

Modèle linéaire et mélange corré

I Cadre et hypothèses

- Variable d'intérêt $y \in \mathbb{R}$ (variable endogène)
 - Variable explicative $x \in \mathbb{R}^P$ (variable exogène) élancée
 - But: prédictre y en fonction de x
+ précisément on cherche $E(y|x)$
 - Modèle linéaire: $E(y|x) = x^* \beta$, avec $\beta_0 \in \mathbb{R}^P$
- Soit en point $\varepsilon = y - E(y|x)$
- $$y = x^* \beta + \varepsilon \text{ avec } E(\varepsilon|x) = 0$$
- Hypothèse forte de structure!
- Pb: inférence sur β_0 à partir de $(y_i = x_i^* \beta + \varepsilon_i)_{1 \leq i \leq n}$
 - Hypothèse H_ε sur les résidus

$$H_\varepsilon = \begin{cases} E(\varepsilon_i | x_i) = 0 & (1) \\ \text{Var}(\varepsilon_i) = \sigma^2 & (2) \\ \text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \text{ si } i \neq j & (3) \end{cases}$$

- (1) résidus centrés conditionnellement à x (indispensable)
- (2) variance finie + (3) déconseilable. (ε étant llnc)

- $H_\varepsilon \Rightarrow E(y_i) = x_i^* \beta$

$$\text{Var}(y_i) = \sigma^2$$

$$\text{Cov}(y_i, y_j) = 0 \text{ si } i \neq j$$

⚠ conditionnel à x_i implicite.

Pb: inférence sur β et σ^2 (les paramètres du modèle)

- Notation matricielle

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} x_1^* \\ \vdots \\ x_n^* \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\text{Observation} \quad Y = X\beta + \varepsilon$$

$$H_0 \Rightarrow E(\varepsilon) = 0 \quad \text{Cov}(\varepsilon) = \sigma^2 I_m$$

$$E(Y) = X\beta \quad \text{Cov}(Y) = \sigma^2 I_m$$

Rq: Si on note $X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix} = (x_1, \dots, x_p)$

$$X\beta = \sum x_i \beta_i = \sum \beta_i x_i$$

■ Identifiabilité: si il existe une solution unique à $X\beta = y_0$.

si l'on veut un résultat unique en β :

⇒ X est de rang $\geq p$ (le maximum possible $X \in \mathbb{R}^{m \times p}$)

Rq: X de rang $p \Rightarrow \text{Im } X$ est un sous-espace de \mathbb{R}^m de dimension p .

• si pas d'identifiabilité, pb mal posé entour de β

⇒ chgt de paramétrisation (approche normale, Ridgeless) ...

■ Matrice X a un caractère artificiel: ordre des nouvelles
ordre des observations

$$Y = X\beta + \varepsilon \quad \text{et} \quad Y = \tilde{X}\tilde{\beta} + \varepsilon$$

mt 2 paramétrisations équivalentes si l'on peut passer de β_0 à $\tilde{\beta}_0$
de manière unique.

Ex: $\tilde{X} = XZ^{-1}$ et $\tilde{\beta} = Z\beta$, avec Z matrice

Même prédition mais interprétabilité et sécén. peuvent être différents

Ex: régression affine simple

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_m \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$

Pour un échantillon $y_i = \beta'_0 + \beta_1(x_i - \bar{x}) + \varepsilon_i$

$$(\beta'_0 - \beta_1 \bar{x} = \beta_0)$$

Régression affine multiple

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i$$

$$X = \begin{pmatrix} 1 & x_{i1} & \dots & x_{ip} \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix}$$

$\Rightarrow p+1$ paramètres (+ σ^2 !)

Analyse de la variance

$$y_{ijk} = m_i + \varepsilon_{ijk} \quad \begin{cases} 1 \leq k \leq n_i \\ 1 \leq i \leq p \end{cases} \quad (n_i > 0)$$

$$X = \begin{pmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{pmatrix}_{n \times np} \quad \beta = \begin{pmatrix} m_1 \\ \vdots \\ m_p \end{pmatrix}$$

p variables

Analyse de la variance

$$y_{ijk} = a_i + b_i x_{ijk} + \varepsilon_{ijk} \quad \begin{cases} 1 \leq q \leq n_i \\ 1 \leq m_i \leq p \end{cases}$$

$$X = \begin{pmatrix} 1 & \dots & 1 & x_{11} & \dots & x_{1q} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 1 & \dots & 1 & x_{m1} & \dots & x_{mq} \end{pmatrix}$$

$$\beta = \begin{pmatrix} a_1 \\ \vdots \\ a_p \\ b_1 \\ \vdots \\ b_q \end{pmatrix}$$

2p variables

Modèle de courbe: $y_i = \sum \alpha_k f_k(x_i) + \varepsilon_i$. On appelle f_k comme

$$X = \begin{pmatrix} f_1(x_i) & \dots & f_k(x_i) & \dots & f_p(x_i) \end{pmatrix}$$

$$\beta = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_p \end{pmatrix} \quad p \text{ variables}$$

II Moindre carré (ordinaire) MCO
(least square) \hat{y}_i de rang plus (identifiabilité)

Idée: chercher $\hat{\beta}$ qui minimise la somme des erreurs quadratiques (SCR en français SSR en anglais)

$$SCR(\hat{\beta}) = \|Y - X\hat{\beta}\|^2 = \sum (y_i - x_i^* \hat{\beta})^2$$

Term d'attache aux données qui priviliege des prédictions proches des valeurs observées

$$R_y \cdot E(SCR(\beta)) = E(\|\varepsilon\|^2) = \text{Tr}(Cov(\varepsilon)) = \sigma^2 n$$

on voit qu'en remplaçant le norme 2 par une autre norme (ou une autre fonction) ($\| \cdot \|_1$ ou même $\| \cdot \|_1$ et $\| \cdot \|_2$ à la Hahn) mais aussi facile et bonne propriété théorique

Propriétés

- 1) $\nabla SCR(\hat{\beta}) = -2X^*(Y - X\hat{\beta})$

- 2) $\frac{\partial^2 SCR}{\partial \hat{\beta}^2}(\hat{\beta}) = 2X^*X$

- 3) SCR est strictement convexe.

Preuve

$$SCR(\hat{\beta}) = \|Y - X\hat{\beta}\|^2$$

$$SCR(\hat{\beta} + h) = \|(Y - X\hat{\beta}) - Xh\|^2$$

$$\begin{aligned} &= \|Y - X\hat{\beta}\|^2 - 2 \langle Xh, Y - X\hat{\beta} \rangle + \|Xh\|^2 \\ &= \|Y - X\hat{\beta}\|^2 - 2 \langle h, X^*(Y - X\hat{\beta}) \rangle + \|Xh\|^2 \end{aligned}$$

$$\Rightarrow \nabla SCR(\hat{\beta}) = -2X^*(Y - X\hat{\beta})$$

$$\Rightarrow \frac{\partial^2 SCR}{\partial \hat{\beta}^2}(\hat{\beta}) = 2X^*X$$

$$\text{or } \forall \alpha \in \mathbb{R}^p \quad \alpha^*(2X^*X)\alpha$$

$$= 2(X\alpha)^*X\alpha = 2\|X\alpha\|_2^2 \geq 0$$

et même > 0 car X de rang plus

$\Rightarrow 2X^*X$ définit positive \Rightarrow SCR strictement convexe

$$\text{Cor: } \hat{\beta} = \underset{\beta}{\operatorname{argmin}} \text{SCR}(\tilde{\beta}) = \underset{\beta}{\operatorname{argmin}} \|Y - X\beta\|^2$$

existe et est unique.

$$\text{De plus } \hat{\beta} = (X^* X)^{-1} X^* Y$$

Prouve: le premier point découle de la théorie connue

le second s'obtient en utilisant la condition d'optimalité du min

$$\nabla \text{SCR}(\hat{\beta}) = 0 = -2X^*(Y - X\hat{\beta})$$

$$\Rightarrow X^* X \hat{\beta} = X^* Y$$

or $X^* X$ définie positive $\Rightarrow X^* X$ inversible

$$\Rightarrow \hat{\beta} = (X^* X)^{-1} X^* Y$$

Propriétés Estimation de β par moindres carrés

L'estimation des moindres carrés $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \text{SCR}(\tilde{\beta})$ est obtenue

$$\text{par } \hat{\beta} = (X^* X)^{-1} X^* Y$$

$$\text{et vérifie } 1) E(\hat{\beta}) = \beta$$

$$2) \text{Var}(\hat{\beta}) = \sigma^2 (X^* X)^{-1}$$

Prouve: la première propriété est une réduction du 1^{er} résultat

$$E(\hat{\beta}) = E((X^* X)^{-1} X^* Y)$$

$$= (X^* X)^{-1} X^* E(Y) = (X^* X)^{-1} X^* X \beta$$

$$= \beta$$

$$\cdot \text{Var}(\hat{\beta}) = \text{Var}((X^* X)^{-1} X^* Y)$$

Si $\tilde{\beta} = Z\beta$ alors

par les moindres carrés

$$\hat{\beta} = Z\tilde{\beta}.$$

$$= (X^* X)^{-1} X^* \text{Var}(Y) ((X^* X)^{-1} X^*)^*$$

$$= (X^* X)^{-1} X^* \sigma^2 I_m \times (X^* X)^{-1}$$

$$= \sigma^2 (X^* X)^{-1}$$

Prop (Gauss-Markov) Pour les estimations linéaires sans biais de β ,

l'estimation des moindres carrés $\hat{\beta}$ est celle de moindre moindre

(BLUE Best Unbiased Linear Estimator)

Préuve. $\hat{\beta}^* = L Y$

$$E(\hat{\beta}^*) = \beta \Rightarrow L X \beta = \beta \quad \forall \beta$$

$$\Rightarrow L X = \text{Id}$$

$$\text{Cov}(\hat{\beta}^*) = L \text{Cov}(Y) L^* = \sigma^2 L L^*$$

$$= \sigma^2 (M + \Delta) (M + \Delta)^* \quad \text{où } M = (X^* X)^{-1} X^* \\ \text{et } \Delta = L - M$$

$$= \sigma^2 [M^* M^* + \Delta \Delta^* + M \Delta^* + \Delta M^*]$$

$$\text{or } \Delta M^* = \underbrace{(L - M)}_{0 \text{ car } L \in M \text{ no lin}} X (X^* X)^{-1}$$

$$0 \text{ car } L \in M \text{ no lin} \Rightarrow L X = M X = \text{Id}$$

$$= 0 = M \Delta^*$$

$$\text{Cov}(\hat{\beta}^*) = \sigma^2 (X^* X)^{-1} + \sigma^2 \underbrace{\Delta \Delta^*}_{\text{sd p}} \Rightarrow \text{Cov}(\hat{\beta}) \geq \text{Cov}(\hat{\beta})$$

Prediction: $\hat{Y} = X \hat{\beta}$ prédictio de $E(y|x)$ sur les x_i observes
 $\hat{y} = x^* \hat{\beta}$ prédictio de $E(y|x)$ pour un nouvel x

Prop: $E(\hat{Y}) = X \beta \quad \text{Cov}(\hat{Y}) = \sigma^2 X (X^* X)^{-1} X^*$

$$E(\hat{y}) = x^* \beta \quad \text{Var}(\hat{y}) = \sigma^2 x^* (X^* X)^{-1} x$$

Préuve: immédiat!

Parce que $\hat{y} = x^* \hat{\beta}$ est optimal pour la
norme norme les estimateurs linéaires $|E(C^* Y)| = C^* X \beta = C^* \hat{y}$
 $\text{Var}(C^* Y) = \sigma^2 C^* C \geq \sigma^2 C^* X (X^* X)^{-1} X^* C = \sigma^2 x^* (X^* X)^{-1} x = \text{Var}(\hat{y})$

Interprétation géométrique $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \|Y - X \beta\|^2$

$$\Rightarrow \hat{Y} = \underset{\substack{Y \\ = X \hat{\beta}}}{\operatorname{argmin}} \|Y - \hat{Y}\|^2$$

$$= \underset{\substack{Y \\ \in \text{Im } X}}{\operatorname{argmin}} \|Y - \hat{Y}\|^2$$

$\Rightarrow \hat{Y}$ est la projection orthogonale de Y sur $\text{Im } X$
(lien avec l'espace orthogonale L^\perp et sa norme orthogonale linéaire)

$$\hat{Y} = P_{\text{Im } X} Y \quad \text{où } P \text{ est une matrice de projection } \perp$$

$$\text{Rq } \hat{Y} = X\beta = \underbrace{X(X^*X)^{-1}X^*}_? Y$$

$H = X(X^*X)^{-1}X^*$ que Matrix = matrice d'yeux = projecteur?

$$H^* = H$$

$$HH = X(X^*X)^{-1}X^* X(X^*X)^{-1}X^* = X(X^*X)^{-1}X^* = H$$

$$H = X(X^*X)^{-1}X^* \Rightarrow \text{Im } H \subset \text{Im } X$$

$$\forall z \in \text{Im } X \exists \beta, z = X\beta \Rightarrow Hz = Hz = X(X^*X)^{-1}X^*X\beta = X\beta = z$$

\Rightarrow on retrouve par le calcul la propriété précédente.

Résidu : $\hat{\epsilon} = Y - \hat{Y}$ "estimation" des ϵ ;

$$\text{Puf: } E\begin{pmatrix} \hat{Y} \\ \hat{\epsilon} \end{pmatrix} = \begin{pmatrix} X\beta \\ 0 \end{pmatrix} \text{ et } \text{Cov}\begin{pmatrix} \hat{Y} \\ \hat{\epsilon} \end{pmatrix} = \sigma^2 \begin{pmatrix} H = P_{\text{Im } X} & 0 \\ 0 & I-H = P_{\text{Im } X}^\perp \end{pmatrix}$$

Prem . . $E(\hat{Y}) = X\beta$ a déjà été vérifié

$$E(\hat{\epsilon}) = E(Y - \hat{Y}) = E(Y) - E(\hat{Y}) = 0$$

$$\hat{Y} = HY = P_{\text{Im } X} Y$$

$$\hat{\epsilon} = Y - \hat{Y} = (I - H) Y = (I - P_{\text{Im } X}) Y = P_{\text{Im } X^\perp} Y$$

$$\begin{pmatrix} \hat{Y} \\ \hat{\epsilon} \end{pmatrix} = \begin{pmatrix} P_{\text{Im } X} \\ P_{\text{Im } X^\perp} \end{pmatrix} Y$$

$$\Rightarrow \text{Cov}\begin{pmatrix} \hat{Y} \\ \hat{\epsilon} \end{pmatrix} = \begin{pmatrix} P_{\text{Im } X} & P_{\text{Im } X}^* \\ P_{\text{Im } X^2} & \end{pmatrix} \begin{bmatrix} \sigma^2 J_m \\ 0 \end{bmatrix} \begin{pmatrix} P_{\text{Im } X}^* & P_{\text{Im } X^\perp}^* \\ P_{\text{Im } X^\perp} & P_{\text{Im } X^2} \end{pmatrix}$$

$$= \sigma^2 \begin{bmatrix} P_{\text{Im } X} P_{\text{Im } X}^* & P_{\text{Im } X} P_{\text{Im } X^\perp}^* \\ P_{\text{Im } X^2} P_{\text{Im } X}^* & P_{\text{Im } X^2} P_{\text{Im } X^\perp}^* \end{bmatrix}$$

$$= \sigma^2 \begin{bmatrix} P_{\text{Im } X} & 0 \\ 0 & P_{\text{Im } X^\perp} \end{bmatrix}$$

$$\text{Car } E\left(\begin{pmatrix} \hat{\beta} \\ \hat{\epsilon} \end{pmatrix}\right) = \begin{pmatrix} \beta \\ 0 \end{pmatrix} \text{ et } \text{Cor}\left(\begin{pmatrix} \hat{\beta} \\ \hat{\epsilon} \end{pmatrix}\right) = \begin{pmatrix} (X^*X)^{-1} & 0 \\ 0 & P_{\text{Im}X^\perp} \end{pmatrix}$$

Here: Il suffit de noter que $\hat{\beta} = (X^*X)^{-1} X^* \hat{Y}$.

$$\text{Prouv: } \hat{\sigma}^2 = \frac{\text{SCR}(\hat{\beta})}{m-p} = \frac{\|Y - X\hat{\beta}\|^2}{m-p} \text{ est un estimateur sans biais de } \sigma^2$$

$$\text{Here: } \text{SCR}(\hat{\beta}) = \|Y - X\hat{\beta}\|^2 = \|\hat{\epsilon}\|^2$$

$$\text{on connait } E(\hat{\epsilon}) = 0 \quad E(\|\hat{\epsilon}\|^2) = E \text{Tr}(\hat{\epsilon} \hat{\epsilon}^*) \\ = \text{Tr}(E(\hat{\epsilon} \hat{\epsilon}^*))$$

$$\text{et connait } E(\hat{\epsilon}) = 0 \quad E(\|\hat{\epsilon}\|^2) = \text{Tr}(\text{Cor}(\hat{\epsilon}))$$

$$\text{or } \text{Tr}(\text{Cor}(\hat{\epsilon})) = \text{Tr}(\sigma^2 P_{\text{Im}X^\perp}) \\ = \sigma^2 (m-p) \quad \text{car } \dim \text{Im}X^\perp = m-p$$

Ex: Régression linéaire simple

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

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

$$\text{avec } X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \text{ et } \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$

$$X^*X = \begin{pmatrix} m & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} = m \begin{bmatrix} 1 & \bar{x} \\ \bar{x} & \bar{x}^2 \end{bmatrix}$$

$$(X^*X)^{-1} = \frac{1}{m \bar{x}^2 - (\bar{x})^2} \begin{bmatrix} \bar{x}^2 - \bar{x}x \\ -\bar{x} & 1 \end{bmatrix} =$$

$$X^*Y = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix} = m \begin{pmatrix} \bar{y} \\ \bar{xy} \end{pmatrix}$$

$$\hat{\beta} = (X^* X)^{-1} X^* Y = \frac{1}{\bar{x}^2 - \bar{x}^2} \begin{bmatrix} \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} \begin{pmatrix} \bar{y} \\ \bar{xy} \end{pmatrix}$$

$$= \frac{1}{\bar{x}^2 - \bar{x}^2} \begin{bmatrix} \bar{x}^2 \bar{y} - \bar{x} \bar{xy} & \bar{x}^2 \bar{y} + \bar{x} (\bar{x} \bar{y}) \\ -\bar{x} \bar{y} & \bar{xy} \end{bmatrix}$$

$$\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \begin{bmatrix} \bar{y} - \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x}^2 - \bar{x}^2} \bar{x} \\ \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x}^2 - \bar{x}^2} \end{bmatrix} = \begin{bmatrix} \bar{y} - \frac{\text{cov}(x,y)}{\text{var}(x)} \bar{x} \\ \frac{\text{cov}(x,y)}{\text{var}(x)} \end{bmatrix}$$

$$\rho(x,y) = \frac{\text{cov}(x,y)}{\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}} = \frac{\text{cov}(x,y)}{\sigma(x) \sigma(y)}$$

$$\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \left(\begin{array}{c} \bar{y} - \frac{\rho(x,y) \sigma(y)}{\sigma(x)} \bar{x} \\ \frac{\rho(x,y) \sigma(y)}{\sigma(x)} \end{array} \right)$$

$$\hat{y} = \bar{y} + \frac{\rho(x,y) \sigma(y)}{\sigma(x)} (x - \bar{x}) \quad \frac{\hat{y} - \bar{y}}{\sigma(x)} = \rho(x,y) \frac{x - \bar{x}}{\sigma(x)}$$

$$\hat{Y} = \bar{y} \mathbf{1} - \frac{\rho(x,y) \sigma(y)}{\sigma(x)} (X - \bar{x} \mathbf{1})$$

$$Y - \hat{Y} = (Y - \bar{y} \mathbf{1}) - \frac{\rho(x,y) \sigma(y)}{\sigma(x)} (X - \bar{x} \mathbf{1})$$

$$\frac{1}{n} \|Y - \hat{Y}\|^2 = \frac{1}{n} \|Y - \bar{y} \mathbf{1}\|^2 - 2 \frac{\rho(x,y) \sigma(y)}{\sigma(x)} \frac{1}{n} \langle Y - \bar{y} \mathbf{1}, X - \bar{x} \mathbf{1} \rangle + \rho^2(x,y) \sigma^2(y)$$

$$\begin{aligned} &= \sigma^2(y) - 2 \frac{\rho(x,y) \text{cov}(x,y) \sigma(y)}{\sigma(x)} + \rho^2(x,y) \sigma^2(y) \frac{\sigma^2(x)}{E(\|X - \bar{x}\|^2)} \\ &= \sigma^2(y) \left(1 - \rho^2(x,y) \right) \end{aligned}$$

$$SCR = \|Y - \hat{Y}\|^2 = m \sigma^2(y) (1 - p^2(x_0))$$

$$\widehat{\sigma^2} = \frac{1}{m-2} \|Y - \hat{Y}\|^2 = \frac{m}{m-2} \sigma^2(u) (1 - p^2(x_0))$$

$$\text{Cov} \left(\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} \right) = (X^* X)^{-1}$$

$$= \frac{1}{m \text{Var}(x)} \begin{bmatrix} \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix}$$

$$\Rightarrow \text{Var}(\hat{\beta}_1) = \frac{1}{m \text{Var}(x)} \rightarrow x \text{ dispersé}$$

$$\text{Var}(\hat{\beta}_0) = \frac{1}{m} \left[1 + \frac{\bar{x}^2}{\bar{x}^2 - \bar{x}^2} \right] \rightarrow \bar{x} \text{ proche de } 0$$

By estimation non corréalis si $\bar{x} = 0$

III le modèle linéaire gaussien

on pose de $H_\varepsilon = \begin{cases} E(\varepsilon \varepsilon^*) = I \\ \text{Cov}(\varepsilon) = \sigma^2 I_m \end{cases}$

à $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$

Hypothèse blanche plus forte \Rightarrow résultats plus forts
corrélation des erreurs

Modèle: $Y \sim \mathcal{N}(X\beta, \sigma^2 I_m)$

Prop: $\hat{\beta} = (X^* X)^{-1} X^* Y, \quad \hat{Y} = X \hat{\beta} = P_{\text{Im}X} Y$

et $\hat{\varepsilon} = Y - \hat{Y} = P_{\text{Im}X^\perp} Y$

$\hat{\beta} \sim \mathcal{CN}(\beta, \sigma^2 (X^* X)^{-1}), \quad \hat{Y} \sim \mathcal{CN}(X\beta, \sigma^2 P_{\text{Im}X})$

et $\hat{\varepsilon} \sim \mathcal{CN}(0, \sigma^2 P_{\text{Im}X^\perp})$

$(\hat{\beta}, \hat{Y}) \perp\!\!\!\perp \hat{\varepsilon}$

Premier point est un rappel

le second utilise le fait que l'image par une application

linéaire d'un vecteur gaussien est un vecteur gaussien

et que dans-ci et entièrement corrélatif pour sa moyenne

et on montre que l'on a bien calculé

• le troisième point s'obtient en notant que $\begin{pmatrix} \hat{\beta} \\ \hat{\epsilon} \end{pmatrix}$ est un vecteur gaussien dont la covariance a une structure par blocs $\begin{pmatrix} P & 0 \\ 0 & Q \end{pmatrix}$

$$\text{Pois } \|\hat{\epsilon}\|^2 \sim \sigma^2 \chi^2(m-p)$$

Prouve : $\hat{\epsilon} \sim \mathcal{N}(0, P_{m \times m})$

avec $\exists B$ matrice orthogonale tel que $B P_{m \times m} B^* = \begin{pmatrix} I_{m-p} & 0 \\ 0 & 0 \end{pmatrix}$

$$\Rightarrow B\hat{\epsilon} \sim \mathcal{N}\left(0, \sigma^2 \begin{pmatrix} I_{m-p} & 0 \\ 0 & 0 \end{pmatrix}\right)$$

$$\Rightarrow (B\hat{\epsilon})_i = \sigma e_i, 1 \leq i \leq m-p \quad \text{avec } e_i \text{ ind}$$

$$\|B\hat{\epsilon}\|^2 = \sum_i (B\hat{\epsilon})_i^2 = \sum_{i=1}^{m-p} \sigma^2 |e_i|^2$$

$$\sim \sigma^2 \chi^2(m-p)$$

Q) propriété générale des gaussiennes

$$\epsilon_1 \sim \mathcal{N}(0, P) \text{ et } \epsilon_2 \sim \mathcal{N}(0, Q) \text{ avec } P \neq Q \text{ deux matrices orthogonales distinctes}$$

$$\Rightarrow \|E_1\|^2 \sim \chi^2(\text{rg } P) \quad \|E_2\|^2 \sim \chi^2(\text{rg } Q)$$

$$\text{et } \|E_1\|^2 \perp \|E_2\|^2$$

■ lien avec le maximum de vraisemblance

$$\tilde{Y} \sim \mathcal{N}(X\beta, \sigma^2 I_m)$$

$$\Rightarrow dP(\tilde{Y} - Y; \beta, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{m/2}} e^{-\frac{\|\tilde{Y} - Y - X\beta\|^2}{2\sigma^2}}$$

\Rightarrow La log vraisemblance des observations est donnée

$$L_m(\beta, \sigma^2) = -\frac{m}{2} [\log 2\pi + \log \sigma^2] - \frac{1}{2\sigma^2} \|\tilde{Y} - Y - X\beta\|^2$$

\Rightarrow Estimation du maximum de vraisemblance

$$\frac{\partial L_m}{\partial \beta} = X^*(Y - X\tilde{\beta}) \Rightarrow \hat{\beta}^{MV} = (X^*X)^{-1}X^*Y = \hat{\beta}$$

$$\frac{\partial L_m}{\partial \sigma^2} = -\frac{m}{2\sigma^2} + \frac{1}{2\sigma^4} \|\tilde{Y} - Y - X\beta\|^2$$

$$\Rightarrow \hat{\sigma}^{MV} = \sqrt{\frac{\|\tilde{Y} - Y - X\hat{\beta}\|^2}{m}} = \sqrt{\frac{m-p}{m}} \hat{\sigma}^2$$

lien avec l'estimation en espérance !

Pr^oof Parmi les estimateurs sans biais, $\hat{\beta}$ est l'estimateur de moindre moindre (BLUE Best Unbiased Estimator)

Preuve La matrice de Fisher est donnée (\approx Fisher information)

$$\begin{aligned} J(\beta) &= E\left(-\frac{\partial \ln(\rho\sigma^2)}{\partial \beta}\right)^* \frac{\partial \ln(\rho\sigma^2)}{\partial \beta} \\ &= E\left(-\frac{\partial^2 \ln(\rho\sigma^2)}{\partial \beta^2}\right) \stackrel{\text{car cela fait un minimum}}{\leq} E\left[\frac{1}{\sigma^2} X^* X\right] / \text{SSR}(\beta) \\ &= \frac{1}{\sigma^2} (X^* X) \end{aligned}$$

Le théorème de Cramer-Rao indique qu'un estimateur sans biais T de β doit satisfaire $\text{Var}(T) \geq (J(\beta))^{-1} = \sigma^2 (X^* X)^{-1}$

L'estimateur des moindres carrés atteint cette borne, il est donc optimal

Prédiction $\hat{y} = x^* \hat{\beta} \sim \mathcal{N}(x^* \beta, x^* (X X)^{-1} x)$.

IV Tests dans le modèle linéaire gaussien \rightarrow voir plus bas

V Analyse asymptotique du modèle linéaire

$$y_i = x_i^* \beta + \epsilon_i \quad 1 \leq i \leq n$$

que se passe t'il quand n grandit? $X_m = \begin{pmatrix} x_1^* \\ x_2^* \\ \vdots \\ x_m^* \end{pmatrix}$ joue un rôle important

Soit H_E Soit H_E , $\beta \hat{\beta}_m \xrightarrow{P} \beta$ dis que $\text{Tr}[(X_m^* X_m)^{-1}] \rightarrow 0$

Preuve $E(\hat{\beta}_m) = \beta$ et $\text{Var}(\hat{\beta}_m) = (X_m^* X_m)^{-1}$

$$\hat{\beta} \parallel \hat{\beta}_m - \beta \parallel^2 = \text{Tr}((X_m^* X_m)^{-1}) \rightarrow 0$$

Cor Soit H_E , $\text{Tr}((X_m^* X_m)^{-1}) \rightarrow 0 \Rightarrow \hat{\beta}_m \xrightarrow{P} \beta$

Rq $\text{Tr}((X_m^* X_m)^{-1}) \rightarrow 0 \Rightarrow (X_m^* X_m)^{-1} \rightarrow 0$ car matrice de covariance

Ex Si X_i iid de loi Z telle que $\varphi = E(Z^* Z)$ dp

$$\Rightarrow -\frac{1}{n} X_m^* X_m \rightarrow Q \quad \text{en la loi faible des grandeurs}$$

$$\Rightarrow n(X_m^* X_m)^{-1} \rightarrow Q^{-1}$$

$$\Rightarrow (X_m^* X_m)^{-1} \rightarrow 0$$

Sous H_0 + ε_i indépendant + condition (celloques)

• $\frac{1}{m} X_m^* X_m \rightarrow Q$ définit positive ("Logique")

$$H_{\text{ex}}^+ : \sum \|x_i\|^{2+\delta} |\varepsilon_i|^{2+\delta} = O(m^{1+\delta/2}) \quad (\text{"Technique"})$$

(ok si ε_i iid $E(|\varepsilon_i|^{2+\delta}) < +\infty$ et $\sum \|x_i\|^{2+\delta} \leq O(m^{1+\frac{\delta}{2}})$)

(ok si $\|x_i\| < c$ ou $E(\|x_i\|^{2+\delta}) < +\infty$)

Prm Sous H_{ex}^+ , $\sqrt{m}(\hat{\beta}_m - \beta) \xrightarrow{D} \mathcal{N}(0, \sigma^2 Q^{-1})$

Preuve

$$\bullet E(\hat{\beta}_m) = \beta \text{ et } \text{Cov}(\hat{\beta}_m) = \sigma^2 (X_m^* X_m)^{-1}$$

$$\text{dans } E(\sqrt{m}(\hat{\beta}_m - \beta)) = 0 \text{ et } \text{Cov}(\sqrt{m}(\hat{\beta}_m - \beta)) = \sigma^2 \left(\frac{1}{m} (X_m^* X_m)^{-1} \right)$$

$$\bullet \sqrt{m}(\hat{\beta}_m - \beta) = \sqrt{m} \left[(X_m^* X_m)^{-1} X_m^* Y - \beta \right]$$

$$= \sqrt{m} (X_m^* X_m)^{-1} [X_m^* Y - X_m^* X_m \beta]$$

$$= \sqrt{m} (X_m^* X_m^{-1})^T X_m^* [Y - X_m \beta]$$

$$= \sqrt{m} (X_m^* X_m)^{-1} X_m^* \varepsilon_m$$

$$= \underbrace{\left(\frac{1}{m} X_m^* Y_m \right)^{-1}}_{Q^{-1}} \underbrace{\frac{1}{\sqrt{m}} X_m^* \varepsilon_m}_{Z_m}$$

$$\begin{aligned} a^* Z_m &= \frac{1}{\sqrt{m}} a^* X_m^* \varepsilon_m = \frac{1}{\sqrt{m}} (X_m a)^* \varepsilon_m \\ &= \frac{1}{\sqrt{m}} \sum_{i=1}^m (x_i^* a) \varepsilon_i \times \frac{1}{\sqrt{m}} \end{aligned}$$

$\beta_i \perp \!\!\! \perp$

$$S_m^2 = \sum_{i=1}^m E(|\beta_i|^2)$$

$$\bullet E(\beta_i) = 0 \quad K_m = \sum_{i=1}^m E(|\beta_i|^{2+\delta}) = \sigma^2 a^* Q a$$

$$S_m^2 = \sum_{i=1}^m \left(\frac{1}{m} (x_i^* a)^2 \right) E(\varepsilon_i)^2 = \sigma^2 \frac{1}{m} a^* X_m^* X_m a \rightarrow \sigma^2 a^* Q a > 0$$

$$K_m = \sum \frac{1}{m^{1+\delta/2}} (k_i^* a)^{2+\delta} E|\varepsilon_i|^{2+\delta} = O(1) \Rightarrow \frac{K_m}{(S_m)^{2+\delta}} \rightarrow 0$$

Th de Lyapounov s'applique

$$\frac{1}{S_m} \sum_{i=1}^n \beta_i \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1)$$

$$\frac{1}{S_m} \frac{a^* z_m}{\sqrt{m}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1)$$

$$\begin{aligned} S_m^2 &\rightarrow \sigma^2 a^* Q_m a \rightarrow \sigma^2 a^* Q a \\ \Rightarrow a^* z_m &\xrightarrow{\mathcal{L}} \mathcal{N}(0, \sigma^2 a^* Q a) \\ \Rightarrow z_m &\xrightarrow{\mathcal{L}} \mathcal{N}(0, \sigma^2 Q) \end{aligned}$$

$$\text{or } \sqrt{m} (\hat{\beta}_m - \beta) = Q_m^{-1} z_m$$

$$\text{avec } Q_m^{-1} \rightarrow Q^{-1}$$

$$\Rightarrow \sqrt{m} (\hat{\beta}_m - \beta) \rightarrow \mathcal{N}(0, \sigma^2 Q^{-1})$$

P En pratique, on utilise Q_m au lieu de Q

$$\sqrt{m} Q_m^{1/2} (\hat{\beta}_m - \beta) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \sigma^2 I_p)$$

■ Convergence pour $\widehat{\sigma}^2$

P : Sous $H_{\epsilon_X}^+$, $\widehat{\sigma}^2 \xrightarrow{\mathcal{P}} \sigma^2$
+ $\epsilon_i \text{ iid}$

$$\text{Preuve : } \widehat{\sigma}^2 = \frac{\|Y - \hat{Y}\|^2}{m-p} = \frac{m}{m-p} \frac{\|(I - P_{ImX})Y\|^2}{m}$$

$$= \frac{m}{m-p} \frac{\|(I - P_{ImX})\epsilon\|^2}{m}$$

$$= \frac{m}{m-p} \frac{\epsilon^* (I - P_{ImX}) \epsilon}{m}$$

$$= \frac{m}{m-p} \left[\frac{\epsilon^* \epsilon}{m} - \frac{\epsilon^* P_{ImX} \epsilon}{m} \right]$$

\downarrow
1 $\downarrow \widehat{\sigma}^2$

(loi des grands nombres)

$$\begin{aligned}
 \frac{\mathbf{E}^* P_{Im \times E}}{m} &= \frac{\mathbf{E}^* X_m (X_m^* X_m)^{-1} X_m^* E}{m} \\
 &= \frac{1}{\sqrt{m}} (\bar{X}_m^* E)^* (X_m^* X_m)^{-1} \frac{1}{\sqrt{m}} (X_m^* E) \\
 &= Z_m (X_m^* X_m)^{-1} Z_m \\
 &= \frac{1}{m} Z_m Q_m^{-1} Z_m \\
 &= \frac{1}{m} \| \underbrace{Q_m^{-1/2}}_{P_2} Z_m \|_2^2 \xrightarrow{P_2} 0 \\
 &\quad \Leftrightarrow N(0, \sigma^2 I_p)
 \end{aligned}$$

Or si $E(|\varepsilon_i|^4) < +\infty$,

$$\sqrt{m} (\hat{\sigma}^2 - \sigma^2) \rightarrow N(0, \text{Var}(|\varepsilon_i|^2))$$

Preuve

$$\begin{aligned}
 \sqrt{m} (\hat{\sigma}^2 - \sigma^2) &= \sqrt{m} \left[\frac{\mathbf{E}^* \mathbf{E}}{m} - \sigma^2 \right] + \frac{P}{\sqrt{m}} \frac{\mathbf{E}^* \mathbf{E}}{m} + \frac{\mathbf{E}^* P_{Im \times E}}{m-p} \\
 &\quad \downarrow \text{N}(0, \text{Var}(|\varepsilon_i|^2)) \quad \frac{P}{\sqrt{m}} \xrightarrow{P} 0 \quad \frac{\mathbf{E}^* P_{Im \times E}}{m-p} \xrightarrow{P} 0
 \end{aligned}$$

IV Test dans le modèle linéaire gaussien

Test simples

Test de Student

$$(H_0) \quad \beta_q = a \quad \text{contre} \quad (H_1) \quad \beta_q \neq a$$

$$\text{Prat : } \text{Si } H_0 \quad \frac{\hat{\beta}_q - a}{\hat{\sigma} \sqrt{(X^* X)_{qq}}} \sim T(n-p) \quad \text{Student}$$

Preuve sous H_0 $\beta_q \sim N(a, \sigma^2 (X^* X)_{qq})$ et $\hat{\sigma}^2 \sim \sigma^2 X^* (n-p)$

$$\text{Or : Test } \left| \frac{\hat{\beta}_q - a}{\hat{\sigma} \sqrt{(X^* X)_{qq}}} \right| > Q_{T(n-p)}(1 - \frac{\alpha}{2}) = \tilde{Q}_{T(n-p)}\left(\frac{\alpha}{2}\right)$$

$\beta_q \Leftrightarrow$ avec J_C $F_T(Q_\alpha) = J_C^\alpha \quad \tilde{Q}(z) = Q(1-\alpha)$

Réponse

$$\begin{aligned}
 &2 Q_{T(n-p)} \left(\frac{|\hat{\beta}_q - a|}{\hat{\sigma} \sqrt{(X^* X)_{qq}}} \right) \\
 &= 2(1 - F_{T(n-p)}(\frac{|\hat{\beta}_q - a|}{\hat{\sigma} \sqrt{(X^* X)_{qq}}})) \\
 &\geq \alpha \text{ fonction de la zone de rejet}
 \end{aligned}$$

• Test sur une contrainte linéaire (Généralité)

$$H_0: b^* \beta = a \quad \text{contre} \quad H_1: b^* \beta \neq a$$

Pr sous $\frac{b^* \hat{\beta} - a}{\hat{\sigma} \sqrt{b^* (\hat{X}^T \hat{X})^{-1} b}} \sim T(n-p)$

\Rightarrow Test similaire au précédent

• Test sur σ^2

$$H_0: \sigma^2 = \sigma_0^2 \quad \text{contre} \quad H_1: \sigma^2 \neq \sigma_0^2$$

Pr sous $H_0 = (n-p) \frac{\hat{\sigma}^2}{\sigma_0^2} = \frac{1}{\sigma_0^2} \sum |y_i - x_i^T \hat{\beta}|^2 \sim \chi^2(n-p)$

$$= \frac{SSR(\beta)}{\sigma_0^2}$$

\Rightarrow

Test $\frac{SSR(\beta)}{\sigma_0^2} \geq Q_{\chi^2(n-p)}(1-\alpha)$

ou $\frac{SSR(\beta)}{\sigma^2} \leq Q_{\chi^2(n-p)}(\alpha)$

$$\begin{aligned} p \text{ values} &= \min \left(2 \left(1 - F_{\chi^2} \left(\frac{SSR(\beta)}{\sigma_0^2} \right) \right) \right. \\ &\quad \left. 2 \left(F_{\chi^2} \left(\frac{SSR(\beta)}{\sigma^2} \right) \right) \right) \end{aligned}$$

\Rightarrow Intervalle de conf.

o Test de Fisher

$$M_1: Y \sim N(X\beta, \sigma^2 I_m) \quad \beta \in \mathbb{R}^p$$

$$M_0: Y \sim N(Z\gamma, \sigma^2 I_m) \quad \gamma \in \mathbb{R}^q \quad q < p$$

on dit que M_0 est un sous modèle de M_1
(ou que les modèles sont imbriqués)

$$\text{si } \text{Im } Z \subset \text{Im } X$$

Ea : $\beta = H \gamma$ cf analyse de la norme

- importance de $n = p-q$ entre les deux sous espaces de dimension q à pourcentages.

Prédiction par ordre de complexité

$$\begin{array}{c} Y \quad \hat{Y}^{M_1} \quad \hat{Y}^{M_0} \\ || \quad \quad \quad || \\ P_{\text{Im } X} Y \quad P_{\text{Im } Z} Y \end{array}$$

comme $\text{Im } Z \subset \text{Im } X$ on en déduit

$$\text{que } I - P_{\text{Im } X} \perp P_{\text{Im } X} - P_{\text{Im } Z}$$

Condition gaussien $(I - P_{\text{Im } X})Y \perp \begin{cases} Y - \hat{Y}^{M_1} \\ Y - \hat{Y}^{M_0} \end{cases} \perp \begin{cases} (P_{\text{Im } X} - P_{\text{Im } Z})Y \\ \hat{Y}^{M_1} - \hat{Y}^{M_0} \end{cases}$

sous H_0 : M_0 est nulle

$$Y - \hat{Y}^{M_1} \sim \mathcal{N}(0, \sigma^2 P_{\text{Im } X}^{-1}) \quad (\text{c'est un équivalent au } M)$$

$$\hat{Y}^{M_1} - \hat{Y}^{M_0} \sim \mathcal{N}(0, \sigma^2 P_{\text{Im } Z \text{ dans } \text{Im } X})$$

$$\Rightarrow \|Y - \hat{Y}^{M_1}\|^2 \sim \sigma^2 X^2(m-p) + \perp \quad \| \hat{Y}^{M_1} - \hat{Y}^{M_0} \|^2 \sim \sigma^2 X^2(p-q) + \perp$$

$$\Rightarrow \frac{\|\hat{Y}^{M_1} - \hat{Y}^{M_0}\|^2 / (p-q)}{\|Y - \hat{Y}^{M_1}\|^2 / (m-p)} \sim T(p-q, m-p)$$

Pw1 $\frac{[\text{SSR}(M_0) - \text{SSR}(M_1)] / (p-q)}{\text{SSR}(M_1) / (m-q)} \sim T(p-q, m-q)$

l' erre $\|Y - \hat{Y}^{M_1}\|^2 = \text{SSR}(M_1)$

$$\|Y - \hat{Y}^{M_0}\|^2 = \text{SSR}(M_0) = \|\hat{Y} - \hat{Y}^{M_1}\|^2 + \|Y^{M_1} - \hat{Y}^{M_1}\|^2$$

$$\Rightarrow \|\hat{Y}^{M_1} - \hat{Y}^{M_0}\|^2 = \text{SSR}(M_0) - \text{SSR}(M_1)$$

\Rightarrow Test de Fieller.

O Test de Wald

$$N(\mathbf{x}\beta, \sigma^2 I_m) \text{ avec } \beta \in \mathbb{R}^p$$

Sont $C \in \mathbb{R}^{p \times p}$ de rang $(p-q)$ et $p-q$ colonnes avec $p \leq q \leq p$
et $c \in \mathbb{R}^{p-q}$

$$H_0: C\beta = c \text{ contre } H_1: C\beta \neq c$$

$$\text{Sous } H_0 \quad C\hat{\beta} - c \sim N(0, \sigma^2 C(X^*X)^{-1}C^*)$$

$$\begin{aligned} \underline{\text{Preuve}} \quad Z &\sim N(0, \Gamma) \text{ avec } \Gamma \text{ inversible} \\ &\Rightarrow Z^* \Gamma^{-1} Z \sim \chi_m^2 \end{aligned}$$

$$\begin{aligned} \text{Preuve} \quad \Gamma^{-1} &= P D P \quad \text{avec } \Gamma \text{ inversible et matrice de corrélation} \\ &= P D^{1/2} D^{1/2} P \quad D \geq 0 \end{aligned}$$

$$Z^* \Gamma^{-1} Z = \| D^{1/2} P Z \|^2$$

$$\text{Or } D^{1/2} P Z \sim N(0, D^{1/2} P \Gamma^{-1} P D^{1/2})$$

$$D^{1/2} P \Gamma^{-1} D^{1/2} P^* D^{1/2}$$

I_m

$$\text{Preuve} \quad \frac{[(C\hat{\beta} - c)^* [C(X^*X)^{-1}C^*]^{-1} (C\hat{\beta} - c)]/(p-q)}{\hat{\sigma}^2} \sim T(p-q, m-p)$$

Preuve Lemme + id de $\hat{\beta}$ et $\hat{\sigma}^2$

$$\text{"Corr"} \quad C = \text{Id} \quad \hat{\beta} = \beta$$

$$\frac{[(\hat{\beta} - \beta)^* (X^*X) (\hat{\beta} - \beta)]/p}{\hat{\sigma}^2} \sim T(p, m-p)$$

\Rightarrow Ellipsoïde de confiance

Ex Test du rapport à analyse de la corrélation avec variables

$$\text{Période 1: } \mathbb{R}^2 \quad Y_1 \sim \mathcal{N}(X_1 \beta_1, \sigma_1^2 I_n) \quad \beta_1 \in \mathbb{R}^p \quad n_1 \text{ observations}$$

$$\text{Période 2: } \quad Y_2 \sim \mathcal{N}(X_2 \beta_2, \sigma_2^2 I_n)$$

Test d'hypothèse de rupture

- Test $\sigma_1^2 = \sigma_2^2$

$$\frac{\text{SCR}(\hat{\beta}_1) / (n_1 - p)}{\text{SCR}(\hat{\beta}_2) / (n_2 - p)} \sim F(n_1 - p, n_2 - p)$$

- Modèle "complet"

Test $\beta_1 = \beta_2$ sur la méthode de Fisher

• Test de Fisher du modèle

$$H_0: E(Y) = \beta_0 \mathbf{1} \quad H_1: E(Y) = X\beta$$

sous H_0 : $\frac{[\text{SCR}(\hat{\beta}_0) - \text{SCR}(\hat{\beta}_1)] / (p-1)}{\text{SCR}(\hat{\beta}_1) / (n-p)} \sim F(p-1, n-p)$

$$= \frac{\|X\hat{\beta} - \bar{Y}\|^2 / (p-1)}{\|X\hat{\beta} - \bar{Y}\|^2 / (n-p)}$$

• Coefficient $R^2 = \frac{\|\hat{Y} - \bar{Y}\|^2}{\|Y - \bar{Y}\|^2} \leq 1$

mesure d'ajustement de \hat{Y}

$$R^2 = 1 - \frac{\|Y - \hat{Y}\|^2}{\|Y - \bar{Y}\|^2}$$

on peut utiliser $\langle \hat{Y} - \bar{Y}, Y - \hat{Y} \rangle = 0$, $\hat{Y} - \bar{Y} \perp Y - \hat{Y}$

on a également $\langle \hat{Y} - \bar{Y}, Y - \bar{Y} \rangle = \|\hat{Y} - \bar{Y}\|^2$

$$\Rightarrow R^2 = \frac{\langle \hat{Y} - \bar{Y}, Y - \bar{Y} \rangle^2}{\|\hat{Y} - \bar{Y}\|^2 \|Y - \bar{Y}\|^2} = \rho(Y, \hat{Y})$$

IV On vérifie que Fisher $R_0 \subset \mathcal{L}_1$

$$F = \frac{\frac{m - \dim(\mathcal{L}_1)}{\dim(\mathcal{L}_1) - \dim(\mathcal{L}_0)}}{\frac{R_{\mathcal{L}_1}^2 - R_{\mathcal{L}_0}^2}{1 - R_{\mathcal{L}_1}^2}}$$

VII Moindres carrés généralisés

$$\cdot H'_\varepsilon \quad \begin{cases} E(\varepsilon) = 0 \\ Cn(\varepsilon) = \Gamma^2 \Sigma \text{ avec } \Sigma \neq c \text{ Id} \end{cases}$$

2 cas distincts : $Cn(\varepsilon) = \text{Diag}(\Gamma^2)$

Héritage des méthodes

$\cdot \text{Corr}(\varepsilon_i, \varepsilon_j) \neq 0 \sim \varepsilon_i$ peuvent l'empêcher ?

VI Cas Σ connu

on note $\Sigma^{1/2}$ une racine de Σ qui existe

$$\Sigma = O D O \Rightarrow \Sigma^{1/2} = O D^{1/2} O \text{ où } D^{1/2} \text{ est diagonal}$$

$$\tilde{Y} = \Sigma^{-1/2} Y$$

$$\Rightarrow \tilde{Y} = \Sigma^{-1/2} X \beta + \underbrace{\Sigma^{-1/2} \varepsilon}_{\tilde{\varepsilon}}$$

$$\Rightarrow E(\tilde{\varepsilon}) = 0$$

$$Cn(\tilde{\varepsilon}) = O^2 \text{ Id}$$

\Rightarrow on s'est ramené au cas ordinaire

$$\Rightarrow \hat{\beta} = (\tilde{X}^* \tilde{X})^{-1} \tilde{X}^* \tilde{Y}$$

Interprétation géométrique: $\Sigma^{-1} \Rightarrow$ fonction quadratique définie positive
 \Rightarrow forme sur \mathbb{R}^m

$$\|Y\|_{\Sigma^{-1}}^2 = Y^* \Sigma^{-1} Y$$

$$\begin{aligned}\beta &= \text{origin } \|Y - X\beta\|^2 \\ &= \text{origin } \|Z^{1/2}(Y - X\beta)\|^2 \\ &= \text{origin } \|Y - X\beta\|_{\Sigma^{-1}}^2\end{aligned}$$

\rightsquigarrow Projection Σ^{-1} orthogonale!

Prop: • $E(\hat{\beta}) = \beta$ $\text{Var}(\hat{\beta}) = (X^* \Sigma^{-1} X)^{-1}$

• (Gauss-Markov) $\hat{\beta}$ est le meilleur estimateur linéaire sans biais de β

• $\hat{Y} = X\hat{\beta}$ est le projet Σ^{-1} orthogonal de Y sur $\mathbb{R}^m X$

• $\hat{\sigma}^2 = \frac{\|Y - X\hat{\beta}\|_{\Sigma^{-1}}^2}{n-p} = \frac{\|Y - X\beta\|_{\Sigma^{-1}}^2}{n-p}$ estime sans biais σ^2
 Preuve: immédiate. ($\neq \frac{\|Y - X\beta\|^2}{n-p}$)

Prop: Si Y est gaussien: $\hat{\beta}$ est gaussien $\| \hat{\sigma}^2 \sim \sigma^2 \frac{\chi^2_{(n-p)}}{n-p} \)$

• Pour β coincide avec le maximum de vraisemblance

• $\hat{\beta}$ est optimal parmi les estimateurs sans biais.

Preuve: $L_n(y; \beta, \sigma^2) = (2\pi\sigma^2)^{-n/2} \det(\Sigma)^{-1/2} \exp(Y - X\beta)^* \Sigma^{-1} (Y - X\beta)$

Pg: Pour calculer $\hat{\beta}$, on pose tout sur \hat{Y} et le calcul de $\Sigma^{-1/2}$
 • tout sur la forme et le calcul de Σ
 \rightarrow calcul des SCR par des formules \neq

■ Cas Σ inconnu

~ on utilise les méthodes courtes ordinaires:

$$\hat{\beta} = \arg \min \|y - X\beta\|_2^2$$

$$\Leftrightarrow \hat{\beta} = (X^* X)^{-1} X^* y$$

Pw₁ $E(\hat{\beta}) = \beta$

$$\text{Cov}(\hat{\beta}) = \sigma^2 [X^* X]^{-1} X^* \Sigma X (X^* X)^{-1}$$

⚠ on perd tout le reste ⚡

Pw₁ si $\lambda_M(\Sigma_m) = o(1)$

$$\lambda_m(X^* X) \xrightarrow[m \rightarrow \infty]{} +\infty \quad [\lambda_m((X^* X)^{-1}) = o(1)]$$

$$\Rightarrow \text{Tr}((X^* X)^{-1} X^* \Sigma X (X^* X)^{-1}) \rightarrow 0$$

$\Rightarrow \hat{\beta}$ est asymptotiquement asymptotique

Preuve:

$$\text{Tr}((X^* X)^{-1} X^* \Sigma X (X^* X)^{-1})$$

$$* = \text{Tr}(\Sigma X (X^* X)^{-2} X^*)$$

$$\leq \lambda_M(\Sigma) \text{Tr}(X (X^* X)^{-1} X^*) \leq \lambda_M(\Sigma) \text{tr}((X^* X)^{-1}) \\ \leq \lambda_M(\Sigma) \rho \lambda_M((X^* X)^{-1}) \rightarrow 0$$

: $\Sigma = O^ D O$ avec O orthogonale

$$\text{Tr}(\Sigma V) = \text{Tr}(O^* D O V)$$

$$= \text{Tr}(D O V O^*)$$

$$\leq (\text{Diag } D_{ii}) \text{Tr}(O V O^*)$$

$$\leq \lambda_M(\Sigma) \text{Tr}(V)$$

- Estimation de Σ impossible : n^2 variables
mais estimation "possible" si Σ a une forte proportion simple

VII Etude des résidus

1/ Toute

$$\text{résidu } \hat{\varepsilon}_i = Y_i - \hat{Y}_i \approx \varepsilon_i = Y_i - E(Y)$$

$$E(\hat{\varepsilon}) = 0$$

$$E(\varepsilon) = 0$$

$$\text{var}(\hat{\varepsilon}) = \sigma^2(1-H)$$

$$\text{var}(\varepsilon) = \sigma^2 I$$

$$\hookrightarrow \text{Var } \hat{\varepsilon}_i = \sigma^2 (1-h_{ii})$$

$$\hookrightarrow \text{Var } (\varepsilon_i) = \sigma^2$$

- Modification des résidus pour les rendre de même var.

$$\tilde{\varepsilon}_i = \frac{\hat{\varepsilon}_i}{\sqrt{1-h_{ii}}}$$

$$\tilde{\varepsilon}_i = \frac{\hat{\varepsilon}_i}{\sigma \sqrt{1-h_{ii}}}$$

\Rightarrow Standardisation en utilisant une estimation de la variance

$$\varepsilon_i = \frac{\hat{\varepsilon}_i}{\widehat{\sigma} \sqrt{1-h_{ii}}}$$

$$\text{ou } \varepsilon_i^* = \frac{\hat{\varepsilon}_i}{\widehat{\sigma}_{(i)} \sqrt{1-h_{ii}}}$$

où $\widehat{\sigma}_{(i)}$ estimation du module où l'on enlève la i^e observation

$$\hookrightarrow \varepsilon_i \sim \mathcal{N}(0,1)$$

$$\varepsilon_i^* \sim F(n-1-p)$$

Pb pour calculer $\widehat{\sigma}_{(i)}$ \Rightarrow Formule "magique"

$$\frac{\widehat{\sigma}^2}{\widehat{\sigma}_{(i)}^2} = \frac{n-p}{n-1-p-t_i^2}$$

$$\Rightarrow \varepsilon_i^* = \varepsilon_i \sqrt{\frac{n-p}{n-1-p-\hat{\varepsilon}_i^2}}$$

△ Pas de décorrélator entre les ϵ_i , $t_{c,i}$, t_i . △

En pratique, on fait les tests sur ϵ_i ou t_i :

⇒ Test de signe

- Wilcoxon

basé sur $\sum_i \text{sign}(t_{c,i})$

⇒ Tests

⇒ Test de symétrie / indépendance

- Wolf-Wolfowitz Test des séries

$$\text{sign}(t_{c,1}) = + + - - + + + - + - -$$

⇒ on compte le nombre du "nombre de signes"

⇒ Test d'indépendance

- Durbin-Watson

→ Estimateur de p dans un modèle AR(1)

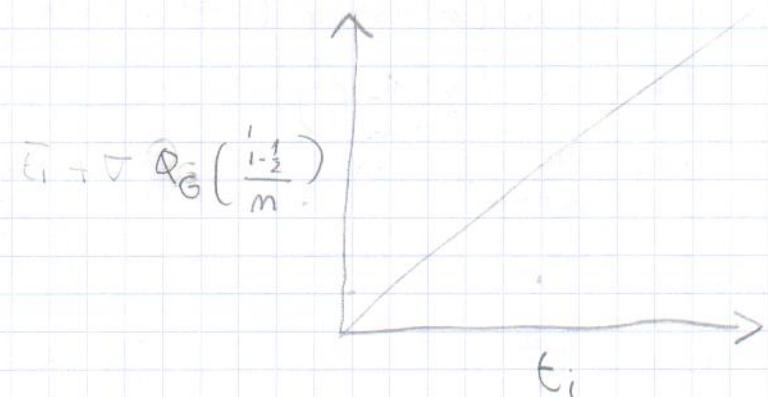
$$\epsilon_i = p \epsilon_{i-1} + \sqrt{1-p^2} \omega_i$$

+ Test de $p=0$

⇒ Test de normalité

- Pearson : Test du χ^2 sur des boîtes équiprobables

- SQ-Q plot de Shapiro-Francis



Shapiro-Francis

: test sur la corrélation

entre $(\tilde{\epsilon}_i)$ et $(Q_0(\frac{i-1}{m}))$

Lilliefors (Kolmogorov-Smirnov)

$$D_+ = \max \left(\frac{1}{m} - F_G \left(\frac{x_{(i)} - \bar{x}}{\hat{\sigma}} \right) \right)$$

$$D_- = \max \left(F_G \left(\frac{x_{(i)} - \bar{x}}{\hat{\sigma}} \right) - \frac{i-1}{m} \right)$$

$$D = \max(D_+, D_-)$$

$$\Rightarrow \|F_m - F_G\|_\infty : \text{Kolmogorov-Smirnov}$$

Groner-von Mises

$$\sum_{i=1}^m \left(F_G \left(\frac{x_{(i)} - \bar{x}}{\hat{\sigma}} \right) - \frac{i-\frac{1}{2}}{m} \right)$$

Anskum-Darling

$$\sum_{i=1}^{m/2} (2i-1) \left[\ln \left(F_G \left(\frac{x_{(i)} - \bar{x}}{\hat{\sigma}} \right) \right) + \ln \left[1 - F_G \left(\frac{x_{(m-i+1)} - \bar{x}}{\hat{\sigma}} \right) \right] \right]$$

plus de poids que le test de Gruber-Mises

Sloping Wilk

$$\frac{1}{\hat{\sigma}^2} \sum_{i=1}^{m/2} a_i \left(t_{(m-i+1)} - t_{(i)} \right)^2 \quad (\text{by rule of Slope-Wilk})$$

Test d'Anderson / Test de Jarque-Bera

Tous basés sur des relations entre les moments des lois gaussiennes.

\Rightarrow Test de séparabilité

2/ Diagnostics

Tente de comprendre "bizarres" (différent) sur les données

■ Leverage = Lévrier

$$\cdot h_{ii} = x_i^* (X^* X)^{-1} x_i = X^* (X X^*)^{-1} X = H$$

≈ distance entre le centre du nuage des (x_i)

muni de la distance de Mahalanobis associée au nuage de points.

$$\cdot \hat{y}_i = h_{ii} y_i + \sum_{i \neq j} h_{ij} y_j$$

Propriété : $\sum h_{ii} = p$

Propriété $\sum h_{ii} = \text{Tr}(H) = p$

Heuristique : si tous les points ont aussi importants $h_{ii} \approx \frac{p}{n}$

Points à étudier si $h_{ii} > C \frac{p}{n}$ $C=2$ ou 3

■ Residu standardisé $t_i = \frac{\hat{\epsilon}_i}{\sqrt{1-h_{ii}}}$

Heuristique : si t_i est grand, il faut comprendre si c'est un point mal placé ou un point atypique

$$|t_i| > t_\alpha \quad t_\alpha \approx 2-3$$

■ Residu standardisé $t_i^* = \frac{\hat{\epsilon}_i}{\sqrt{\hat{\sigma}_m^2 (1-h_{ii})}} = \hat{\epsilon}_i \sqrt{\frac{n-1-p}{n-p-t_i^2}}$

Heuristique similaire

$$|t_i^*| > t_\alpha \quad t_\alpha \approx 2-3 \quad \sim F(m-1-p)$$

Rq si le point est mal placé ou atypique il a une influence fort sur $\hat{\sigma}$ $\Rightarrow |t_i^*| \gg |t_i| \rightarrow$ plus sensible

Rq: on a toujours $|e_i^*| \geq |e_i|$

■ Variabilité de la prédictibilité quand on retire i :

Difference of fits

$$DFITS_i = \frac{\hat{y}_i - \hat{y}_{(i),i}}{\hat{\sigma}_{(i)} \sqrt{h_{ii}}} = e_i^* \sqrt{\frac{h_{ii}}{1-h_{ii}}} \quad (\sim \mathcal{N}(0,1))$$

Relation entre e_i^* et h_{ii}

Tst d'origine $|DFITS_i| > c \sqrt{\frac{p}{m}}$ $c \approx 2-3$

■ Variabilité de la prédictibilité quand on retire i

Distance de Cook

$$D_i = \frac{\|\hat{y} - \hat{y}_{(i),i}\|^2}{\hat{\sigma}^2 p} = \frac{e_i^2}{p} \frac{h_{ii}}{1-h_{ii}} \quad \left(\approx \frac{\chi^2(p)}{p} \right)$$

Tst $D_i > 1$

$$\text{ou } D_i > \frac{4}{n-p}$$

■ Variabilité de l'estimation des coefficients quand on retire i

$$DFBETAS_{j,i} = \frac{\hat{\beta}_j - \hat{\beta}_{(i),j}}{\hat{\sigma}_{(i)} \sqrt{(X^* X)^{-1}_{j,j}}} \quad (\sim \mathcal{N}(0,1))$$

$$\text{DFBETAS}_{j,i} > 1$$

$$DFBETAS_{j,i} > \frac{c}{\sqrt{m}}$$

$$= e_i^* \left[\frac{[(X^* X)^{-1} X^*]_{j,i}}{\sqrt{(X^* X)^{-1}_{j,j} (1-h_{ii})}} \right]$$

\Rightarrow Mesure l'influence conjointe sur les paramètres

■ Variabilité de la matrice

$$V(\hat{\beta}) = \hat{\sigma}^2 \det(X^* X)^{-1} \Rightarrow \text{risque des biais}$$

$$\text{COVRATIO}_i = \frac{V(\hat{\beta}_{(i)})}{V(\hat{\beta})}$$

$$= \left[\frac{m-1-p}{m-p} + \frac{e_i^2}{m-p} \right] \frac{p}{(1-h_{ii})}$$

comparer au rapport à 1

< 1 désirable

> 1 anormal

$$\text{Test } |\text{CorRatio}_{ij} - 1| > \frac{\epsilon^P P}{m} \quad C \sim z \rightarrow$$

3 / Résidus partiels,

outil de diagnostic

$$y_0 = x_0^* \beta + \epsilon_0$$

$$\Rightarrow y_0 = x_{(i)}^* \beta_{(i)} + (x_i \beta_i + \epsilon_i)$$

$$\hat{\epsilon}_{(i)} = y_0 - x_{(i)}^* \beta_{(i)} = x_{(i)} \beta_i + \epsilon_i$$

\Rightarrow on se ramène à 1 pb à une variable

En pratique

$$\hat{\epsilon}_{(i)} = y - x_{(i)} \hat{\beta}_{(i)}$$

avec $\hat{\beta}_{(i)}$ étant soit dans le modèle complet
soit dans le modèle sans X_i
en arrière

$$\text{L'observation de } \hat{\epsilon}_{(i)} = y_0 - x_{(i)}^* \hat{\beta}_{(i)}$$

en fonction de $x_{(i)}$ permet de visualiser
le reste de dépendance en X_i

\Rightarrow Estimation "visuelle" de l'information restante

VIII Sélection de Modèles / Sélection de Variables

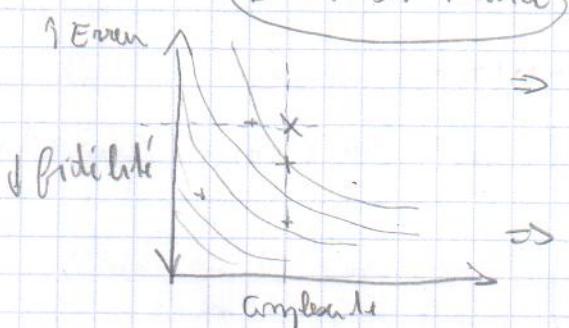
- Même question : comment choisir les modèles à utiliser ?

Objectif ≠ : prediction → sélection de modèles
interprétation → sélection de variables

- Principe de sélection = principe de parcimonie

Principe du raisin d'Occam

Rédezise d'un compromis entre fidélité aux observations et complexité.



⇒ A même complexité, on choisit le modèle le plus fidèle

⇒ A même fidélité, on choisit le modèle le moins complexe

Pb Comment donner un "ordre total" sur R^2 ?

- Ordre partiel : Test de Fisher

$M_0 \subset M_1$

$$\Rightarrow \frac{[RSS(0) - RSS(1)] / (p_1 - p_0)}{RSS(1) / (n - p_1)} > t_\alpha$$

$$RSS(0) + t_\alpha \frac{RSS(1)}{n - p_1} p_0 > RSS(1) + t_\alpha \frac{RSS(1)}{n - p_1} p_1$$

⇒ choix du modèle M_1

$$\text{critérium } C(R) = RSS(M) + t_\alpha \hat{\sigma}^2 p$$

- $R^2 = 1 - \frac{\|y - \hat{y}\|^2}{\|y - \bar{y}\|^2}$ augmente avec la fidélité mais ne diminue pas avec la complexité
⇒ pas de compromis

$$R^2_q = 1 - \frac{\|Y - \hat{Y}\|^2 / (n-p)}{\|Y - \bar{Y}\|^2 / (n-1)}$$

perte en ampleur de la complémenté avec le terme on $(n-p)$

$$R^2_q(\mathcal{M}_0) \leq R^2_q(\mathcal{M}_1)$$

$$\frac{\|Y - \hat{Y}\|^2}{n-p_0} \geq \frac{\|Y - \hat{Y}'\|^2}{n-p_1}$$

$$\widehat{U}_0^2 \geq \widehat{U}_1^2 \rightsquigarrow \text{Varina dans le modèle précédent}$$

\Rightarrow Difficile à utiliser.

■ Principe d'optimisation on moyenne

qualité d'un estimateur mesuré par son erreur quadratique moyenne:

$$\begin{aligned} & E(\|E(Y) - \hat{Y}\|^2) \\ &= \underbrace{\|E(Y) - E(\hat{Y})\|^2}_{\text{biais}^2} + \underbrace{\text{Var}(\|\hat{Y}\|^2)}_{\text{Varina}} \\ &= \|E(Y) - P_{\text{Im}X} E(Y)\|^2 + \text{Var}(\|P_{\text{Im}X} Y\|^2) \\ &= \|E(Y) - P_{\text{Im}X} E(Y)\|^2 + \sigma^2 p \end{aligned}$$

Pb: Terme de biais impossible à mesurer.

\Rightarrow Estimation à l'aide de $\|Y - \hat{Y}\|^2$

$$\begin{aligned} E(\|Y - \hat{Y}\|^2) &= E(\|Y - P_{\text{Im}X} Y\|^2) \\ &= \|E(Y) - P_{\text{Im}X} E(Y)\|^2 + \text{Var}(\|Y - P_{\text{Im}X} Y\|^2) \\ &= \|E(Y) - P_{\text{Im}X} E(Y)\|^2 + \sigma^2 (n-p) \end{aligned}$$

$$\Rightarrow E(\|E(Y) - \hat{Y}\|^2) = E(\|Y - \hat{Y}\|^2) + (2p-n) \sigma^2$$

$$\Rightarrow \text{Utilisation du critère: } \|Y - \hat{Y}\|^2 + (2p-n) \sigma^2$$

\Rightarrow pb estimation de σ^2 .

$$\Rightarrow \text{donc } \hat{\sigma}^2 = \frac{\|Y - \hat{Y}\|^2}{n-p} \Rightarrow \|Y - \hat{Y}\|^2 + (2p-n) \hat{\sigma}^2 = \frac{p}{n-p} \|Y - \hat{Y}\|^2$$

pas négatif si nos de biais

$$= p R^2_q$$

\Rightarrow plus stricte que R^2_{adj}

\Rightarrow Forme classique : C_p de Mallows

$$C_p = \frac{\|Y - \hat{Y}\|^2}{\sigma^2} + 2p - m$$

Si l'estimateur est sans biais $C_p \geq p$

Règle du critère : minimiser C_p

Règle du biais : dans le modèle le plus simple tel que $C_p \leq p$.

■ Vraisemblance et généralisation

$$\ln(L(Y; (\beta, \sigma^2); M))$$

$$= -\frac{m}{2} \log \sigma^2 - \frac{m}{2} \log 2\pi - \frac{1}{2\sigma^2} \|Y - X\beta\|^2$$

Valeur en le maximum de vraisemblance

$$\hat{\beta} = (X^T X)^{-1} X^T Y \quad \hat{\sigma}^2 = \frac{\|Y - X\hat{\beta}\|^2}{m} + \text{MCO}$$

$\ln(\text{ML})$:

$$\ln(L(Y; (\hat{\beta}, \hat{\sigma}^2); M))$$

$$= -\frac{m}{2} \log \hat{\sigma}^2 - \frac{m}{2} [\log 2\pi + 1]$$

$$= -\frac{m}{2} \log \left(\frac{\|Y - \hat{Y}\|^2}{m} \right) - \frac{m}{2} [\log 2\pi + 1]$$

ne dépend que de l'ultime avc données (comme le R^2)

\Rightarrow nécessité d'une généralisation

choix critique: $-2\ln(\text{ML}) + 2\lambda_m |\beta|_1$

Variante sur le choix de λ_m

AIC : $\lambda_m = 2$

BIC : $\lambda_m = \log m$

■ AIC : Am Information Criterion
Akaike Information Criterion

Idée : Bon modèle est un modèle tel que la probabilité $P(X; \hat{\theta}(Y); M)$ soit proche de $P(X)$
en moyenne
 \Rightarrow Mesure pour la divergences de Kullback-Leibler
 $KL(P_0, P(\cdot; \hat{\theta}(x); M))$

$$= E_{X \sim P_0} \left[\ln \frac{P_0(x)}{P(\cdot; \hat{\theta}, M)} \right]$$

$$= - E_{X \sim P_0} \left[-\ln P(X; \hat{\theta}(x); M) \right] + E_{X \sim P_0} \left[\ln P_0 \right]$$

$$E_{X \sim P_0} [KL(P_0, P(\cdot; \hat{\theta}(x); M))]$$

$$= - E_{X \sim P_0} E_{Y \sim P_0} \left[\ln P(X; \hat{\theta}(y); M) \right]$$

$$+ \underbrace{E_{X \sim P_0} \left[\ln P_0(x) \right]}_{\text{cste}}$$

Pb. Comment estimer $E_{Y \sim P_0} E_{X \sim P_0} [-\ln P(X; \hat{\theta}(y); M)]$
à partir de $-\ln P(Y; \hat{\theta}(Y); M)$

Idée de la preuve

$$\hat{\theta}_0 = \underset{\theta}{\operatorname{argmin}} E_{X \sim P_0} [-\ln P(X; \theta; M)]$$

$$E_{Y \sim P_0} E_{X \sim P_0} [-\ln P(X; \hat{\theta}(y); M)] - E_X [-\ln P(X; \theta_0; M)] \\ \simeq d/2 \quad (\chi^2)$$

$$E_Y [-\ln P(Y; \theta_0; M)] - [-\ln P(Y; \theta_0(Y); M)] \\ \simeq d/2$$

$\Rightarrow \chi^2$ qui tend vers ...

$$E_{\theta \sim \Omega} E_{x \sim P_0} [-\ln P(x; \theta(y); M)]$$

$$\simeq -\ln P(y; \theta(y); M) + d = AIC(M)/2$$

* Validation croisée.

$$AIC(M) = -2 \ln P(y; \theta(y); M) + 2d$$

• BIC : Bayesian Information Criterion

Idée: Modélisation bayésienne et choix du modèle le plus probable étant donnée y .

$$\begin{aligned} P(Y; \theta; M) &= P(Y|\theta; M) P(\theta|M) P(M) \\ &= P(M, \theta|Y) P(Y) \end{aligned}$$

$$\begin{aligned} P(M|Y) &= \int P(M, \theta|Y) dP_\theta \\ &= \int \frac{P(Y|\theta; M) P(\theta|M) P(M)}{P(Y)} P(\theta|M) d\theta \\ &= \frac{P(M)}{P(Y)} \int P(Y|\theta; M) P^2(\theta|M) d\theta \end{aligned}$$

Si $P(M) = \text{stc}$ il suffit de déterminer l'unique

$$\int P(Y|\theta; M) P^2(\theta|M) d\theta = \int e^{-m \left[\frac{1}{n} \sum \ln P(Y|\theta; M) - 2 \ln P(\theta|M) \right]} d\theta$$

$$\lim_{m \rightarrow \infty} \int e^{-m L(u)} du \approx e^{-m L(u^*)} \left(\frac{2\pi}{m} \right)^{\frac{d}{2}} | -L''(u^*) |^{-1/2}$$

$$L = \frac{1}{m} \left[-\ln P(Y|\theta; M) - 2 \ln P(\theta|M) \right]$$

$$-L''(\theta) = \underbrace{+ \frac{1}{m} \frac{\partial^2}{\partial \theta^2} \ln P(Y|\theta; M)}_{\left[\frac{\partial^2}{\partial \theta^2} \ln P(\theta|M) \right]} + \frac{2}{n} \frac{\partial^2}{\partial \theta^2} \ln P(\theta|M)$$

$$\hookrightarrow A(\theta_M) = O\left(\frac{1}{m}\right)$$

$$- \ln P(\mathbf{u} | \gamma)$$

$$\approx -\ln P(\gamma | \mathbf{M}, \theta^*) - 2 \ln P(\theta^* | \mathbf{u})$$

$$+ \frac{d}{2} \log n - \frac{d}{2} \log 2\pi - \ln |\mathcal{A}_n(\theta^*)|$$

$$\theta^* = \text{argmin } L(\theta) = \text{argmin} \left[-\frac{1}{m} \ln P(\gamma | \theta, \mathbf{u}) - \frac{2}{m} \ln P(\theta, \mathbf{u}) \right]$$

$$\rightarrow \hat{\theta} = \text{argmin} \left[\frac{1}{m} \ln P(\gamma | \theta, \mathbf{u}) \right]$$

$$\Rightarrow -\ln P(\mathbf{u} | \gamma) \approx -\ln P(\gamma | \mathbf{M}, \hat{\theta}) + \frac{d}{2} \log n$$

$-2 \ln P(\theta^* | \mathbf{u}) \frac{d}{2} \log 2\pi - \frac{1}{2} \ln |\mathcal{A}_n(\theta^*)|$

$O(1)$

$$\Rightarrow \text{BIC}(\mathbf{u}) = -2 \ln P(\gamma | \mathbf{u}, \hat{\theta}) + \text{Argin d}$$

$$\approx -\ln P(\mathbf{u} | \gamma)$$

Rq : ne dépend pas de $P(\theta | \mathbf{u})$

Prise en compte l'a priori sur les modèles

en ajoutant un terme $-\ln P(\mathbf{u})$

■ Compréhension

$$\begin{array}{ccc} R_q & (P) \cdot \text{AIC} & \text{BIC} \\ \xleftarrow{\text{comparaison des modèles dans}} & & \end{array}$$

Rq : AIC $\rightarrow C_p$ si σ^2 est connu.

lin avec l'ordre sur les modèles emportés.

■ Explorations exhaustives

→ Essayer tous les modèles

et choisir celui qui minimise la valeur croissante.

1b : Algorithmique couteuse → avec exploration avec le nb de modèles

■ Exploration intelligente.

→ Explorer l'ensemble des modèles en ayant / tenant
des variables explicatives en fonction du précédent.

→ Stratégie "gloutonne" : moins coûteuse mais difficile
de montrer l'optimalité.

■ Partial Least Squares

on ajoute les régressions au fur et à mesure en fonction
de leur corrélation avec le rendu initial

on choisit ensuite des variables explicatives le meilleur

■ Méthode arête généralisée : si $X^T X$ n'est pas inversible

$$\|Y - X\beta\|^2 + \lambda \|\beta\|_p$$

Si $p=0$: \approx critère de sélection de modèles

Si $p=2$: critère de régularisation de type Tikhonov

$p=1$: sélection et algorithme efficace

⇒ LASSO : sélection + résultat théorique

■ Validation croisée : Test des modèles avec une approche
similaire à AIC en découplant les données en apprendre
et vérifier.

Modèle linéaire généralisé

\Rightarrow Généralisation du modèle linéaire

Famille de loi P_0 : gaussien, binomiel nbin

$$\eta_i = x_i^* \beta$$

$$\Rightarrow y_i \sim P_0 \Rightarrow g(E(y_i)) = x_i^* \beta$$

\uparrow fonction de link

ex: P_0 = gaussien, $g(x) = x \Rightarrow$ Modèle linéaire gaussien

P_0 Bernoulli, $g(x) = \ln \frac{x}{1-x} \Rightarrow$ logit

P_0 Poisson $\rightarrow g(x) = \ln x$

\Rightarrow Famille "exponentielle", Maximum de vraisemblance, Test

I Famille exponentielle

y_i variables indépendantes admettant des distributions issues d'une même structure exponentielle.

\Rightarrow Densité ayant une forme particulière

$$f(y_i; \theta_i, \phi) = e^{\underbrace{\frac{y_i \theta_i - r(\theta_i)}{\phi}}_{\text{fonction de masse naturelle de la variable}} + \underbrace{w(y_i, \phi)}_{\text{ne dépend pas de } \theta_i}}$$

Rq $E(y_i; \theta, \phi) = r'(\theta)$

$$\text{Var}(y_i; \theta, \phi) = \phi r''(\theta)$$

$$\begin{aligned} \frac{\partial \ln f}{\partial \theta} &= \frac{y_i - r'(\theta)}{\phi} E\left(\frac{\partial \ln f}{\partial \theta}\right) = 0 \\ \frac{\partial^2 \ln f}{\partial \theta^2} &= -\frac{r''(\theta)}{\phi} - E\left(\frac{\partial^2 \ln f}{\partial \theta^2}\right) \\ &= E\left(\left(\frac{\partial \ln f}{\partial \theta}\right)^2\right) \end{aligned}$$

Rq L'estimation du maximum de vraisemblance

de $y_1, \dots, y_n \sim f(\cdot; \theta; \phi)$ est $\hat{\theta}$ défini

$$\text{par } r'(\hat{\theta}) = \frac{1}{n} \sum y_i$$

\Rightarrow Structure simple

⚠ estimation de f plus complexe ...

Ex • Gaussian

$$f(y; m, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y-m)^2}$$

$$= e^{\left(\frac{ym - \frac{m^2}{2}}{\sigma^2} + \left[-\frac{y^2}{2\sigma^2} - \frac{1}{2}\ln(2\pi) - \frac{1}{2}\ln\sigma^2\right]\right)}$$

$\theta = m$

$$n(m) = \frac{m^2}{2} \quad \phi = \sigma^2$$

$$n'(m) = m \Rightarrow E(y) = m$$

$$n''(m) = 1 \Rightarrow V(y) = \sigma^2$$

• Bernoulli

$$f(z, \pi) = \pi^z (1-\pi)^{1-z} = 1-\pi \left(\frac{\pi}{1-\pi}\right)^z$$

$$= e^{z \ln \frac{\pi}{1-\pi} + \ln(1-\pi)}$$

$$\theta = \ln \frac{\pi}{1-\pi} \Rightarrow e^\theta = \frac{\pi}{1-\pi} \Rightarrow (1-\pi)e^\theta = \pi$$

$$\rho = 1 \qquad \qquad \qquad \Rightarrow \pi = \frac{e^\theta}{1+e^\theta}$$

$$f(z, \theta) = e^{z\theta - \ln(1+e^\theta)}$$

$$n(\theta) = -\ln(1+e^\theta) \quad (= \ln(1-\pi))$$

$$n'(\theta) = \frac{e^\theta}{1+e^\theta} \quad (= \pi) \quad (= E(z))$$

$$n''(\theta) = \frac{e^\theta(1+e^\theta) - e^\theta e^\theta}{(1+e^\theta)^2} = \frac{e^\theta}{(1+e^\theta)^2} - \frac{e^\theta}{1+e^\theta} \cdot \frac{1}{1+e^\theta} \quad (= \pi \times (1-\pi))$$

$$= V(y)$$

• Poisson

$$f(z, \lambda) = \frac{\lambda^z e^{-\lambda}}{z!} = e^{z \ln \lambda - \lambda + \ln z!}$$

$$\theta = \ln \lambda \quad \phi = 1$$

$$n(\theta) = e^\theta \quad (= \lambda)$$

$$n'(\theta) = e^\theta \quad (= \lambda = E(z))$$

$$n''(\theta) = e^\theta \quad (= \lambda = E(z))$$

Parmi, Gamma ...

- By $\eta''(\theta) = \text{Var}_\theta(y) \Rightarrow \eta$ strictement concave
 $\Rightarrow E(y) = b'(\theta)$ est en liaison avec θ
- Gaussian $E(y) = m = \theta$
- Bernoulli $E(y) = \pi = \frac{e^\theta}{1+e^\theta} \Leftrightarrow \theta = \ln \frac{\pi}{1-\pi} = \ln \frac{E(y)}{1-E(y)}$ liaison canonique
- Poisson $E(y) = \lambda = e^\theta \Leftrightarrow \theta = \ln \lambda = \ln E(y)$

II Modèles linéaires généralisés

- 1. Famille exponentielle pour y_i indépendants
- 2. Existence de régressions x_i et d'un coefficient β tel que
 $g(E(y_i|x_i)) = x_i^* \beta$ où g est une fonction croissante

- Ex • Gaussian : $g = I$ (= fonction unie) \Rightarrow modèle linéaire gaussien
- Bernoulli : $\bullet g = \ln \frac{x}{1-x}$ (= fonction canonique) \Rightarrow modèle logit
- $\bullet g^{-1}(x) = \frac{e^x}{1+e^x}$
- $\bullet g = \Phi^{-1}(x)$ (= quantile gaussien) \Rightarrow modèle probit
- $g^{-1}(x) = \Phi(x)$
- $\bullet g(x) = \ln(-\ln(1-x)))$ \Rightarrow modèle log log
- $g^{-1}(x) = 1 + e^{-\frac{1}{x}}$
- Poisson $\bullet g(\theta) = \ln \theta$ \Rightarrow modèle régression avec lien logarithmique
- $g^{-1}(x) = e^x$
- $\Rightarrow g(\eta'(\theta_i)) = x_i^* \beta \Leftrightarrow \theta_i = (\eta')^{-1}(g^{-1}(x_i^* \beta))$
- Ré si g fonction croissante $g(\eta'(\theta)) = \theta = x_i^* \beta$
 $\Rightarrow \theta_i = x_i^* \beta$

D) Prédit $\hat{\beta} \Rightarrow E(y) = g^{-1}(x^* \beta)$

D) Vraisemblance

$$l_m(y_i; \theta, \beta, \phi)$$

$$= \sum_{i=1}^m l_m(y_i; \theta(x_i, \beta), \phi)$$

$$= \sum_{i=1}^m \left[\frac{y_i \theta_i(\beta, x_i) - \eta'(\theta_i(x_i, \beta))}{\phi} + w(y_i, \phi) \right]$$

$$\frac{\partial l_m}{\partial \beta_j} = \sum \frac{\partial l_m}{\partial \theta_i} \frac{\partial \theta_i}{\partial \beta}$$

$$= \sum \frac{y_i - \eta'(\theta_i(x_i, \beta))}{\phi} \frac{\partial \theta_i(x_i, \beta)}{\partial \beta}$$

$$\frac{\partial \theta_i}{\partial \beta_j} = \frac{1}{N''(\theta_i)} x_{ij} (g')'(x_i^* \beta)$$

$$\frac{\partial l_m}{\partial \beta_j} = \sum_{i=1}^n \frac{(y_i - g'(x_i^* \beta)) x_{ij}}{\phi N''(\theta_i) g'(x_i^* \beta)} (g')'(x_i^* \beta)$$

\Rightarrow Equations du maximum de vraisemblance

$$\frac{\partial l_m}{\partial \beta_j} = 0 \quad \forall j$$

\Rightarrow Résolution unique possible

Rq : si $g = (N)^{-1}$ $(g')' = N''$

$$\Rightarrow \frac{\partial l_m}{\partial \beta_j} = \sum_{i=1}^n \frac{(y_i - N'(x_i^* \beta)) x_{ij}}{\phi}$$

$$\Rightarrow X^*(Y - N) = 0 \quad \text{où} \quad N = \begin{pmatrix} N'(x_1^* \beta) \\ N'(x_2^* \beta) \\ \vdots \\ N'(x_n^* \beta) \end{pmatrix} = \begin{pmatrix} E(y_1) \\ E(y_2) \\ \vdots \\ E(y_n) \end{pmatrix}$$

III Qualités d'agnostique et tests

~~1) Qualités d'agnostique~~

• Démonstration : $D = -2(L - L_0)$ ($\sim \chi^2$)

↑ ↗

vraisemblance dans le modèle
au maximum de solué à n paramètres,
vraisemblance dans le modèle真是的
à p paramètres

Asymptotiquement sous H_0 le modèle est vrai : $D \sim \chi^2(n-p)$

• Pearson

$$\text{Sous } H_0 \quad \chi^2 = \sum_{i=1}^m \frac{(y_i - g^{-1}(x_i^* \hat{\beta}))^2}{\phi v''(v)^{-1} g^{-1}(x_i^* \hat{\beta})} \sim \chi^2(n-p)$$

~~2) Tests~~

• Rapport de vraisemblance

Ratio dans modèles entraînés M_0 contre

$$\text{Sous } H_0, \quad D_1 - D_0 = -2(L_1 - L_0) \sim \chi^2(p_1 - p_0)$$

si ϕ est connu

$\sim F(p_1 - p_0, n-p)$
Gammique

• Test de Wald

$H_0: C\beta = d$ avec $r = \text{nb de contraintes}$.

$$\Rightarrow \text{Sous } H_0: (C\beta - d)^* \left(C(C^* W C)^{-1} C^* \right)^{-1} (C\beta - d) \sim \chi^2(r)$$

où W est diagonal et négatif. $W_{ii} = \frac{1}{\phi v''(v)^{-1} g^{-1}(x_i^* \beta)} \left[(g^{-1})'(x_i^* \beta) \right]^2$

Pb Estimer f à partir de l'observation

$$\text{ok } y_i = f(x_i) + \varepsilon_i$$

on a un comment faire si $f(x_i) = \beta_0 + \beta_1 x_i$

on plus généralement si $f(x_i) = \beta_0 \sum_{k=1}^p \beta_k x_k$

Exemp $\beta_k(x_i) = x_i^k$

⇒ Approche paramétrique

Hyp $f \in \mathcal{F}_0$ avec $\Theta \subset \mathbb{R}^p$

$$\Rightarrow \hat{f}_{\theta} \text{ avec } \hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum (y_i - f(\theta))^2$$

on a également étudié le cas de la puissance formelle

Approche non paramétrique

Hyp $f \in \mathcal{C}^0 \Rightarrow$ Top de solution unique !

on va chercher \hat{f}_m dans un modèle \mathcal{M}_m

dans la complexité augmente avec $m \rightarrow \mathcal{C}^0$

Ex. $\mathcal{M}_m = \text{polynôme de degré } k(m)$ avec $k(m) \rightarrow +\infty$

Sélection du degré par régression

$L^2(t_0, t_1)$ pour validation croisée

• $\mathcal{M}_m = \text{polynôme trigonométrique}$

↓ Même principe

$$L^2(t_0, t_1)$$

Biais = régulier
relat grande

Variance diminue

- Birmin = décaillage en bouteille
 - ⇒ utiliser un modèle linéaire simple dans chaque bouteille
 - ⇒ constat / hypothèse / ...
 - ⇒ difficulté taille des bouteilles ; biais vs variance.

- Pb : taille des bouteilles ⇒ discertinité
 - ⇒ solution de la régénér. bouteille
 - ⇒ régénér. avec les plus petits
 - ⇒ Pb discertinité entre avec l'entraînement et la taille des pots
 - ⇒ Noyau à support coquard artificiel / sceller. tout / bouteilles
 - ⇒ Régénér. bouteille préalable à la régénér. robuste
 - ⇒ Pb chose de la bouteille bouteille / biais du potager / bouteille / Young

- Bouteille qui se "porte" : condition de racine
 - ⇒ sphère
 - ⇒ chose des mouds ⇒ sélection de modèles.
 - ⇒ Problématique bouteille varia (formulation variationnelle?)

- Ondulations : exc de la bouteille de Ptolemy.
 - ⇒ "Pb de la taille des bouteilles"
 - ⇒ modèles embûchés + Ptolemy erreure du $L^2(\Omega)$