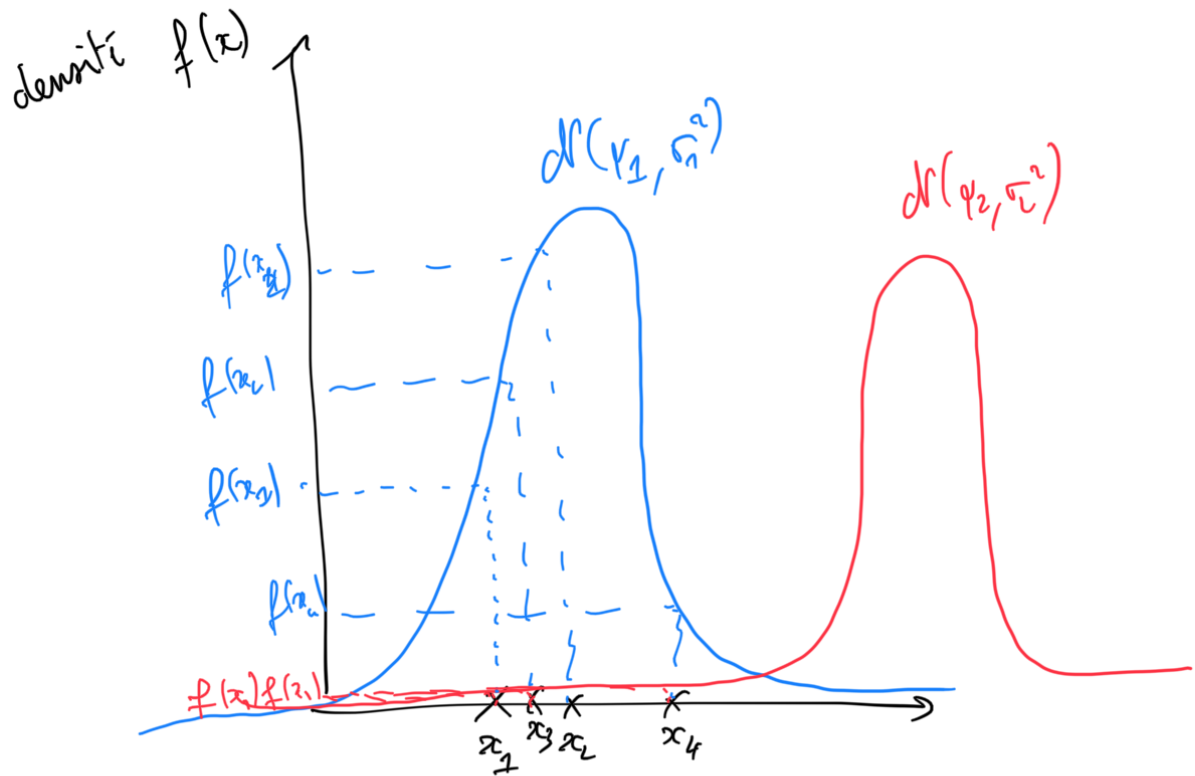
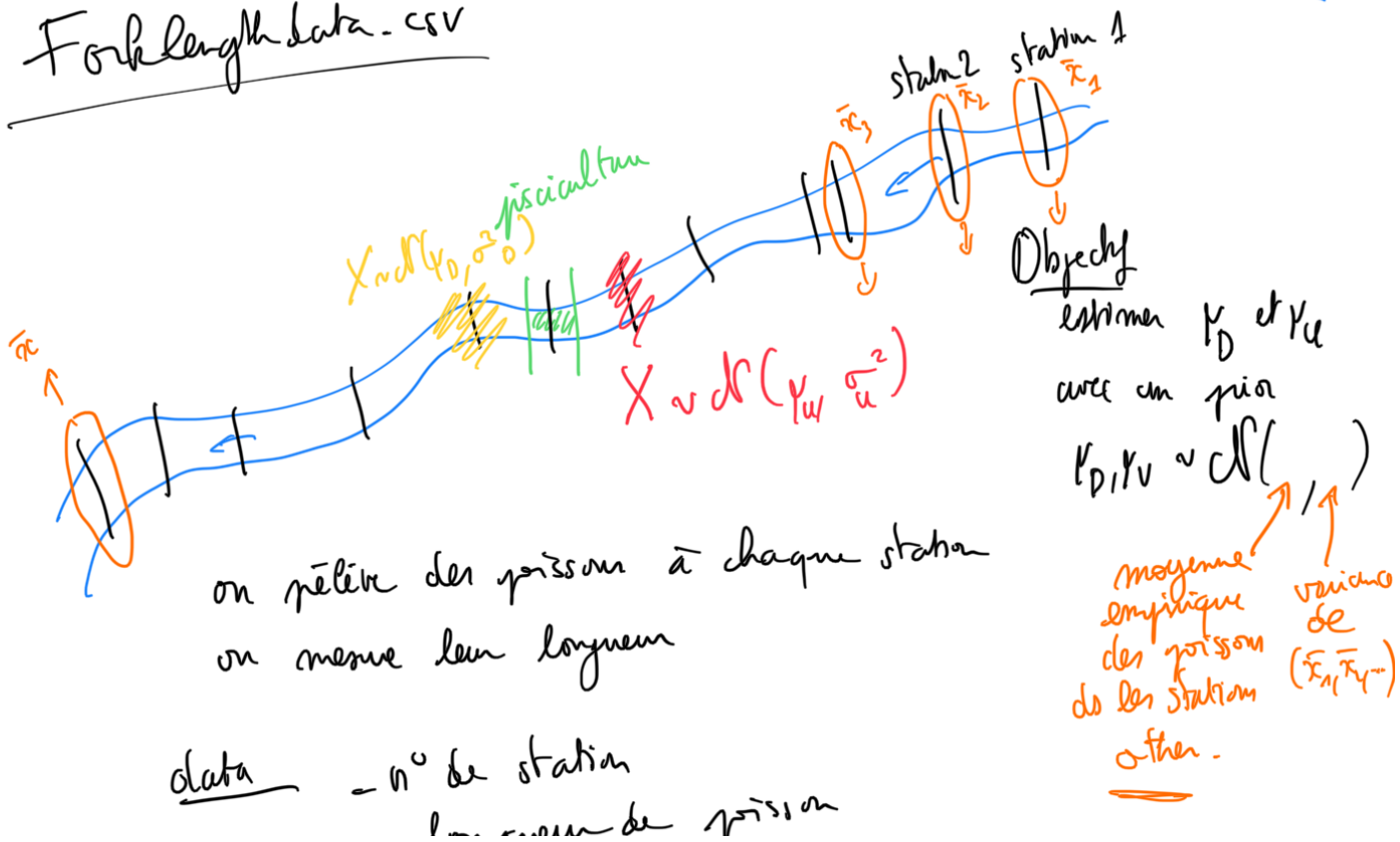


Vraisemblance



$$\underbrace{f(x_1) \times \dots \times f(x_n)}_{\text{vraisemblance}} > f(x_1) \times \dots \times f(x_n)$$

Fork length data - csv



- non given
- info sur la station: location
 - "other"
 - by the farm
 - downstream
 - upstream

En appliquant la formule pour la loi a posteriori, on obtient

$$\begin{cases} X_D \sim \mathcal{N}(\mu_D, \sigma_D^2) \\ X_U \sim \mathcal{N}(\mu_U, \sigma_U^2) \\ \mu_D, \mu_U \sim \mathcal{N}(m, s^2) \end{cases}$$

$$\begin{cases} m \approx 97 \\ s^2 \approx 70 \end{cases}$$

issu des données "other"

taille de données D

$$\mu_D | \underline{x} \sim \mathcal{N} \left(\frac{s^2 \times \bar{x}_D + m \times \frac{\sigma_D^2}{n}}{s^2 + \frac{\sigma_D^2}{n}}, \frac{s^2 \frac{\sigma_D^2}{n}}{s^2 + \frac{\sigma_D^2}{n}} \right)$$

Calcul de la loi a posteriori

Théorème de Bayes

$$p(\theta | \underline{x}) =$$

$$\frac{p(\underline{x} | \theta) p(\theta)}{p(\underline{x})}$$

$$p(\underline{x} | \theta) p(\theta)$$

$p(\underline{x}) = ?$ = proba des données "quelque soit" la valeur de θ

$$\begin{cases} p(\underline{x}|\theta) \\ p(\underline{x}, \theta) = p(\underline{x}|\theta) p(\theta) \end{cases}$$

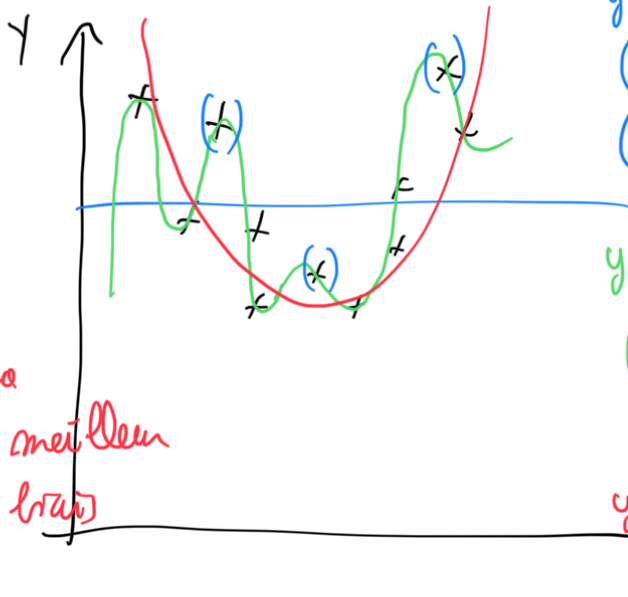
$$p(\underline{x}) = \int_{(\theta)} p(\underline{x}, \theta) d\theta$$

$$= \int_{(\theta)} \underbrace{p(\underline{x}|\theta)}_{\mathcal{L}} \underbrace{p(\theta)}_{\mathcal{P}} d\theta$$

$\underbrace{\hspace{10em}}_{\mathcal{M}}$

Comment comparer 2 modèles en statistique ?

ex: intégration



On recherche
toujours le meilleur
compromis biais
variance

$$y_i = a + bx_i + \epsilon_i$$

- 😊 • grand biais (ϵ_i trop grand)
- 😞 • faible variance : peu sensible aux fluctuations d'échantillon.

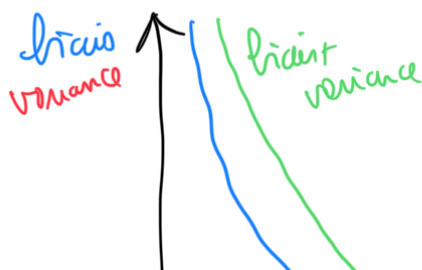
$$y_i = a + \sum_{j=1}^n b_j x_i^j + \epsilon_i$$

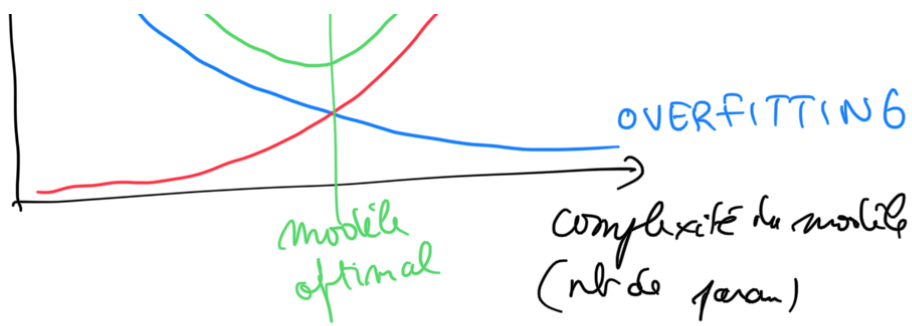
- 😊 • biais faible
- 😞 • forte variance

OVERFITTING

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

- 😊 • faible variance
- 😊 • faible biais



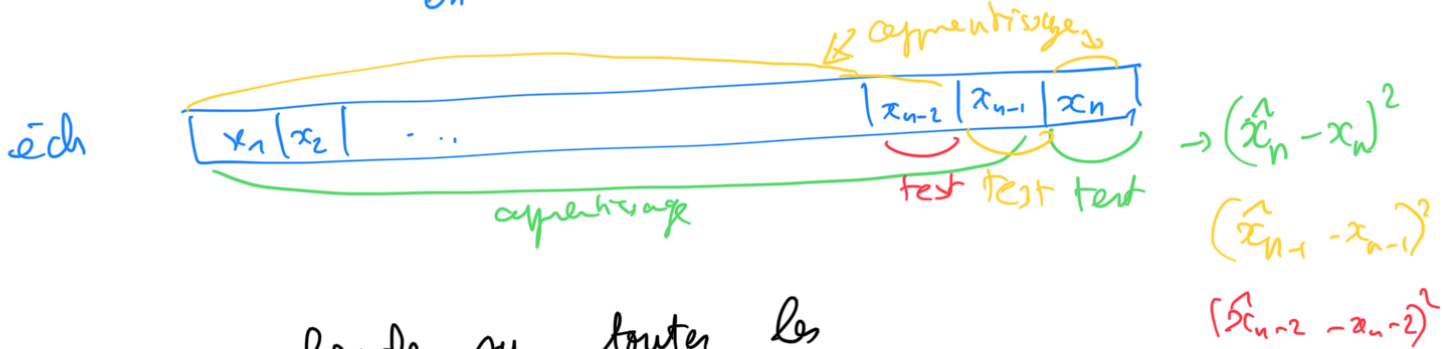


Quel critère

- juste les modèles proba universel
- vraisemblance : NON choisir toujours le modèle le plus complexe
 - somme des carrés résiduels sur l'éch d'apprentissage : NON
 - test du rapport de vraisemblance maximale : pour 2 modèles en compétition
 - vraisemblance pénalisée : AIC OUI
 - somme des carrés résiduels sur un éch test : OUI
- ex: $y_n x_1 + x_2$
 $y_n x_3$

↳ on découpe notre base de données en une partie "apprentissage" ($\frac{2}{3}$ des données) et une partie "test".
on estime le modèle sur la partie apprentissage et on l'évalue sur la partie test.

→ quand on a peu de données, l'alternative est la validation croisée :



on boucle sur toutes les données de sorte que chacune ait servi une et une seule fois de test

VALIDATION CROISÉE LEAVE-ONE-OUT

(validation croisée K fold, ...)

Modèle de régression linéaire

$$\begin{cases} y_i = \alpha + \sum_{j=1}^p \beta_j x_{ij} + \varepsilon_i \\ \text{avec } \varepsilon_i \sim \mathcal{N}(0, \sigma^2) \end{cases}$$

$$\Rightarrow y_i | x_i \sim \mathcal{N}(\alpha + \sum \beta_j x_{ij}, \sigma^2)$$

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \alpha + \sum_j \beta_j x_{1j} \\ \vdots \\ \alpha + \sum_j \beta_j x_{nj} \end{pmatrix}, \begin{pmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{pmatrix} \right)$$

$$Y \sim \mathcal{N} \left(\alpha \underbrace{\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}}_{1_n} + \underbrace{\begin{pmatrix} x_{11} - x_{1p} \\ \vdots \\ x_{n1} - x_{np} \end{pmatrix}}_X \underbrace{\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}}_{\beta}, \sigma^2 \underbrace{\begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}}_{I_n} \right)$$

$$Y \sim \mathcal{N}(\alpha 1_n + X\beta, \sigma^2 I_n)$$

Th

Si X et ε sont gaussiennes,
la meilleure* fonction $f(x)$ qui permet
d'expliquer Y est

$$f(x) = aX + b$$

* au sens $E[(Y - f(x))^2]$ soit minimale

Regression ^{linéaire} en grande dimension

- corrélation dans les $X \Rightarrow V(\hat{\beta})$ très grande \Rightarrow non significative
- $n \approx p$ (ou $n < p$) \Rightarrow "fausse" corrélation qui disparaît

Solutions

① algorithme de sélection de variables de type backward / forward / stepwise (step en R)
 marche bien pour quelques dizaines de variables ($\leq 20, 30, \dots$)

② construire de nouvelles "méta" variables qui résument l'information contenue dans X , tout en étant non corrélées entre elles
 \rightarrow régression en composantes principales
 \rightarrow régression PLS
 Partial Least Square

efficace en prédiction, mais on perd en interprétabilité du modèle (pcr/pls en R)

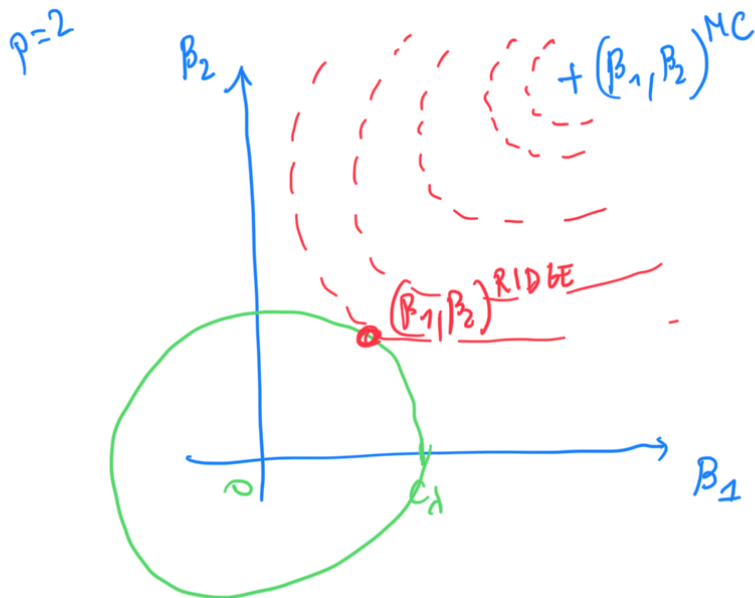
③ méthodes de régression pénalisée

$$Y = X\beta + \epsilon$$

$\hat{\beta}$ maximum de vraisemblance \Leftrightarrow moindres carrés

$$\hat{\beta} \text{ minimise } \sum_i (y_i - (\beta_0 + \sum_j \beta_j x_{ij}))^2$$

RIDGE on cherche $\hat{\beta}$ qui minimise $\sum_i (y_i - (\dots))^2$ sous la contrainte $\sum_j \beta_j^2 \leq C_\lambda$
 $\Leftrightarrow \sum_i (y_i - (\beta_0 + \sum_j \beta_j x_{ij}))^2 + \lambda \sum_j \beta_j^2$



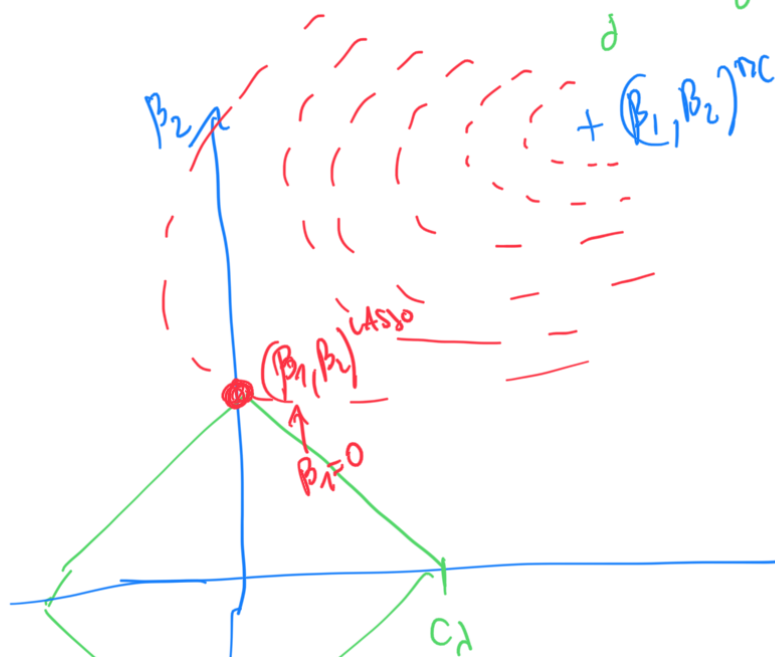
\Leftrightarrow modèle bayésien avec $\beta \sim \mathcal{N}(0, \frac{\sigma^2}{\lambda} I)$

$$\hat{\beta} = (X^T X + \lambda I)^{-1} X^T Y$$

LASSO on cherche $\hat{\beta}$ minimise

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

$$\sum_j |\beta_j| \leq C_\lambda$$



\Rightarrow les β_j des variables "non significatives" sont mis à 0.

meilleure alternative pour la
régularisation en dimension

Rq c'est une régression bayésienne avec

$\beta \sim \text{Laplace}(\cdot)$

