum_{i=1}^N (t^{(i)} - \dots)$$

Approche 4: Regression LASSO

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



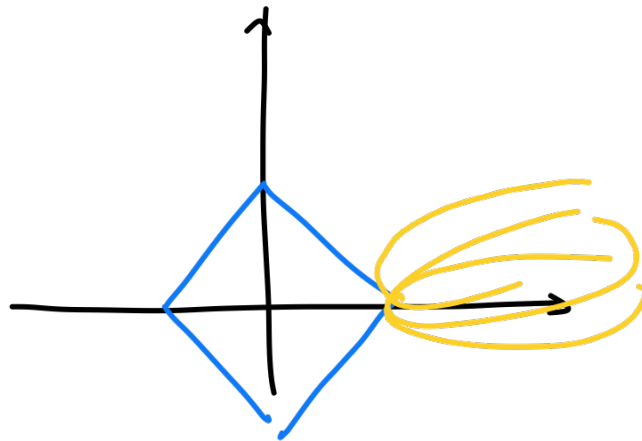
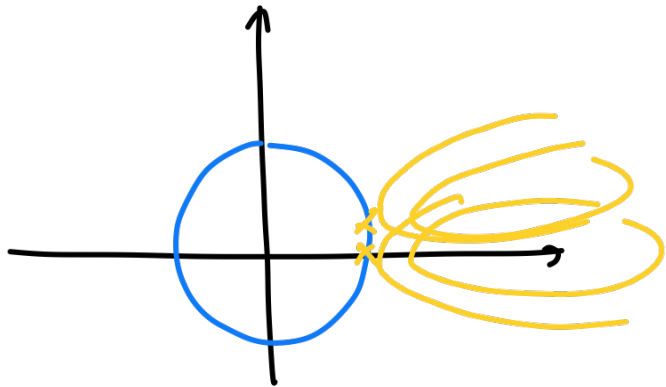
Désavantage:
perte de la dérivabilité



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

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

LASSO



$$\|\vec{\beta}\|_2 = \sqrt{\sum_{j=1}^D \beta_j^2}$$

$$\|\vec{\beta}\|_1 = \sum_{j=1}^D |\beta_j|$$

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

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

→ norme l_p

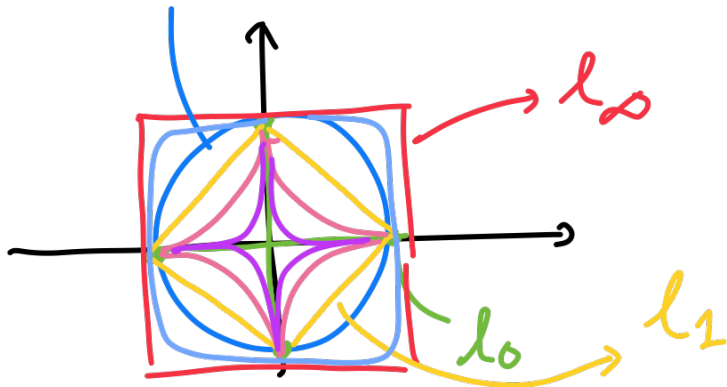
A chacune des normes l_p on associe la boule l_p

parmi lesquelles la boule l_1

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

l_2 la boule l_2

$$\sum_{j=1}^D |\beta_j|^2 \leq t$$



Ridge → pour rappel : Fonction somme des résidus Carrés

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

$$\rightarrow l(\beta) = \frac{1}{N} \vec{E}^T \vec{E} - \frac{1}{N} \vec{E}^T \tilde{X} \beta - \frac{1}{N} \beta^T \tilde{X}^T \vec{E} + \frac{1}{N} \beta^T \tilde{X}^T \tilde{X} \beta + \lambda \beta^T \beta$$

$$\text{grad } l(\beta) = -\frac{2}{N} \tilde{X}^T \vec{E} + \frac{2}{N} \tilde{X}^T \tilde{X} \beta + 2\lambda \beta$$

$$\text{grad } l(\beta) = 0 \Rightarrow (\tilde{X}^T \tilde{X} + \lambda \mathbf{I}) \beta = \tilde{X}^T \vec{E}$$

$$\beta_{\text{Ridge}} = (\tilde{X}^T \tilde{X} + \lambda \mathbf{I})^{-1} \tilde{X}^T \vec{E}$$

Problème avec OLS c'est que lorsqu'on a de petites valeurs propres, l'inverse sera grande.

$$\underset{=}{\tilde{X}}^T \underset{=}{\tilde{X}} v = \xi v \rightarrow (\underset{=}{\tilde{X}}^T \underset{=}{\tilde{X}} + \lambda I) v$$

$$\underset{=}{\tilde{X}}^T \underset{=}{\tilde{X}} \vec{v} + \lambda \vec{v} = \xi \vec{v} + \lambda \vec{v} \\ = (\xi + \lambda) \vec{v}$$

Soit (v, ξ) une

paire vecteur propre
valeur propre de $\tilde{X}^T \tilde{X}$

$\tilde{X}^T \tilde{X}$ est semi-définie positive \rightarrow toutes les valeurs propres
sont positives ou nulles

Même si ϵ est petit pour un n suffisamment grand,
On peut contrôler l'amplitude de l'inverse, et donc des
coefficients de régression.

Supposons que les données sont générés via

$$t^{(i)} = \beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)} + \epsilon^{(i)}, \quad \epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$$

$\epsilon^{(i)}$ indépendants

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

identiquement
distribués

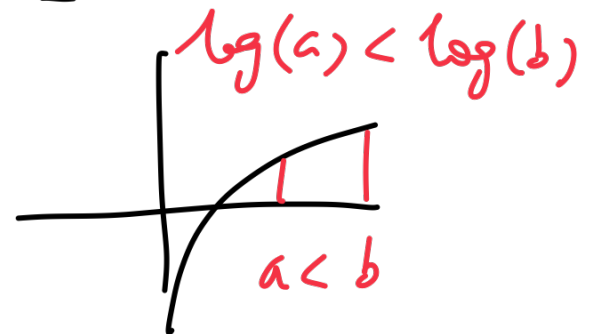
$$t^{(i)} \sim \mathcal{N}(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}, \sigma^2)$$

$$p(t^{(i)} | \beta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}{2\sigma^2}\right)$$

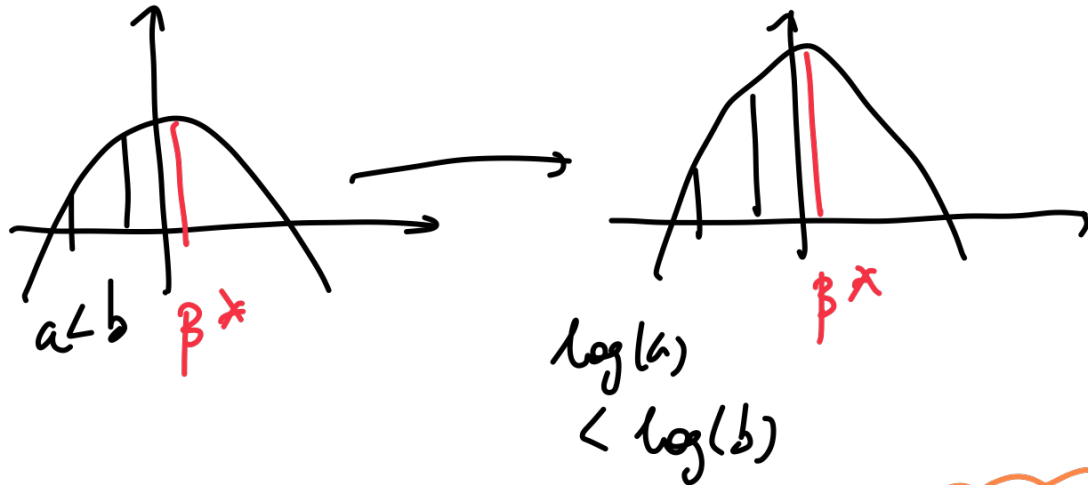
$$p(\{t^{(i)}\}_{i=1}^N | \beta) = \prod_{i=1}^N p(t^{(i)} | \beta)$$

$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}{2\sigma^2}\right)$$

$$\vec{\beta} = \operatorname{argmax}_{\beta} p(\{t^{(i)}\}_{i=1}^N | \beta)$$



$$\vec{\beta} = \arg \max_{\beta} \log \left(\prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left(- \frac{(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}{2\sigma^2} \right) \right)$$



Fonction de log vraisemblance

$$\vec{\beta} = \arg \max_{\beta} \left\{ \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \sum_{i=1}^N \frac{(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}{2\sigma^2} \right\}$$

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

→ le Modèle de régression linéaire entraîné via la minimisation de la fonction de coût définie par la somme des carrés des résidus correspond à l'estimateur de moindres carrés.

Aucune hypothèse sur le Modèle (i.e coefficients) → MLE
 $f(t^{(i)} | \beta)$

Maximum a Posterior: $f(\beta | t^{(i)})$

Règle de Bayes: $f(B|A) = \frac{f(A|B)f(B)}{f(A)}$

$$\text{MAP} : \underbrace{p(\beta | t^{(i)})}_{\text{probabilité a posteriori}} = \frac{p(t^{(i)} | \beta) p(\beta)}{p(t^{(i)})} \rightarrow \text{ne depend pas de } \beta$$

Etant donné $p(\beta | t^{(i)})$ on peut définir β_{MAP} comme l'estimateur qui maximise cette probabilité

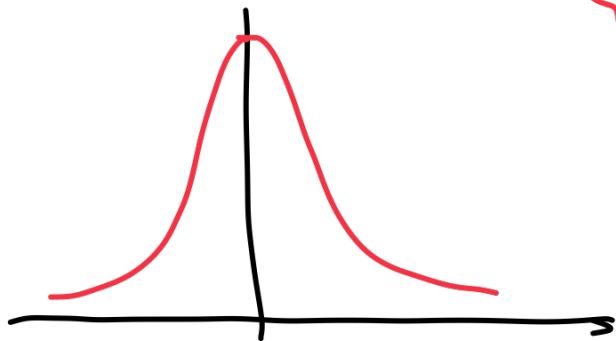
$$\beta_{\text{MAP}} = \operatorname{argmax}_{\beta} \underbrace{p(t^{(i)} | \beta)}_{\text{A priori sur le modèle}} p(\beta)$$

Si $p(\beta)$ est une forme de régularité le β_{MLE} .

→ Un premier exemple d'a priori sur β consisterait à favoriser les coefficients qui sont pour la plupart

Proches de 0 \rightarrow

$$p(\beta) = \prod_{j=1}^D p(\beta_j) = \prod_{j=1}^D \frac{1}{\sqrt{2\pi}\xi} \exp\left(-\frac{\beta_j^2}{2\xi^2}\right)$$



log

$$p(\beta | \{t^{(i)}\}_{i=1}^N) \propto \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}{2\sigma^2}\right)$$
$$\times \prod_{j=1}^D \frac{1}{\sqrt{2\pi}\xi} \exp\left(-\frac{\beta_j^2}{2\xi^2}\right)$$

$$\vec{\beta}_{\text{MAP}} = \underset{\beta}{\operatorname{argmax}} p(\beta | \{t^{(i)}\}_{i=1}^N)$$

$$\vec{\beta}_{\text{MAP}} = \arg \max_{\beta} \log (p(\beta | \{t^{(i)}\}_{i=1}^N))$$

$$= \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}{2\sigma^2}$$

$$+ \sum_{j=1}^D \log \left(\frac{1}{\sqrt{2\pi}\xi} \right) - \frac{\beta_j^2}{2\xi^2}$$

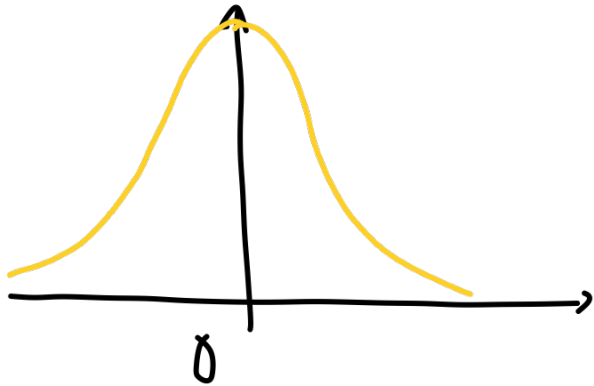
constants par rapport à β

$$\vec{\beta}_{\text{MAP}} = \arg \max_{\beta} - \sum_{i=1}^N \frac{(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}{2\sigma^2} - \sum_{j=1}^D \frac{\beta_j^2}{2\xi^2}$$

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

λ

Formulation Ridge



$$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

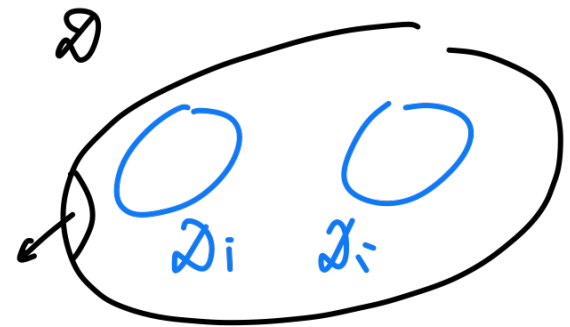


$$\frac{1}{2\lambda} \exp\left(-\frac{|x|}{\lambda}\right)$$

Dans le cas de la distribution de Laplace, le Maximum a posteriori retrouve l'estimateur LASSO

Equilibre "biais-variance"

Erreur quadratique moyenne $\{x, t(x)\}$



$$\text{MSE}(x) = \mathbb{E}_{D_i} \left\{ (t(x) - h_{D_i}(x))^2 \right\}$$

$h_{D_i}(x) =$ modèle entraîné sur

$$= \mathbb{E}_{D_i} \left\{ \left(t(x) - \mathbb{E}_{D_i} h_{D_i}(x) + \mathbb{E}_{D_i} h_{D_i}(x) - h_{D_i}(x) \right)^2 \right\} \text{ les déviés de } D_i$$

$$= \mathbb{E}_{\mathcal{D}_i} \left\{ \left(t(x) - \mathbb{E}_{\mathcal{D}_i} h_{\mathcal{D}_i}(x) \right)^2 \right\} \leftarrow \text{bias}^2$$

$$+ \mathbb{E}_{\mathcal{D}_i} \left\{ \left(\mathbb{E}_{\mathcal{D}_i} h_{\mathcal{D}_i}(x) - h_{\mathcal{D}_i}(x) \right)^2 \right\} \leftarrow \text{Variance}$$

$$+ 2 \mathbb{E}_{\mathcal{D}_i} \left\{ \left(t(x) - \mathbb{E}_{\mathcal{D}_i} h_{\mathcal{D}_i}(x) \right) \left(\mathbb{E}_{\mathcal{D}_i} h_{\mathcal{D}_i}(x) - h_{\mathcal{D}_i}(x) \right) \right\}$$

$$= 2 \left(t(x) - \mathbb{E}_{\mathcal{D}_i} h_{\mathcal{D}_i}(x) \right) \mathbb{E}_{\mathcal{D}_i} \left\{ \left(\mathbb{E}_{\mathcal{D}_i} h_{\mathcal{D}_i}(x) - h_{\mathcal{D}_i}(x) \right) \right\}$$

$$\mathbb{E}_{\mathcal{D}_i} h_{\mathcal{D}_i}(x) - \mathbb{E}_{\mathcal{D}_i} h_{\mathcal{D}_i}(x)$$

Complexité augmente \rightarrow variance \nearrow
biais \searrow

Complexité diminue \rightarrow variance \searrow
biais \nearrow

