

# APPRENTISSAGE SUPERVISE

## Modèle linéaire gaussien & Introduction au modèle linéaire généralisé

HPC – BIG DATA  
Mars 2023

# Introduction

---

1. Les statistiques dans le cursus HPC - Big Data
2. Intérêt de la modélisation statistique
3. Exemple de la prévision opérationnelle du temps

# Introduction / Apprentissage statistique

## 1. Apprentissage supervisé :

On cherche à expliquer une variable  $Y$  (prédicte) au moyen d'autres variables (prédicteurs), 2 situations possibles :

- $Y$  est quantitative → modèle de régression
- $Y$  est qualitative → modèle de discrimination

## 2. Apprentissage non supervisé :

Absence de prédicte, on cherche à regrouper les données ayant des caractéristiques proches (recherche de classes) ou à analyser les données dans un sous-espace plus adapté :

→ méthodes de classification (clustering)

→ méthodes d'analyse factorielle (ACP)

## Objectifs du cours

---

1. Connaître le cadre théorique du modèle linéaire gaussien
2. Savoir exploiter le modèle linéaire gaussien
3. Maîtriser les procédures de sélection automatique des prédicteurs et de validation du modèle
4. Avoir des notions sur le modèle linéaire généralisé et en particulier la régression logistique

# Prérequis

---

1. Algèbre linéaire / Calcul matriciel
2. Probabilités / Vecteurs aléatoires gaussiens
3. Estimateurs / Maximisation de la vraisemblance
4. Tests statistiques / Utilisation de la p-value d'un test

# Plan du cours

---

1. Rappels
2. Modèle linéaire gaussien
3. Inférence
4. Prédicteurs qualitatifs
5. Validation
6. Modèle linéaire généralisé (GLM)

# Bibliographie

---

- **Statistique : la théorie et ses applications, Lejeune, Springer** (prg 1ère année)
- **Régression avec R, Cornillon & al., Springer**
- **An introduction to statistical learning, James & al., Springer** (→ moodle)

*Pour aller plus loin :*

- An introduction to generalized linear models, Dobson & Barnett, CRC Press
- An R companion to applied regression, Fox & Weisberg, Sage
- The elements of statistical learning, Hastie & al., Springer (→ moodle)

# Plan du cours

---

## **1. *Rappels***

2. Modèle linéaire gaussien

3. Inférence

4. Prédicteurs qualitatifs

5. Validation

6. Modèle linéaire généralisé (GLM)



## Résultats utiles

- Vecteur aléatoire dans  $\mathbb{R}^n$**

$Y_i$  ( $i = 1, \dots, n$ ) = variables aléatoires réelles

$\xi_i$  =  $i^{\text{ème}}$  vecteur de la base canonique

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ Y_n \end{pmatrix} = Y_1 \xi_1 + Y_2 \xi_2 + \dots + Y_n \xi_n \qquad \mathbb{E}[Y] = \mu = \begin{pmatrix} \mathbb{E}[Y_1] \\ \mathbb{E}[Y_2] \\ \mathbb{E}[Y_3] \\ \vdots \\ \mathbb{E}[Y_n] \end{pmatrix}$$

**$Y$  vecteur aléatoire de  $\mathbb{R}^n$  d'espérance  $\mathbb{E}[Y]$**

Matrice de covariance de  $Y$  :  $V[Y] = \Gamma$  = matrice  $n \times n$  symétrique, réelle,  
de terme général  $\Gamma_{ij} = \text{COV}(Y_i, Y_j)$

## Résultats utiles

Soit  $A$  matrice  $m \times n$  non aléatoire, et soit  $Z = AY$ .

$$E[Z] = A E[Y] \quad \text{et} \quad V[Z] = A V[Y] {}^tA$$

- Loi normale dans  $\mathbb{R}^n$

$Y$  vecteur gaussien ; toute combinaison linéaire de ses composantes suit une loi de Gauss à une dimension.

Si  $\Gamma = V[Y]$  est inversible,  $Y$  admet pour densité :

$$f(Y) = \frac{1}{\sqrt{(2\pi)^n \det(\Gamma)}} \exp\left(-\frac{1}{2} {}^t(Y - \mu) \Gamma^{-1} (Y - \mu)\right)$$

$A$  matrice  $m \times n$ ,  $Z = AY$ , et  $Y \sim N_n(\mu, \Gamma) \Rightarrow Z \sim N_m(A\mu, A \Gamma {}^tA)$

## Résultats utiles

- Loi normale sphérique

La loi normale **sphérique** :  $\Gamma = \sigma^2 \mathbf{I}_n$

**Composantes décorrélées (et donc indépendantes dans ce cadre gaussien)**

**Composantes de même variance  $\sigma^2$**

- métrique euclidienne :

$\langle \mathbf{X}, \mathbf{Y} \rangle = \sum_i \mathbf{X}_i \mathbf{Y}_i$  produit scalaire canonique

$\|\mathbf{Y}\|^2 = \langle \mathbf{Y}, \mathbf{Y} \rangle = \sum_i \mathbf{Y}_i^2 = {}^t\mathbf{Y}\mathbf{Y}$

- densité sphérique :

$$\frac{1}{(\sigma\sqrt{2\pi})^n} \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{Y} - \boldsymbol{\mu}\|^2\right)$$

- propriété *importante* :  $\mathbf{Y} \sim \mathbf{N}_n(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_n) \Rightarrow \|\mathbf{Y} - \boldsymbol{\mu}\|^2 / \sigma^2 \sim \chi^2_n$

## Résultats utiles

- **Projection orthogonale**

**$Q$  sous-espace de  $\mathbb{R}^n$  de dimension  $q$ .**

**$\Pi_Q Y$  projeté orthogonal de  $Y$  sur  $Q$ .**

**$\Pi_Q Y$  définie par  $\Pi_Q Y \in Q$  et  $Y - \Pi_Q Y \perp$  à tout vecteur de  $Q$ .**

**$Q$  engendré par les  $q$  vecteurs colonnes de la matrice  $T$ , avec  $\text{rg}(T)=q$   
 $\Rightarrow \Pi_Q = T ({}^t T T)^{-1} {}^t T$**

**$\Pi_Q$  est une matrice symétrique idempotente :  $\Pi_Q = {}^t \Pi_Q = \Pi_Q^2$**

**$\text{tr}(\Pi_Q) = \dim Q = q$**

# Plan du cours

---

1. Rappels

**2. *Modèle linéaire gaussien***

3. Inférence

4. Prédicteurs qualitatifs

5. Validation

6. Modèle linéaire généralisé (GLM)

## Modèle linéaire gaussien : formulation

$$\mathbf{Y} = \mathbf{T} \boldsymbol{\beta} + \mathbf{e} \quad \text{avec} \quad \mathbf{e} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

= comp. déterministe + comp. aléatoire

- $\mathbf{Y}$  : vecteur aléatoire gaussien de  $\mathbb{R}^n$ , contenant les  $n$  observations de la **variable expliquée** (= prédictand)
- $\mathbf{T}$  : matrice non aléatoire de **données explicatives**, de dimension  $n \times q$ , contenant les  $n$  observations sur les  $q$  variables explicatives (= prédicteurs)
- $\boldsymbol{\beta}$  : vecteur des  $q$  **paramètres du modèle** (à estimer,  $q$ =dimension du modèle)
- $\mathbf{e}$  : vecteur aléatoire gaussien de  $\mathbb{R}^n$  contenant les **résidus du modèle**.  
On suppose que le vecteur  $\mathbf{e}$  suit une loi normale sphérique de vecteur espérance nul et de matrice de covariance  $\sigma^2 \mathbf{I}_n$ , avec  $\sigma^2$  à estimer.

## Formulation

En notant  $x_{ij}$  la  $i^{\text{ème}}$  coordonnée du vecteur  $X_j$ ,

pour  $i=1,\dots,n$  on peut écrire :

$$Y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_q x_{iq} + e_i$$

Soit, sous forme matricielle :

$$\begin{pmatrix} Y_1 \\ \cdot \\ \cdot \\ \cdot \\ Y_n \end{pmatrix} = \begin{pmatrix} \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \\ X_1 & X_2 & \dots & X_q \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \end{pmatrix} \begin{pmatrix} \beta_1 \\ \cdot \\ \beta_q \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ \cdot \\ \cdot \\ e_n \end{pmatrix}$$

Avec  $T=(X_1|X_2|\dots|X_q) : Y = T \beta + e$

$T$  = matrice du plan d'expérience = "*design matrix*"

## Formulation

**Y est donc un vecteur aléatoire suivant une loi gaussienne sphérique, dont le vecteur espérance  $\mu$  est une fonction linéaire des covariables :**

$$\mathbf{E}[\mathbf{Y}] = \boldsymbol{\mu} = \mathbf{T}\boldsymbol{\beta} \quad \text{et} \quad \mathbf{V}[\mathbf{Y}] = \mathbf{V}[\mathbf{T}\boldsymbol{\beta} + \mathbf{e}] = \mathbf{V}[\mathbf{e}] = \sigma^2 \mathbf{I}_n$$

Les **résidus  $\mathbf{e}_i$**  ( $i=1\dots n$ ) sont des variables aléatoires *normales, centrées, indépendantes* et de *variance constante* (homoscédasticité).

$\mathbf{X}_1, \mathbf{X}_2 \dots \mathbf{X}_q$  vecteurs contenant les  $n$  observations des  $q$  variables explicatives.

$\mathbf{X}_j$  ( $j=1, \dots, q$ ) peut être le vecteur des observations d'une variable numérique (prédicteur quantitatif), un vecteur d'indicateurs (prédicteur qualitatif), un vecteur de coordonnées égales à 1 (estimation du terme constant du modèle).



# Formulation

## Remarques

### 1. **Modèle linéaire = linéarité selon $\beta$**

Les prédicteurs peuvent être eux-mêmes des fonctions non linéaires d'autres variables, de multiples prédicteurs élaborés peuvent être testés, par exemple introduire des termes d'interaction peut enrichir le modèle.

Exemple : le modèle défini par  $Y = \beta_1 \ln(X_1) + \beta_2 \exp(X_2) + \beta_3 X_1 \cdot X_2 + e$  entre dans le cadre du modèle linéaire.

### 2. Vocabulaire : **régression multiple et modèle linéaire**

Les prédicteurs sont considérés comme non aléatoires lorsque l'on raisonne en terme de modèle linéaire.

En régression, les prédicteurs sont considérés comme étant aléatoires et le conditionnement apparaît explicitement, les 2 approches étant équivalentes.

$$E[Y] \rightarrow E[Y|X_1, \dots, X_q]$$

$$V[Y] \rightarrow V[Y|X_1, \dots, X_q]$$

## Estimation dans le modèle linéaire gaussien

- **Estimation de  $\beta$  par la méthode MCO (= Moindres Carrés Ordinaires)**

$\hat{\beta}$  = estimation de  $\beta$

$\varepsilon$  = vecteur des résidus estimés :  $\varepsilon = Y - T \hat{\beta}$

On cherche à **minimiser**  $\sum e_i^2 = \|e\|^2 = \|Y - T\beta\|^2$

→  $\hat{\beta} = \operatorname{argmin} ( \|Y - T\beta\|^2 )$

$$\hat{\beta} \text{ tel que : } \left( \frac{\partial \|Y - T\beta\|^2}{\partial \beta} \right)_{\beta=\hat{\beta}} = 0$$

## Estimation

Par dérivation vectorielle, il vient :  $-2^tT(Y - T\beta) = 0$  pour  $\beta = \hat{\beta}$

D'où  $\hat{\beta}$  estimateur des moindres carrés de  $\beta$ , avec  $T$  de rang  $q$  :

$$\hat{\beta}_{MCO} = \left( {}^tT T \right)^{-1} {}^tT Y$$

Propriété :  $\hat{\beta}$  est un **estimateur sans biais** de  $\beta$  :  $E[\hat{\beta}] = \beta$

### Théorème de Gauss-Markov :

$\hat{\beta}$  est de tous les estimateurs linéaires sans biais de  $\beta$  (de la forme  $AY$ ) celui de **variance minimale**.

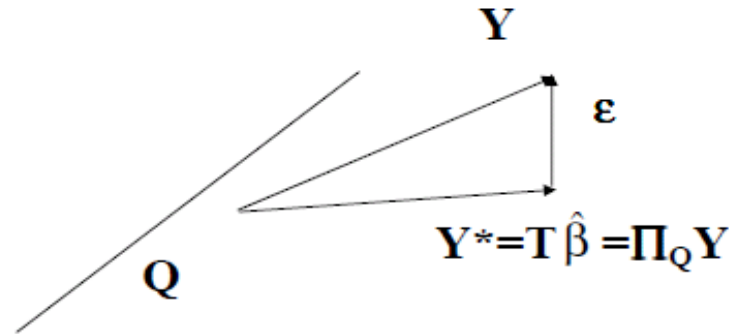
→ BLUE : Best Linear Unbiased Estimator

## Estimation

- Interprétation géométrique**

Soit  $Q$  le plan de  $\mathbb{R}^n$  engendré par les vecteurs colonnes de la matrice  $T$ .  
Le projecteur  $\Pi_Q$  sur  $Q$  s'écrit :  $\Pi_Q = T ({}^t T T)^{-1} {}^t T$

On note  $Y^*$  le vecteur des valeurs ajustées de  $Y$  par le modèle  
(valeurs ajustées = prévisions calculées sur l'archive d'apprentissage) :  
 **$Y^*$  est la projection orthogonale de  $Y$  sur  $Q$**



$$Y^* = T \hat{\beta} = T ({}^t T T)^{-1} {}^t T Y = \Pi_Q Y$$

$$Y - Y^* = \epsilon = Y - T \hat{\beta} = \text{vecteur des résidus estimés du modèle.}$$

## Estimation

- **Loi de  $\hat{\beta}$**

**Y étant un vecteur aléatoire,  $\hat{\beta}$  fonction de Y est aussi un vecteur aléatoire.**

**Dans le cadre du modèle linéaire gaussien,  $\hat{\beta}$  transformé linéaire d'un vecteur gaussien est lui-même gaussien et on a :**

$$\hat{\beta} \sim N_q \left( \beta, \sigma^2 \left( {}^t T T \right)^{-1} \right)$$

**Exercice :** en utilisant  $\hat{\beta} = ({}^t T T)^{-1} {}^t T Y$  et le fait que  $({}^t T T)^{-1}$  est une matrice symétrique, retrouver les expressions de  $E[\hat{\beta}]$  et  $V[\hat{\beta}]$ .

## Estimation

- Le vecteur  $\hat{\beta}$  suit une loi normale multidimensionnelle de dimension  $q$ , d'espérance  $\beta$  et de matrice de covariance  $\sigma^2({}^t\mathbf{T}\mathbf{T})^{-1}$
- Cette matrice de variance-covariance n'est en général pas diagonale : les estimations  $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_q$  ne sont alors pas indépendantes.

Soit  $\hat{\beta}_k$  ( $k=1, \dots, q$ ) la  $k^{\text{ième}}$  coordonnée de  $\hat{\beta}$

$\hat{\beta}_k$  suit donc une loi normale :

- d'espérance  $\beta_k$
- d'écart-type  $\sigma_k = \sigma \sqrt{({}^t\mathbf{T}\mathbf{T})_{kk}^{-1}}$

avec  $({}^t\mathbf{T}\mathbf{T})_{kk}^{-1}$  :  $k^{\text{ième}}$  élément diagonal de la matrice  $({}^t\mathbf{T}\mathbf{T})^{-1}$

# Estimation

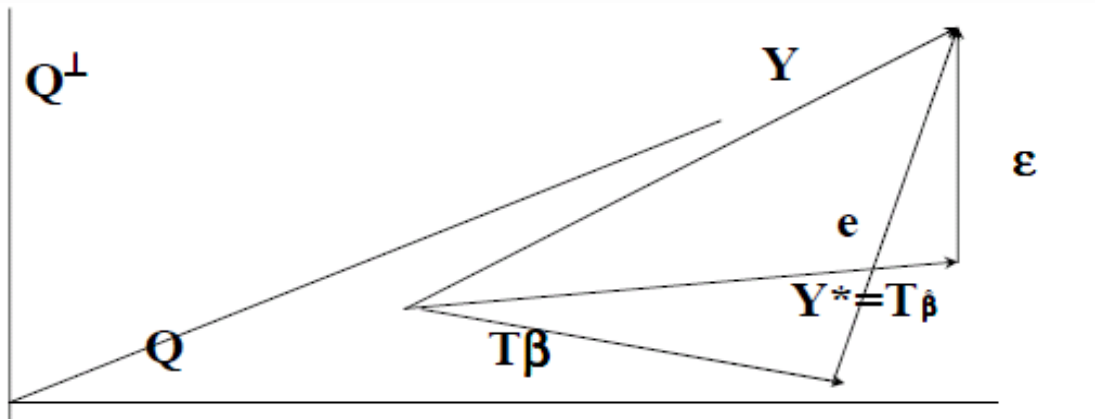
## • Estimation de $\sigma^2$

Soient  $\varepsilon$  le vecteur des résidus estimés ( $\varepsilon = Y - T \hat{\beta}$ ) et  $e$  le vecteur des résidus théorique ( $e = Y - T\beta$ ).

$\varepsilon = Y - Y^* = Y - T \hat{\beta}$  est orthogonal à  $Q$  :  $\varepsilon \in Q^\perp$

$\Pi_Q = T ({}^t T T)^{-1} {}^t T$  est le projecteur sur  $Q$

$\Pi_{Q^\perp} = I - \Pi_Q$  est le projecteur sur  $Q^\perp$  ( $I$  désignant la matrice identité).



## Formulation

On montre que :  $E[ \|Y - T \hat{\beta}\|^2 ] = \sigma^2 (n - q)$

On en déduit donc un estimateur de  $\sigma^2$  :

$$\hat{\sigma}^2 = \frac{\|Y - T\hat{\beta}\|^2}{n - q} = \frac{\|Y - Y^*\|^2}{n - q} = \frac{\|\varepsilon\|^2}{n - q} = \frac{\sum_{i=1}^n (Y_i - Y_i^*)^2}{n - q}$$

On montre que  $\hat{\sigma}^2$  est l'estimateur sans biais de variance minimale de  $\sigma^2$ .

$\varepsilon$  = projection orthogonale de  $e$  (et de  $Y$ ) sur  $Q^\perp$

$$\frac{\|Y - T\hat{\beta}\|^2}{\sigma^2} \sim \chi_{n-q}^2 \quad (\text{Th. De Cochran})$$



# Exercice

1. A partir de l'expression de la densité normale multidimensionnelle,

$$f(Y) = \frac{1}{\sqrt{(2\pi)^n \det(\Gamma)}} \exp\left(-\frac{1}{2} {}^t(Y - \mu) \Gamma^{-1} (Y - \mu)\right)$$

donner la fonction de vraisemblance associée à un  $n$  vecteur aléatoire  $Y$  dans le cadre du modèle linéaire gaussien.

2. Calculer les estimateurs  $\hat{\beta}_{MV}$  et  $\hat{\sigma}_{MV}^2$  obtenus par la méthode de maximisation de la vraisemblance.  
Montrer en particulier que :  $\hat{\beta}_{MV} = \hat{\beta}_{MCO}$
3. L'estimateur  $\hat{\sigma}_{MV}^2$  est-il biaisé ?

# Plan du cours

---

1. Rappels
2. Modèle linéaire gaussien
- 3. *Inférence***
4. Prédicteurs qualitatifs
5. Validation
6. Modèle linéaire généralisé (GLM)

## Tests dans le modèle linéaire gaussien

- **Test du caractère significatif de l'une des composantes de  $\beta$**

$$\hat{\beta}_k \sim N(\beta_k, \sigma_k) \text{ avec } \sigma_k = \sigma \sqrt{(\mathbf{t} \mathbf{T} \mathbf{T})_{kk}^{-1}}$$

$$\frac{(\hat{\beta}_k - \beta_k)}{\hat{\sigma}_k} = \frac{(\hat{\beta}_k - \beta_k)}{\sqrt{\frac{\|Y - T\hat{\beta}\|^2}{n - q} \cdot (\mathbf{t} \mathbf{T} \mathbf{T})_{kk}^{-1}}} \sim t_{n-q}$$

$t_{n-q}$  désignant la variable de Student à  $n-q$  degrés de liberté.

On peut alors en déduire l'intervalle de confiance de  $\beta_k$  ainsi que le test du caractère significatif de  $\beta_k$ .

## Prévision d'une nouvelle valeur du prédicteur

Soit  $\mathbf{t}_k$  une nouvelle observation des  $q$  prédicteurs :  $\mathbf{t}_k = (\mathbf{x}_{k1}, \dots, \mathbf{x}_{kq})$

$$Y_k^* = \hat{\beta}_1 \mathbf{x}_{k1} + \hat{\beta}_2 \mathbf{x}_{k2} + \hat{\beta}_3 \mathbf{x}_{k3} + \dots + \hat{\beta}_q \mathbf{x}_{kq} = \mathbf{t}_k \hat{\beta}$$

$Y_k^*$  transformé linéaire de  $\hat{\beta} \Rightarrow Y_k^* \sim N(\mathbf{t}_k \beta, \sigma^2 \mathbf{t}_k (\mathbf{T}\mathbf{T})^{-1} \mathbf{t}_k)$

Par studentisation, puisque  $\sigma^2$  doit être estimé, et  $Y_k$  aléatoire, on trouve :

$$\frac{Y_k - Y_k^*}{\hat{\sigma} \sqrt{1 + \mathbf{t}_k (\mathbf{T}\mathbf{T})^{-1} \mathbf{t}_k}} \sim t_{n-q}$$

$\Rightarrow$  intervalle de prévision de niveau  $(1 - \alpha)$  sur  $Y_k$  :

$$Y_k^* - t_{n-q, 1-\alpha/2} \hat{\sigma} \sqrt{1 + \mathbf{t}_k (\mathbf{T}\mathbf{T})^{-1} \mathbf{t}_k} < Y_k < Y_k^* + t_{n-q, 1-\alpha/2} \hat{\sigma} \sqrt{1 + \mathbf{t}_k (\mathbf{T}\mathbf{T})^{-1} \mathbf{t}_k}$$

avec  $t_{n-q, 1-\alpha/2}$  = quantile d'ordre  $1-\alpha/2$  de la loi de Student à  $n-q$  degrés de liberté

## Coefficient de corrélation multiple R

**R** est le coefficient de corrélation entre valeurs observées et ajustées par le modèle statistique, c'est à dire entre les vecteurs  $Y$  et  $Y^*$ .

Son carré  $R^2$ , appelé **coefficient de détermination**, mesure la qualité du modèle et s'interprète en termes de variance expliquée :

$$R^2 = \frac{\text{variance expliquée}}{\text{variance totale de } Y} = \frac{SCE}{SCT} = 1 - \frac{SCR}{SCT}$$

Formule de décomposition de la variance :

$$\frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - Y_i^*)^2 + \frac{1}{n} \sum_{i=1}^n (Y_i^* - \bar{Y})^2$$

**Variance totale (SCT) = Variance résiduelle (SCR) + Variance expliquée (SCE)**

## Test de significativité du $R^2$ (test de Fisher global)

Soit  $H_0 : R^2 = 0$  (hypothèse de non-régression)  
 $H' : R^2 > 0$

$$F = \frac{R^2}{1 - R^2} \cdot \frac{n - q}{q - 1} = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - Y_i^*)^2} \cdot \frac{n - q}{q - 1}$$

Sous  $H_0$ ,  $F$  suit une loi de Fischer-Snédecour à  $q-1$  et  $n-q$  degrés de liberté.

**Au niveau  $\alpha$ , on rejette  $H_0$  si  $F > f_{q-1, n-q, \alpha}$  (p-value  $< \alpha$ )**

Ce test compare le modèle complet au modèle ne comportant que le terme constant et vérifiant  $E[Y] = \beta_0$

## Sélection de variables : test de Fisher

### Principe de parcimonie :

A qualité d'information égale, on préfère toujours un modèle simple à un modèle plus compliqué.

#### - Test de Fisher :

**$R^2$  augmente « mécaniquement » avec la dimension du modèle**, même si l'information apportée n'est pas pertinente. En effet  $R^2$  ne tient pas compte de la dimension du sous espace de projection  $Q$ .

→ On teste donc l'accroissement de  $R^2$  :

Modèle complet à  $k+p$  prédicteurs, et modèle restreint à  $k$  prédicteurs

$$H_0 : R_k^2 = R_{k+p}^2 \quad F = \frac{n - k - p}{p} \cdot \frac{R_{k+p}^2 - R_k^2}{1 - R_{k+p}^2}$$

$$H' : R_k^2 < R_{k+p}^2$$

Sous  $H_0$ , on montre que  $F$  suit une loi de Fisher à  $p$  et  $n-k-p$  degrés de liberté.

**On rejettera donc  $H_0$  au niveau  $\alpha$  si  $F > f_{p,n-k-p,\alpha}$  (  $p$ -value  $< \alpha$  )**

Remarque : autres critères de sélection courants :

AIC, BIC,  $C_p$  de Mallows,  $R^2$  ajusté (cf.TP)

# Sélection de variables

## Algorithmes de sélection automatique des prédicteurs :

**Procédure optimale :**

**Comparaison de toutes les combinaisons possibles de prédicteurs.**

**Mais il faut alors estimer  $2^q$  modèles,**

**...soit près de  $10^{30}$  modèles avec une centaine de prédicteurs potentiels.**

- **Sélection descendante** : on retire du modèle complet les variables inutiles une à une, avec une règle d'arrêt (baisse significative du  $R^2$ , AIC minimal...)
- **Sélection ascendante** : on introduit dans le modèle à chaque étape la variable la plus utile, avec une règle d'arrêt (accroissement de  $R^2$  non significatif par exemple) : à privilégier si nb important de prédicteurs.
- **Sélection avec remise en cause** : on introduit dans le modèle à chaque étape une variable supplémentaire, et on teste si toutes les variables sont utiles, avec une règle d'arrêt.



# Plan du cours

---

1. Rappels
2. Modèle linéaire gaussien
3. Inférence
- 4. *Prédicteurs qualitatifs***
5. Validation
6. Modèle linéaire généralisé (GLM)

## Prédicteurs qualitatifs : ANOVA

L'analyse de variance (ANOVA en anglais) est une technique qui permet de **quantifier l'effet de variables *qualitatives* sur une variable quantitative**.

- Analyse de variance à un facteur

Soit  $Y$  la variable quantitative à expliquer. On considère ici l'effet moyen d'une variable qualitative (appelée facteur) à  $p$  modalités.

Le modèle s'écrit :  $Y_{ik} = \mu + \alpha_i + e_{ik}$

Avec :  $i=1,\dots,p$   $k=1,\dots,n_i$   $\sum_{i=1}^p n_i = n$   $e \sim N(0, \sigma^2 I_n)$

$\mu$  est la valeur moyenne de  $Y$ ,  $i$  la  $i^{\text{ème}}$  modalité du facteur  $\alpha$ ,  $\alpha_i$  l'effet moyen de la  $i^{\text{ème}}$  modalité de  $\alpha$ , et  $Y_{ik}$  la  $k^{\text{ème}}$  observation de  $Y$  dans la catégorie  $i$ .

## Prédicteurs qualitatifs : ANOVA

**Modèle linéaire où les variables prédictors sont remplacées par le tableau d'indicatrices des modalités du facteur  $\alpha$ .**

$$\begin{pmatrix} Y_{11} \\ Y_{12} \\ Y_{13} \\ Y_{21} \\ Y_{22} \\ Y_{23} \\ Y_{31} \\ Y_{32} \end{pmatrix} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix} + \begin{pmatrix} e_{11} \\ e_{12} \\ e_{13} \\ e_{21} \\ e_{22} \\ e_{23} \\ e_{31} \\ e_{32} \end{pmatrix}$$

**PROBLEME : T n'est pas de plein rang.**

## Prédicteurs qualitatifs : ANOVA

→ On pose alors une contrainte d'identification :

$$\sum_{i=1}^p \alpha_i = 0 \quad \Leftrightarrow \quad \alpha_p = -\sum_{i=1}^{p-1} \alpha_i$$

(Remarque : d'autres contraintes peuvent être utilisées, le choix de la contrainte n'impacte pas le modèle lui-même mais l'interprétation des paramètres finalement estimés → cf. TD & TP)

Sous cette dernière contrainte, le modèle alors identifiable s'écrit :

$$\begin{pmatrix} Y_{11} \\ Y_{12} \\ Y_{13} \\ Y_{21} \\ Y_{22} \\ Y_{23} \\ Y_{31} \\ Y_{32} \end{pmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & -1 & -1 \\ 1 & -1 & -1 \end{bmatrix} \cdot \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \end{pmatrix} + \begin{pmatrix} e_{11} \\ e_{12} \\ e_{13} \\ e_{21} \\ e_{22} \\ e_{23} \\ e_{31} \\ e_{32} \end{pmatrix}$$

## Prédicteurs qualitatifs : ANCOVA

- **Analyse de covariance (ANCOVA)**

**Analyse de covariance = mélange de prédicteurs qualitatifs et quantitatifs**

- Pas de difficultés particulières
- Une contrainte sera alors imposée pour chaque prédicteur qualitatif
- Par exemple, si la variable  $Y$  dépend des variables  $X_1$ ,  $X_2$  et des facteurs  $\alpha$  et  $v$ , on peut écrire :

$$Y_{ijk} = \mu + \beta_1 X_{1ijk} + \beta_2 X_{2ijk} + \alpha_i + v_j + e_{ijk}$$

# Plan du cours

---

1. Rappels
2. Modèle linéaire gaussien
3. Inférence
4. Prédicteurs qualitatifs
- 5. *Validation***
6. Modèle linéaire généralisé (GLM)

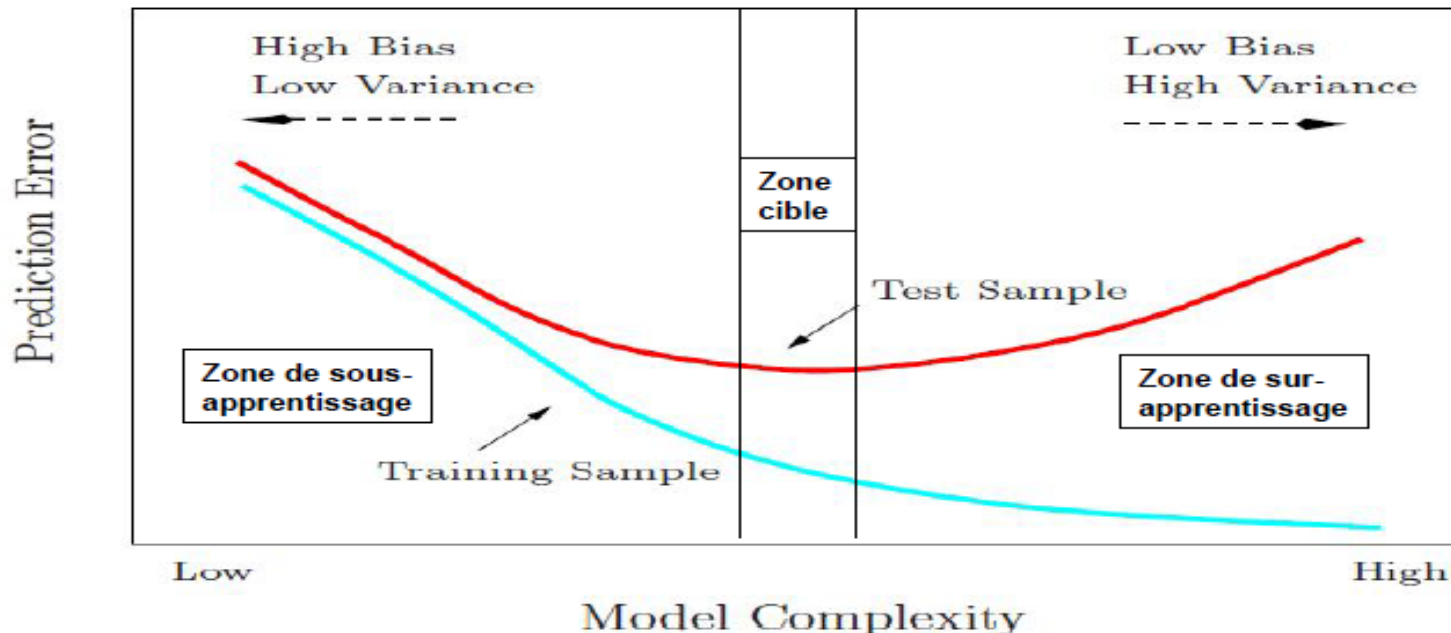
# Diagnostics et validation

## Généralités :

- L'apprentissage du modèle se fera sur un sous-échantillon dédié, une partie des données ayant été préalablement isolée de l'archive afin d'évaluer les performances du modèle en conditions opérationnelles (phase de test).
- Un modèle trop complexe manquera de robustesse (incapacité à généraliser à de nouvelles données), on parle alors de sur-apprentissage (cf. TP).

## Illustration du phénomène de sur-apprentissage

Graphé extrait de l'ouvrage The elements of statistical learning, Hastie & al.



## Diagnostics et validation

- **Multicolinéarité et calcul explicite de  $({}^t\text{TT})^{-1}$**

Si  ${}^t\text{TT}$  est mal conditionnée, son déterminant est proche de zéro.



variance élevée des estimateurs

→ estimateurs instables, conclusions peu fiables

Diagnostic : **facteur d'inflation de la variance (VIF)** :  $V_j = \frac{1}{1 - R_j^2}$

$R_j^2$  coefficient de détermination multiple de la régression de la  $j^{\text{ème}}$  variable avec toutes les autres. De fortes valeurs du VIF permettent d'identifier et de retirer les variables posant problème ( $VIF > 10$ ), on peut également combiner ces variables pour créer un prédicteur élaboré.

### **Régression en composantes principales :**

Permet de résoudre les problèmes de multicolinéarité lorsque l'on tient à conserver toutes les variables.



# Diagnostics et validation

- Etude des résidus

**L'étude des résidus après estimation est *FONDAMENTALE*.**

**Elle permet de *vérifier les hypothèses de base* du modèle linéaire, de repérer les données pathologiques, et de mettre en évidence les inadéquations au modèle.**

**Avant estimation : analyses uni et bivariées, afin d'identifier des problèmes sur les distributions de chacune des variables (dissymétrie, valeurs atypiques) ou sur les liaisons des variables prises deux par deux (non linéarité de la liaison → création d'un prédicteur élaboré plus adapté aux données).**

**Après estimation : détection des violations d'hypothèses (linéarité, indépendance, homoscedasticité) ou de points influents dans le contexte multidimensionnel (cf. TP).**

# Diagnostics et validation

## 1. Normalité non respectée

- estimations de  $\beta$  et  $\sigma^2$  restent sans biais
- $\hat{\beta}$  n'est plus exactement normal, ni optimal
- les tests de Student et Fisher sont biaisés

**Les théorèmes centraux-limites généralisés relativisent l'influence de la non-normalité des données :**

**Echantillon suffisamment grand + loi des données pas trop dissymétrique  
→ propriétés de normalité de  $\hat{\beta}$  et exactitude des tests à peu près conservées.**

### **Vérification :**

**diagrammes quantiles/quantiles, tests de Kolmogorov-Smirnov, Shapiro...**

# Diagnostics et validation

## 2. Non-indépendance

- Les estimateurs restent non biaisés mais ne sont plus optimaux
- Concerne les données issues de séries chronologiques
- Décorrélation totale : souvent difficile à obtenir

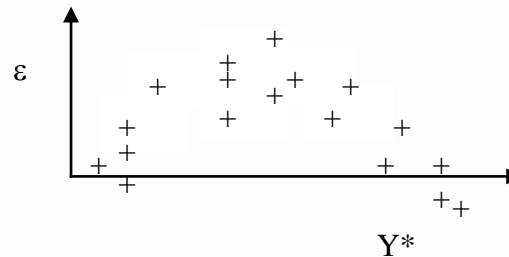
**Tant que la corrélation temporelle reste faible, ce n'est pas très gênant.**

**En cas de problèmes importants :  
modèle autorégressif ou estimation par Moindres Carrés Généralisés.**

**Vérification : corrélogramme des résidus, test de Durbin-Watson**

# Diagnostics et validation

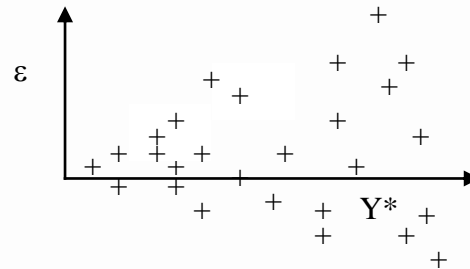
## 3. Inadéquation



Le graphe des résidus en fonction de  $Y^*$  permet de vérifier les propriétés d'adéquation et d'homoscédasticité du modèle.

La relation entre  $Y$  et  $Y^*$  n'est pas linéaire. L'analyse du problème doit être améliorée, en proposant d'autres prédicteurs ou en les transformant, ce qui peut être fait sans précautions particulières.

## 4. Hétéroscédasticité



La structure de la variance est la pathologie la plus problématique. Seules les propriétés de non biais et de convergence de  $\hat{\beta}$  sont conservées.

Nécessite un **changement de variable**, ou l'emploi de techniques de modèle linéaire généralisé.

(Vérification : graphique ou tests de Bartlett, Levene)

# Diagnostics et validation

## 5. Observations influentes

Certaines observations peuvent avoir une plus grande influence que d'autres, par **effet levier**.

On utilise généralement la **distance de Cook** pour détecter de telles observations (lorsque  $D_i > 1$ , la  $i$ ème observation est considérée comme pathologique). Une distance entre  $\hat{\beta}$  et  $\hat{\beta}_i$  est calculée,  $\hat{\beta}_i$  étant le vecteur estimé en retirant l'observation  $i$ .

Ces observations (points leviers ou valeurs aberrantes) influencent grandement les estimations et doivent donc être retirées.

Cette distance est disponible dans la plupart des logiciels statistiques.

# Plan du cours

---

1. Rappels
2. Modèle linéaire gaussien
3. Inférence
4. Prédicteurs qualitatifs
5. Validation
- 6. *Modèle linéaire généralisé (GLM)***

## Modèle linéaire généralisé et fonction de lien

- **Modèle linéaire gaussien :**

L'espérance du prédicteur  $Y$  est modélisée linéairement et

conditionnellement aux  $p$  prédicteurs : 
$$E[Y / X_1, \dots, X_p] = \beta_0 + \sum_{j=1}^p \beta_j X_j$$

- **Modèle linéaire généralisé :** (généralisation à des prédicteurs non gaussiens)

Une fonction de l'espérance du prédicteur  $Y$  est modélisée linéairement

et conditionnellement aux  $p$  prédicteurs : 
$$g(E[Y / X_1, \dots, X_p]) = \beta_0 + \sum_{j=1}^p \beta_j X_j$$

La fonction  $g$ , nommée fonction de lien canonique, dépend de la loi suivie par le prédicteur. La théorie du modèle linéaire généralisé montre qu'il existe une fonction de lien pour chaque loi de la famille exponentielle.

*Remarque : le modèle linéaire gaussien est donc un cas particulier de modèle linéaire généralisé pour lequel la fonction de lien est la fonction identité.*

## Famille exponentielle

- La famille exponentielle regroupe de nombreuses loi courantes : Normale, Exponentielle, Gamma, Bêta, Poisson, Bernoulli, Binomiale...
- La loi de  $Y$  appartient à la famille exponentielle si sa densité peut s'exprimer sous la forme :

$$f(Y) = \exp \left[ \frac{Y\theta - b(\theta)}{a(\phi)} + c(Y, \phi) \right]$$

avec  $a$ ,  $b$ ,  $c$  trois fonctions arbitraires,

$\theta$  paramètre naturel (ou canonique) de la loi, fonction de  $E[Y]$  :

$\theta = g(E[Y])$  avec  $g$  = fonction de lien canonique,

$\phi$  paramètre d'échelle (absent si la loi ne comporte qu'un paramètre).



## Modèle généralisé particulier : la régression logistique

Dans le cas de la régression logistique, le predictand est une variable binaire  $Y$  qui suit une loi de Bernoulli de paramètre :  $P(Y=1)=E[Y]$

La fonction de lien canonique associée à la loi de Bernoulli est la fonction logit :  $g(x) = \ln(x/(1-x))$  ( $\rightarrow$  résultat à montrer en exercice)

Un modèle de régression logistique aura donc pour expression :

$$g(E[Y / X_1, \dots, X_q]) = \ln \left( \frac{P(Y = 1 / X_1, \dots, X_q)}{1 - P(Y = 1 / X_1, \dots, X_q)} \right) = \beta_0 + \sum_{j=1}^q \beta_j X_j$$

L'estimation des paramètres se fera par maximisation de la vraisemblance, on obtiendra ensuite une estimation de l'espérance conditionnelle ainsi (fonction logistique = fct réciproque de la fct logit) :

$$\hat{P}(Y = 1 / X_1, \dots, X_p) = g^{-1}(\hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j X_j)$$

$$= \frac{1}{1 + \exp \left[ -(\hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j X_j) \right]}$$

## Intérêt de la régression logistique

Un modèle de régression logistique va permettre de modéliser l'occurrence d'un phénomène (orage, brouillard...) mais également le dépassement d'un seuil d'un prédicteur quantitatif.

Un tel modèle permet donc la discrimination des deux modalités d'un prédicteur binaire : occurrence (codée 1) et non occurrence (codée 0).

Le modèle calcule des estimations de probabilités d'occurrence du phénomène étudié, il va donc falloir post-traiter ces probabilités prévues pour définir la prévision binaire (occurrence ou non occurrence prévue).

→ Un seuil de probabilité va donc devoir être exploité pour passer d'une prévision probabiliste à une prévision déterministe (cf. TP et courbe ROC)