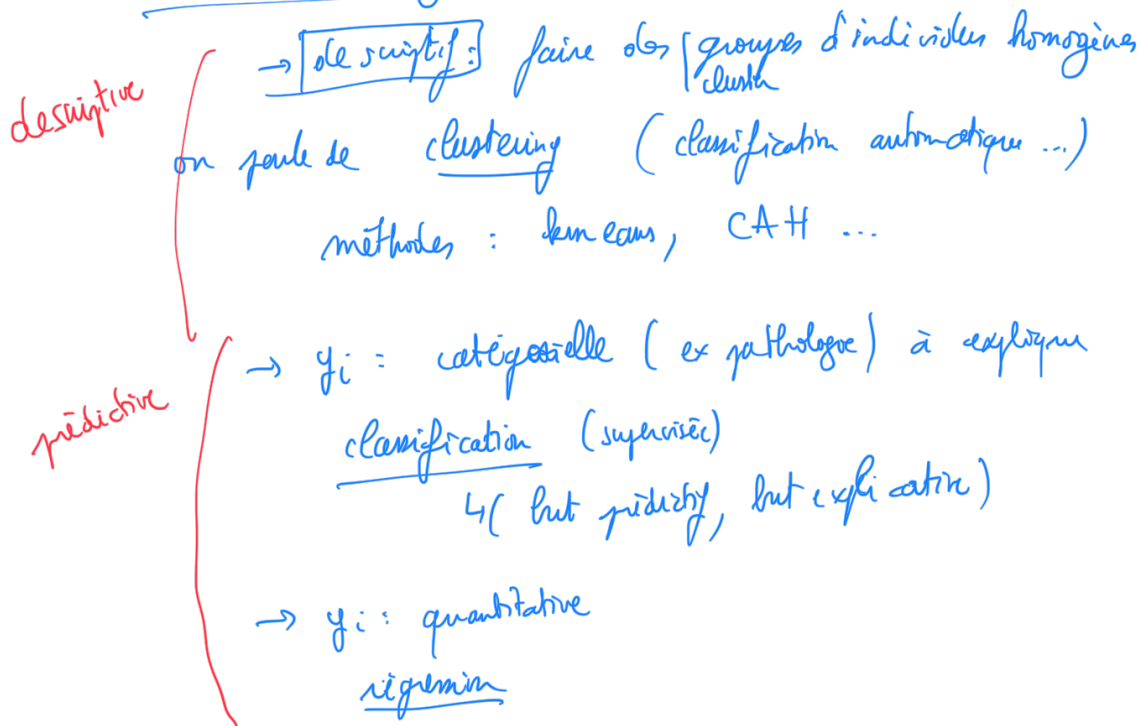


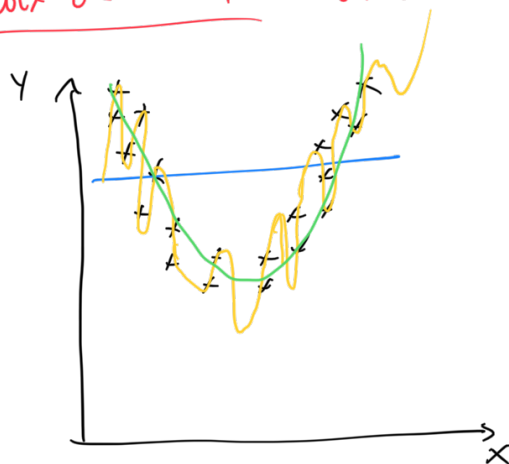
[19/09] → Jeu de données: $(x_{ij})_{\substack{1 \leq i \leq n \\ 1 \leq j \leq p}}$

x_{ij} = valeur de la variable j pour l'observation i

Que veut-on faire de ces données



Choix de modèles on veut comparer des modèles de régression



$$y_i = a + b x_i + \varepsilon_i$$

- ☹️ résidus ε_i très grands, un grand biais
- ☹️ peu sensible aux fluctuations d'éch
- ☺️ ajustement: une variance faible

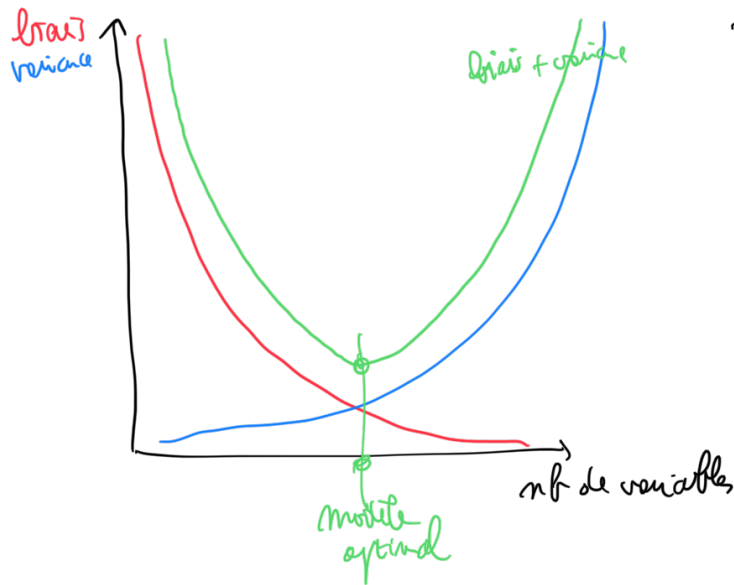
$$y_i = a + b_2 x_i^2 + b_3 x_i^3 + \dots + b_1 x_i + \varepsilon_i$$

- ☹️ résidus très faibles = biais faible
- ☹️ variance très grande
- ☹️ OVER-ADJUSTMENT

$$y_i = a + b_2 x_i^2 + b_1 x_i + \varepsilon_i$$

biais faible, variance faible

] bon compromis: biais réduite



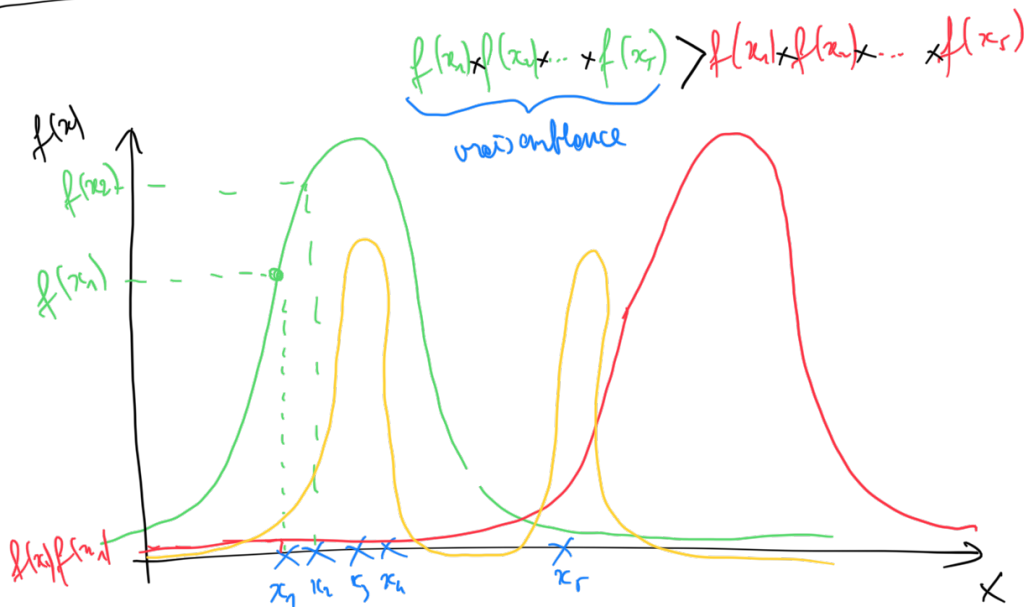
En pratique on va comparer les modèles

↳ critères AIC, BIC

↳ évaluer la performance sur des données indépendantes (qui n'ont pas servi à estimer les modèles)

- apprentissage / test
- validation croisée.

↳ R^2 ajusté
(R^2 croît en fonction du nb de variables)



⚠ À modèle fixe, on estime généralement les paramètres par maximum de vraisemblance

Mais on ne l'utilise pas pour comparer des modèles, car la vraisemblance croît en fonction du nombre des paramètres.

Régression généralisée

modèle $y_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

Max de vraisemblance \approx Moindres carrés (Ordinary Least Square OLS)

on cherche $\beta = (\beta_0, \dots, \beta_p)$ qui minimise

$$\sum_i \left(y_i - \left(\beta_0 + \sum_{j=1}^p \beta_j x_{ij} \right) \right)^2$$

Si p grand, et/ou si variables sont corrélées ($1/\hat{\beta}_j$ grande)
on observe des $\hat{\beta}_j$ qui peuvent prendre de grande valeur
suivant les fluctuations d'échantillonnage

$$\hat{\beta}^{OLS} = (X^T X)^{-1} X^T Y$$

On va donc résoudre un problème amarré en empêchant les $\hat{\beta}_j$ de prendre des valeurs trop grandes

\Rightarrow on cherche $\beta = (\beta_0, \dots, \beta_p)$ qui minimise
on cherche des β qui

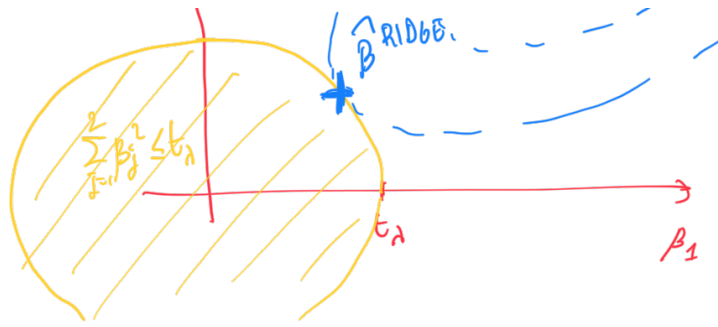
$$\underbrace{\sum_{i=1}^n \left(y_i - \left(\beta_0 + \sum_j \beta_j x_{ij} \right) \right)^2}_{\text{donne de petits résidus}} + \underbrace{\lambda \sum_j \beta_j^2}_{\text{sans prendre de valeurs trop grande}}$$

\Leftrightarrow ou de façon équivalente, on cherche β

qui minimise $\left\| \sum_i (y_i - (\beta_0 + \sum_j \beta_j x_{ij}))^2 \right\|$ sous
la contrainte $\sum_j \beta_j^2 \leq t_\lambda$

où λ (t_λ) sont des hyper-paramètres à régler.





La solution est donnée par $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

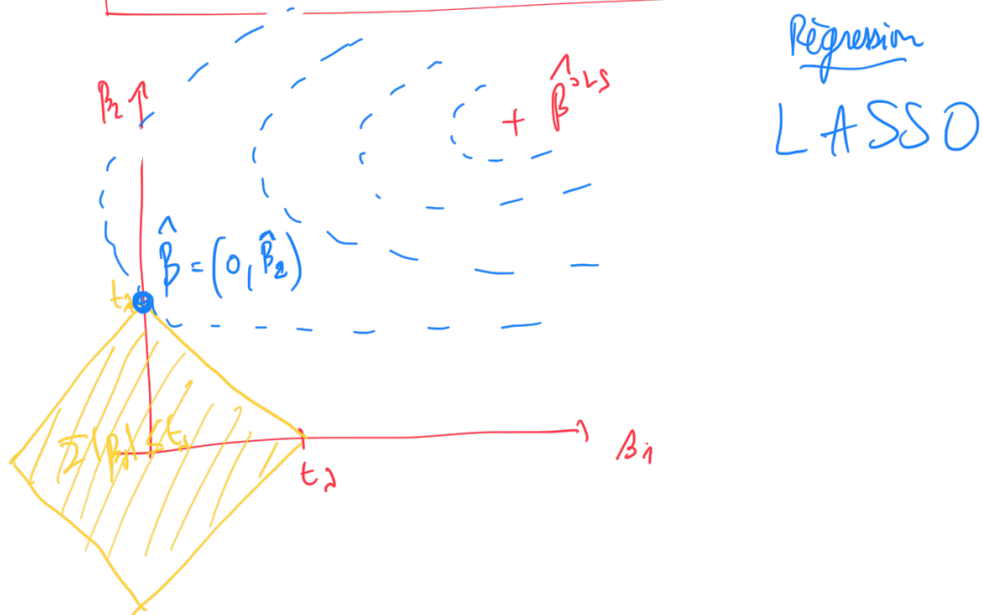
$$\hat{\beta}^{RIDGE} = (X^t X + \lambda I)^{-1} X^t y$$

Régression régularisée / pénalisée.

Changeons la pénalité $\sum \beta_j^2$ par $\sum |\beta_j|$

⇒ On cherche β qui minimise

$$\sum_i (y_i - (\beta_0 + \sum_j \beta_j x_{ij}))^2 \quad \text{sous la contrainte } \sum_j |\beta_j|$$



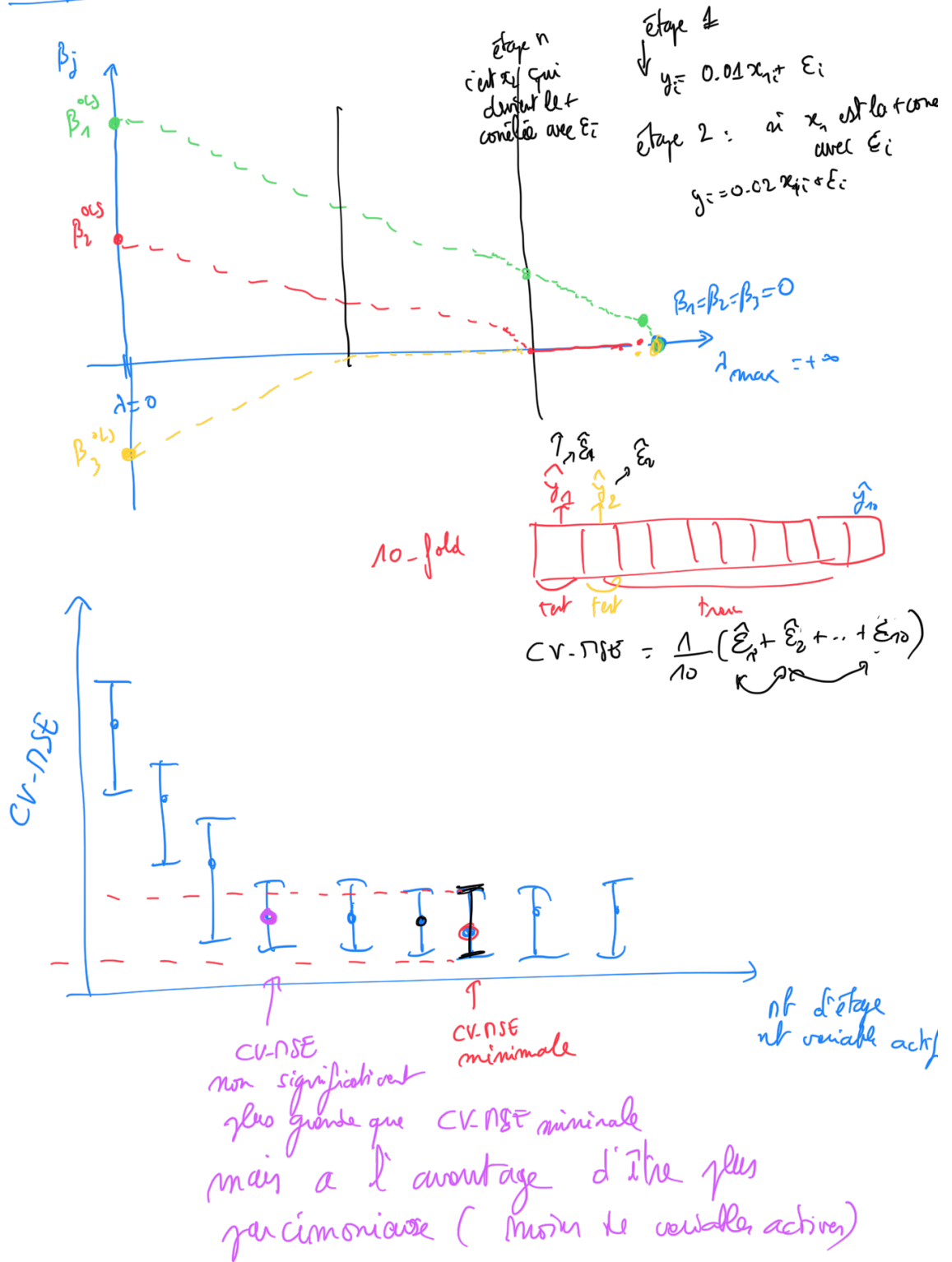
Pb → pas de solution analytique $\hat{\beta}^{LASSO} = ??$

→ on utilisera un algorithme d'optimisation

... qui nous donnera une

même/unique qui nous donne la solution optimale de β LASSO.

Algorithme CARLS



20/09

- Lorsque j'ai beaucoup d'individus (lignes) dans ma base

de données, c'est un avantage d'un point de vue statistique

oui ☐ non ☐

- En présence de variables continues et catégorielles

☐ je cherche à utiliser des méthodes qui permettent nativement de prendre en compte ces deux types de données

☐ je discrétise les variables continues en variables catégorielles

☐ je transforme les données catégorielles en données quantitatives

- La présence de corrélation dans les covariables induit:

☐ un biais dans les estimations $\hat{\beta}$ de β

☐ un accroissement de $V(\hat{\beta})$

- Si j'ai un très grand nombre de variables (10 000) et peu d'individus (100), je peux réaliser:

☐ une régression linéaire classique

☐ une sélection de variables backward

☐ une sélection de variables forward

☐ une régression ridge

☐ une régression LASSO

Et parmi celles que je peux réaliser, les classer par ordre de "préférence"

- Si je suis plus intéressé à expliquer les variables importantes dans le modèle plutôt qu'à avoir une bonne qualité de prédiction, je peux réaliser:

□ eine regression LASSO

- une régression PLS

- une sélection stepwise de variable.