

M1 Statistique

2^e partie

I Vecteurs gaussiens

II Modèle linéaire gaussien

III Estimation non-paramétrique

I Vecteurs gaussiens

1- Propriétés des vecteurs gaussiens

Définition

Une variable aléatoire Z est dite de loi normale centrée réduite, ce qu'on note $Z \sim \mathcal{N}(0,1)$ si sa loi a pour densité $x \mapsto \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

Une v.a. Y est dite de loi normale de paramètres (μ, σ^2) si il existe une v.a. $Z \sim \mathcal{N}(0,1)$ telle que $Y = \mu + \sigma Z$.

Lorsque $\sigma = 0$, on dit que Y est une v.a. gaussienne (normale) dégénérée

Remarque

Si $\sigma > 0$, la densité de la v.a. $Y \sim \mathcal{N}(\mu, \sigma^2)$ est

$$x \mapsto \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Théorème Soit $Y \sim \mathcal{N}(\mu, \sigma^2)$

Pour tout $\xi \in \mathbb{C}$, on a $\varphi_Y(\xi) = e^{i\xi\mu - \frac{\sigma^2 \xi^2}{2}} = \mathbb{E}(e^{i\xi Y})$

Définition

Un vecteur aléatoire X à valeurs dans \mathbb{R}^d est un vecteur gaussien si toute combinaison linéaire de ses composantes est une v.a. gaussienne.

Le vecteur des moyennes de X est le vecteur $(\mathbb{E}(X_1), \dots, \mathbb{E}(X_d))$ et sa matrice de covariance est

$$\begin{aligned}\Sigma_{i,j} &= \text{Cov}(X_i, X_j) = \mathbb{E} \left((X_i - \mathbb{E}(X_i)) (X_j - \mathbb{E}(X_j)) \right) \\ &= \mathbb{E}(X_i X_j) - \mathbb{E}(X_i) \mathbb{E}(X_j)\end{aligned}$$

Théorème

La loi d'un vecteur gaussien est caractérisée par son vecteur des moyennes m et sa matrice de covariance Σ . On note cette loi $\mathcal{N}(m, \Sigma)$, et si $X \sim \mathcal{N}(m, \Sigma)$, alors

$$\xi \in \mathbb{C}^d \quad \varphi_X(\xi) = \mathbb{E}(e^{i\xi \cdot X}) = \exp\left(i m \cdot \xi - \frac{1}{2} \xi^t \Sigma \xi\right)$$

preuve

La loi d'un vecteur aléatoire est caractérisée par sa fonction caractéristique. Soit X un vecteur gaussien, on pose m son vecteur des moyennes et Σ sa matrice de covariance.

Pour $\lambda \in \mathbb{R}^d$, on calcule $\mathbb{E}(e^{i\lambda \cdot X})$

Puisque X est un vecteur gaussien, $\lambda_1 X_1 + \dots + \lambda_d X_d = \lambda \cdot X$ est une v.a. gaussienne. De plus, par linéarité

$$\begin{aligned}\mathbb{E}(\lambda \cdot X) &= \lambda_1 \mathbb{E}(X_1) + \lambda_2 \mathbb{E}(X_2) + \dots + \lambda_d \mathbb{E}(X_d) \\ &= \lambda \cdot m\end{aligned}$$

$$\begin{aligned}\text{Var}(\lambda \cdot X) &= \text{Var}(\lambda_1 X_1 + \dots + \lambda_d X_d) \\ &= \sum_{i=1}^d \text{Var}(\lambda_i X_i) + 2 \sum_{i=1}^d \sum_{j=i+1}^d \text{Cov}(\lambda_i X_i, \lambda_j X_j)\end{aligned}$$

$$\begin{aligned} \text{Var}(\lambda \cdot X) &= \sum_{i=1}^d \lambda_i^2 \text{Var}(X_i) + 2 \sum_{i=1}^d \sum_{j=i+1}^d \lambda_i \lambda_j \text{Cov}(X_i, X_j) \\ &= \sum_{i=1}^d \sum_{j=1}^d \lambda_i \lambda_j \Sigma_{i,j} \\ &= {}^t \lambda \Sigma \lambda \end{aligned}$$

Autrement dit, $\lambda \cdot X \sim \mathcal{N}(\lambda \cdot m, {}^t \lambda \Sigma \lambda)$

$$\text{donc } \mathbb{E}(e^{i \lambda \cdot X}) = e^{i \lambda \cdot m - \frac{1}{2} {}^t \lambda \Sigma \lambda}$$

La fonction caractéristique de X ne dépend que de m et de Σ , donc la loi de X est caractérisée par m et par Σ .

La fonction caractéristique s'étend de façon unique en une fonction holomorphe, on a également

$$\mathbb{E}(e^{i \lambda \cdot X}) = e^{i \lambda \cdot m - \frac{1}{2} {}^t \lambda \Sigma \lambda} \quad \text{pour } \lambda \in \mathbb{C}^d$$

Propriétés

* Σ est une matrice symétrique positive

$$\Sigma_{i,j} = \text{Cov}(X_i, X_j) = \text{Cov}(X_j, X_i) = \Sigma_{j,i}$$

$$\text{Pour tout } \lambda \in \mathbb{R}^d, {}^t \lambda \Sigma \lambda = \text{Var}(\lambda \cdot X) \geq 0$$

* Σ est diagonalisable en base orthonormée

* Si (X_1, \dots, X_d) sont des v.a. gaussiennes **indépendantes** alors $X = (X_1, \dots, X_d)$ est un vecteur gaussien

* X gaussien \Rightarrow Toutes les coordonnées de X sont gaussiennes

Exemple

$$Y \sim \mathcal{N}(0, 1), \quad Z \sim \mathcal{B}\left(\frac{1}{2}\right) \quad Y \perp\!\!\!\perp Z$$

On pose $X = (Y, (2Z-1)Y)$
 $X_1 \sim \mathcal{N}(0, 1)$

$$\begin{aligned} \mathbb{E}\left(e^{i\lambda(2Z-1)Y}\right) &= \mathbb{E}\left(\mathbb{E}\left(e^{i\lambda(2Z-1)Y} \mid Y\right)\right) \\ &= \mathbb{E}\left(\frac{1}{2} e^{i\lambda Y} + \frac{1}{2} e^{-i\lambda Y}\right) \\ &= \frac{1}{2} e^{-\frac{\lambda^2}{2}} + \frac{1}{2} e^{-\frac{(-\lambda)^2}{2}} = e^{-\frac{\lambda^2}{2}} \end{aligned}$$

donc $X_2 \sim \mathcal{N}(0, 1)$

mais X n'est pas gaussien

$$X_1 + X_2 = 2Z Y \quad \text{n'est pas gaussienne}$$

$$\mathbb{P}(X_1 + X_2 = 0) = \frac{1}{2}$$

Proposition

Soit X un vecteur gaussien

Si $\text{Cov}(X_i, X_j) = 0$, alors $X_i \perp\!\!\!\perp X_j$

preuve

Le vecteur (X_i, X_j) est un vecteur gaussien

de matrice de covariance $\Sigma = \begin{pmatrix} \text{Var}(X_i) & 0 \\ 0 & \text{Var}(X_j) \end{pmatrix}$

$$\text{donc } \mathbb{E} \left(e^{i(\lambda X_i + \mu X_j)} \right) = e^{i\lambda \mathbb{E}(X_i) + i\mu \mathbb{E}(X_j) - \frac{1}{2} \lambda^2 \text{Var}(X_i) - \frac{1}{2} \mu^2 \text{Var}(X_j)}$$

$$= \mathbb{E} \left(e^{i\lambda X_i} \right) \mathbb{E} \left(e^{i\mu X_j} \right)$$

donc $X_i \perp\!\!\!\perp X_j$

Proposition

Si $X \sim \mathcal{N}(m, \Sigma)$ alors pour tout $A \in \mathbb{R}^{d', d}$
 et pour tout $b \in \mathbb{R}^{d'}$, $AX + b \sim \mathcal{N}(\underline{Am + b}, \underline{A\Sigma^t A})$

preuve

Toute combinaison linéaire des coordonnées de $AX + b$
 est une combinaison linéaire des coordonnées de X
 donc est gaussienne

On calcule son vecteur des moyennes. Soit $i \leq d'$

$$\begin{aligned} \mathbb{E} \left((AX + b)_i \right) &= \mathbb{E} \left((AX)_i + b_i \right) \\ &= \mathbb{E} \left((AX)_i \right) + b_i \\ &= \mathbb{E} \left(\sum_{k=1}^d A_{ik} X_k \right) + b_i \\ &= \sum_{k=1}^d A_{i,k} \mathbb{E} (X_k) + b_i \\ &= \sum_{k=1}^d A_{i,k} m_k + b_i = (Am + b)_i \end{aligned}$$

On calcule la matrice de covariances.
 Soit $i, j \leq d$ et on calcule

$$\begin{aligned}
 & \text{Cor}((AX+b)_i, (AX+b)_j) \\
 &= \text{Cor}((AX)_i + b_i, (AX)_j + b_j) \\
 &= \text{Cor}((AX)_i, (AX)_j) \\
 &= \text{Cor}\left(\sum_{k=1}^d A_{i,k} X_k, \sum_{l=1}^d A_{j,l} X_l\right) \\
 &= \sum_{k=1}^d \sum_{l=1}^d A_{i,k} A_{j,l} \text{Cor}(X_k, X_l) \\
 &= \sum_{k=1}^d \sum_{l=1}^d A_{i,k} A_{j,l} \Sigma_{k,l} \\
 &= \sum_{l=1}^d (A\Sigma)_{i,l} A_{l,j} = (A\Sigma^t A)_{i,j}
 \end{aligned}$$

Exemple

$$X \sim \mathcal{N}(0, I_2)$$

$$Y = \underbrace{AX+b}_{\text{avec } A = \begin{pmatrix} 1 & 2 \\ 0 & 3 \\ 2 & 1 \end{pmatrix} \quad b = \begin{pmatrix} 0 \\ -1 \\ 0 \end{pmatrix}}$$

$$Y \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 5 & 6 & 4 \\ 6 & 9 & 3 \\ 4 & 3 & 5 \end{pmatrix}\right)$$

$$\begin{pmatrix} 1 & 0 & 2 \\ 2 & 3 & 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 2 \\ 0 & 3 \\ 2 & 1 \end{pmatrix}$$

Algorithme de génération d'un vecteur gaussien

Soit $m \in \mathbb{R}^d$ et Σ une matrice symétrique positive

$$Y \sim \mathcal{N}(m, \Sigma)$$

Σ est diagonalisable en base orthonormée, donc

$$\Sigma = P D {}^t P$$

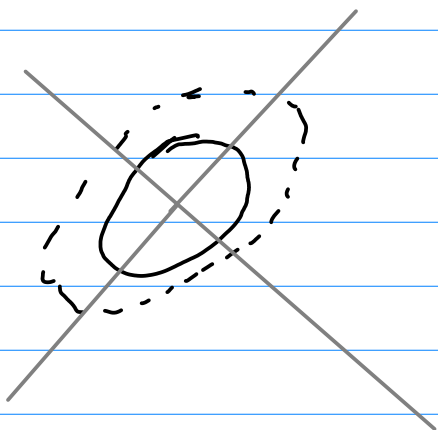
$$\text{Donc } X = {}^t P (Y - m) \sim \mathcal{N}(0, {}^t P \Sigma P) \sim \mathcal{N}(0, D)$$

La matrice de covariance de X est diagonale,
donc les coordonnées X_1, \dots, X_d sont indépendantes

Par conséquent, si on pose Z_1, \dots, Z_d
des v.a. iid de $\mathcal{P}_i \mathcal{N}(0, 1)$, on définit ensuite

$$X_i = \sqrt{D_i} Z_i$$

puis on pose $Y = P X + m$



Théorème

Soit $X \sim \mathcal{N}(m, \Sigma)$.

* Si Σ est inversible, alors

X a pour densité $\frac{1}{(2\pi)^{d/2} \sqrt{\det \Sigma}} \exp\left(-\frac{1}{2} (x-m)^T \Sigma^{-1} (x-m)\right)$

* Si Σ n'est pas inversible, alors

la loi de $(X-m)$ est supportée par $\text{Ker}(\Sigma)^\perp = \text{Vec}(e_i)$
où (e_i) les vep de Σ
de $\text{rang} > 0$

preuve

o Si Σ est inversible, alors

$$X = PZ + m, \text{ avec } \Sigma = P D^t P$$

et $Z \sim \mathcal{N}(0, D)$

Z a pour densité $\prod_{i=1}^d \frac{1}{\sqrt{2\pi D_i}} e^{-\frac{z_i^2}{2D_i}} \quad D_i > 0$

d'où, par changement de variables, on obtient

la densité de X , où on utilise

$$\prod_{i=1}^d D_i = \det \Sigma$$

o Si Σ n'est pas inversible

alors $X = PZ + m$ p.s. avec $Z \sim \mathcal{N}(0, D)$

et le support de Z est donné par

$$\{(x_1, \dots, x_d) \in \mathbb{R}^d : x_i = 0 \text{ si } D_i = 0\}$$

Proposition

Si (Y, X_1, \dots, X_d) est un vecteur gaussien
alors $E(Y | X_1, \dots, X_d)$ est une fonction affine de X_1, \dots, X_d

$$E(Y | X_1, \dots, X_d) = a + b_1 X_1 + \dots + b_d X_d$$
$$Y - E(Y | X_1, \dots, X_d) \perp (X_1, \dots, X_d)$$

Exemple

calculer $E(Y_1 | Y_2, Y_3)$

$$E(Y_1 | Y_2, Y_3) = a + bY_2 + cY_3$$

avec $E(E(Y_1 | Y_2, Y_3)) = a + bE(Y_2) + cE(Y_3)$

$$Y \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 5 & 6 & 4 \\ 6 & 9 & 3 \\ 4 & 3 & 5 \end{pmatrix}\right)$$

$$\Rightarrow 0 = a - b$$

$$E(E(Y_1 | Y_2, Y_3) Y_2) = aE(Y_2) + bE(Y_2^2) + cE(Y_2 Y_3)$$

$$\Rightarrow 6 = -a + 6b + 3c$$

$$E(E(Y_1 | Y_2, Y_3) Y_3) = aE(Y_3) + bE(Y_2 Y_3) + cE(Y_3^2)$$

$$\Rightarrow 4 = 3b + 5c$$

$$\begin{cases} a - b = 0 \\ -a + 6b + 3c = 6 \\ 3b + 5c = 4 \end{cases} \Rightarrow \begin{cases} a = \dots \\ b = \dots \\ c = \dots \end{cases}$$

Théorème Central Limite

Soit (V_1, \dots, V_n) un n-échantillon de vecteurs aléatoires
de vecteur de moyenne m et de matrice de covariance Σ

$$\text{alors } \sqrt{n}(\bar{V}_n - m) \Rightarrow \mathcal{N}(0, \Sigma)$$

2- Le modèle gaussien

Soit (X_1, \dots, X_n) un n -échantillon de loi $\mathcal{N}(\mu, \sigma^2)$
sous la loi $\mathbb{P}_{\mu, \sigma^2}$

Proposition

L'estimateur du maximum de vraisemblance de (μ, σ^2)
est $(\bar{X}_n, \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2)$

$$\overset{\text{preuve}}{L(X_1, \dots, X_n; \mu, \sigma^2)} = \sum_{i=1}^n \log \left(\frac{e^{-\frac{(X_i - \mu)^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \right)$$

$$= \left(\sum_{i=1}^n -\frac{(X_i - \mu)^2}{2\sigma^2} \right) - \frac{n}{2} \log(2\pi\sigma^2)$$

On calcule le gradient de L

$$\left\{ \begin{array}{l} \frac{dL}{d\mu} = \sum_{i=1}^n -\frac{(\mu - X_i)}{\sigma^2} = -n \frac{\mu}{\sigma^2} + \frac{1}{\sigma^2} \sum_{i=1}^n X_i \\ \frac{dL}{d\sigma^2} = \sum_{i=1}^n \frac{(X_i - \mu)^2}{2(\sigma^2)^2} - \frac{n}{2\sigma^2} \end{array} \right.$$

\Rightarrow L'estimateur du maximum de vraisemblance

$$\text{est } (\hat{\mu}_n, \hat{\sigma}_n^2) \text{ où } \hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_n)^2$$

Remarque

$$\mathbb{E}(\hat{\mu}_n) = \mu$$

$$\begin{aligned} \mathbb{E}(\hat{\sigma}_n^2) &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}((X_i - \bar{X}_n)^2) \\ &= \frac{n-1}{n} \sigma^2 \end{aligned}$$

donc on préfère souvent $S_n = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$

car c'est un estimateur sans biais

Théorème de Cochran

Soit X un vecteur gaussien de loi $\mathcal{N}(m, \sigma^2 I)$
On note $\mathbb{R}^d = E_1 \oplus E_2 \oplus \dots \oplus E_r$ en sous-espaces orthogonaux

Les projections orthogonales $\pi_{E_1} X, \dots, \pi_{E_r} X$
sont des vecteurs gaussiens indépendants.

preuve

On construit une base de \mathbb{R}^d e_1, \dots, e_d
telle que e_1, \dots, e_{d_1} est une base de E_1 ,
...

$e_{d_1+1}, \dots, e_{d_1+d_2}$ est une base de E_2

On note $U = (e_1 | \dots | e_d)$ ($U^t U = I$)

On observe que ${}^t U X \sim \mathcal{N}({}^t U m, \underbrace{{}^t U \sigma^2 I U}_{\sigma^2 I})$

La matrice de cov de ${}^t U X$ est diagonale, donc
ses coordonnées sont indep

donc $(e_1, X), \dots, (e_{d_1}, X)$ est indépendant de

$(e_{d_1+1}, X), \dots, (e_{d_1+d_2}, X)$ - -

$(e_{d_1+d_2+1}, X), \dots, (e_{d_1+d_2+d_3}, X)$

d'où $\pi_{E_1} X, \pi_{E_2} X, \dots, \pi_{E_r} X$ sont indépendants

Corollaire

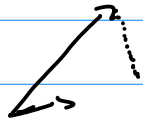
$$\hat{\mu}_n \perp \hat{\sigma}_n^2$$

preuve

$$\text{On décompose } \mathbb{R}^n = \underbrace{\text{Vect}(1, \dots, 1)}_{E_1} \oplus \underbrace{\text{Vect}(1, \dots, 1)^\perp}_{E_2}$$

Posons $u = (1, \dots, 1)$

$$\text{Soit } x \in \mathbb{R}^n, \quad \pi_{E_1} x = \frac{(u, x)}{(u, u)} u = \frac{x_1 + x_2 + \dots + x_n}{n} u$$



$$\pi_{E_2} x = \bar{x}_n u$$

$$\pi_{E_2} x = x - \pi_{E_1} x = x - \frac{x_1 + \dots + x_n}{n} u$$

$$\pi_{E_2} x = \begin{pmatrix} x_1 - \bar{x}_n \\ x_2 - \bar{x}_n \\ \vdots \\ x_n - \bar{x}_n \end{pmatrix}$$

Par théorème de Cochran, $\pi_{E_1} x$ et $\pi_{E_2} x$

sont indépendants.

$$\text{De plus } \hat{\sigma}_n^2 = \frac{1}{n} \|\pi_{E_2} x\|^2$$

$$\hat{\mu}_n = \frac{u \cdot \pi_{E_1}(x)}{n}$$

On en déduit que $\hat{\mu}_n$ et $\hat{\sigma}_n^2$ sont indépendants.

Définition

Loi du χ^2 à d degrés de liberté

X suit une loi du χ^2 à d degrés de liberté

($X \sim \chi^2(d)$) si il existe un vecteur gaussien
 $Y \sim \mathcal{N}(0, I_d)$ tel que $X = \|Y\|_2^2$

$$= Y_1^2 + Y_2^2 + \dots + Y_d^2$$

Propriété Si $X_1 \sim \chi^2(d)$ et $X_2 \sim \chi^2(d')$
et $X_1 \perp X_2$ alors $X_1 + X_2 \sim \chi^2(d+d')$

Proposition
Si $X \sim \chi^2(d)$ alors X admet pour densité

$$z \mapsto \mathbb{1}_{z > 0} \frac{1}{2^{\frac{d}{2}} \Gamma(\frac{d}{2})} z^{\frac{d}{2}-1} e^{-\frac{z}{2}}$$

preuve

$$\mathbb{E}(f(X)) = \mathbb{E}(f(Y_1^2 + \dots + Y_d^2))$$

$$= \int_{\mathbb{R}^d} f(y_1^2 + \dots + y_d^2) \frac{e^{-\frac{y_1^2 + \dots + y_d^2}{2}}}{(2\pi)^{\frac{d}{2}}} dy_1 \dots dy_d$$

$$= c \int_0^{+\infty} r^{d-1} f(r^2) e^{-\frac{r^2}{2}} dr$$

$$z = r^2 \quad dz = 2r dr$$

$$= \frac{c}{2} \int_0^{+\infty} z^{\frac{d}{2}-1} f(z) e^{-\frac{z}{2}} dz$$

Définition

Loi de Student à d d.d.P

La v.a. Z est de loi de Student à d d.d.P
(ce qu'on note $Z \sim \mathcal{C}(d)$) si on peut écrire

$$Z \sim \frac{X}{\sqrt{Y/d}} \quad \text{avec} \quad \begin{cases} X \sim \mathcal{N}(0,1) \\ Y \sim \chi^2(d) \\ X \perp Y \end{cases}$$

Proposition

La densité de la loi $\mathcal{C}(d)$
est donnée par $x \mapsto \frac{1}{\sqrt{\pi d}} \frac{\Gamma(\frac{d+1}{2})}{\Gamma(\frac{d}{2})} \left(1 + \frac{x^2}{d}\right)^{-\frac{d+1}{2}}$

En particulier, la loi $\mathcal{C}(d)$
n'admet pas de moment d'ordre d

Proposition

Soit X un n -échantillon de loi $\mathcal{N}(\mu, \sigma^2)$

$$* \bar{X}_n \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

$$* \frac{n \hat{\sigma}_n^2}{\sigma^2} \sim \chi^2(n-1)$$

$$* \frac{\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma}}{\sqrt{\frac{n \hat{\sigma}_n^2}{\sigma^2(n-1)}}} \sim \mathcal{C}(n-1)$$

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\hat{\sigma}_n} \sim \mathcal{C}(n-1) \sim \sqrt{n} \frac{\bar{X}_n - \mu}{\sqrt{S_n^2}}$$

Remarque

Quand $n \rightarrow \infty$, $\frac{Y_1^2 + \dots + Y_n^2}{n} \rightarrow 1$ p.s.

$$\frac{\chi^2(n)}{n} \rightarrow 1$$

$$\mathcal{E}(n) \Rightarrow \mathcal{N}(0, 1)$$

pour de grandes valeurs de n , on écrit

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sqrt{S_n^2}} \approx \mathcal{N}(0, 1)$$

Intervalles de confiance

Intervalle de confiance de niveau $1 - \alpha$ sur μ

$$\text{Sous } \mathbb{P}_{\mu, \sigma^2} \left(\sqrt{n-1} \left| \frac{\bar{X}_n - \mu}{\sqrt{S_n^2}} \right| > z \right) = \mathbb{P} \left(|\mathcal{E}_{n-1}| > z \right)$$

donc en posant $t_{1-\alpha/2}^{(n-1)}$ le quantile d'ordre $1 - \alpha$ de la loi $\mathcal{E}(n-1)$

$$\text{on construit } \text{IC}_{1-\alpha}(\mu) = \left[\bar{X}_n \pm t_{1-\alpha/2}^{(n-1)} \sqrt{\frac{S_n^2}{n}} \right]$$

De même, pour construire un IC sur la variance

$$\mathbb{P}_{\mu, \sigma^2} \left(\frac{(n-1)S_n^2}{\sigma^2} \in [a, b] \right) = \mathbb{P} \left(\chi^2(n-1) \in [a, b] \right)$$

$$\text{d'où } \text{IC}_{1-\alpha}(\sigma^2) = \left[\frac{(n-1)S_n^2}{\chi_{1-\alpha/2}^{(n-1)}}, \frac{(n-1)S_n^2}{\chi_{\alpha/2}^{(n-1)}} \right]$$

II Le modèle linéaire gaussien

On suppose qu'on dispose d'un vecteur d'observations $Y = (Y_1, \dots, Y_n)$ associé à une matrice de données $X = \begin{pmatrix} X_{1,1} & \dots & X_{1,d} \\ \vdots & & \vdots \\ X_{n,1} & \dots & X_{n,d} \end{pmatrix}$

On suppose qu'il existe une relation linéaire entre X et Y , qu'on essaie de déterminer

1- Définition du modèle linéaire gaussien

Définition vectorielle

Soit Y un vecteur d'observations réelles, on dit que Y suit un modèle linéaire gaussien si Y s'écrit $Y = m + \varepsilon$, où $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$

On suppose que $m \in V$ un sous-espace vectoriel fixé et connu de \mathbb{R}^n

$$(\mathbb{P}_{m, \sigma^2} \sim \mathcal{N}(m, \sigma^2 I), m \in V, \sigma^2 \in \mathbb{R}_+)$$

Définition matricielle

Soit Y un vecteur d'observations réelles, on dit que Y suit un modèle linéaire gaussien si Y s'écrit

$$Y = X\beta + \varepsilon \quad \text{où } \varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$$

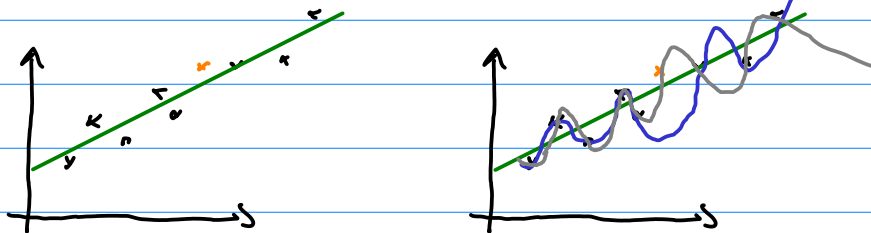
$$\forall j \leq n, Y_j = X_{j,1} \beta_1 + \dots + X_{j,d} \beta_d + \varepsilon_j$$

$$(\mathbb{P}_{\beta, \sigma^2} \sim \mathcal{N}(X\beta, \sigma^2 I), \beta \in \mathbb{R}^d, \sigma^2 \in \mathbb{R}_+)$$

Remarque

* L'espace vectoriel V dans la définition vectorielle correspond à $\text{Im}(X) = \{X\beta, \beta \in \mathbb{R}^d\}$ dans la définition matricielle

* On suppose toujours que $\dim(V) = d \ll n$
Si d est trop grand, on risque le surapprentissage



La matrice X des données est souvent appelée matrice des variables explicatives et Y est le vecteur des variables à expliquer

Exemple

Y_t Le prix d'achat d'une voiture

$X_{t,1}$ Le revenu de l'acheteur

$X_{t,2}$ Le nombre de km parcouru par l'acheteur

Un TLG pour Y serait

$$Y_t = \mu + \beta_1 X_{t,1} + \beta_2 X_{t,2} + \varepsilon_t$$

On pourrait aussi imaginer

$$Y_t = \mu' + \beta_1' X_{t,1}^2 + \beta_2' \sqrt{X_{t,2}} + \varepsilon_t$$

Objectif du modèle Linéaire gaussien

→ Expliquer les données,
expliquer l'existence d'une relation entre Y et X
corrélation n'est pas causalité

→ Prédire les mesures des futures données

→ Objectif statistique
Identifier β et σ^2

Tester l'hypothèse $\{\beta_1 = 0\}$

→ Deux types de données explicatives

Variables quantitatives $X \in \mathbb{R}$

Variables qualitatives $X \in E$ ensemble fini

IP est souvent intéressant

d'identifier E avec $\{(1,000), (0,1,000), (0,0,1,000)\}$

ANOVA

Remarque

Identifiabilité du modèle

$\theta \mapsto P_\theta$ est injective

Le modèle Linéaire gaussien est identifiable
ssi la matrice X est injective (de rang d)

(Pour tout $m \in \text{Im}(X)$, il existe un unique
 $\beta \in \mathbb{R}^d$ tel que $m = X\beta$)

Exemple

On teste une nouvelle combinaison engrais - variété

On construit 4 types de terrains

ancien engrais - ancienne variété

ancien engrais - nouvelle variété

nouvel engrais - ancienne variété

nouvel engrais - nouvelle variété

Y_t Le rendement de la parcelle,

Le rendement moyen de chaque parcelle est donné par un rendement de référence μ et d'un écart de rendement c_1, \dots, c_4 dépendant de la méthode

$$X = \begin{pmatrix} 1 & 1 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 1 & 1 \end{pmatrix} \quad Y = X(\mu, c_1, \dots, c_4) + N(0, I_n)$$

ce modèle est non-identifiable

Pour le rendre identifiable, on peut ajouter une condition, par exemple

- $\mu = 0$, c_i le rendement moyen d'une parcelle de type i
- $c_1 + c_2 + c_3 + c_4 = 0$

2. Estimation de paramètre

On considère " \cdot " le produit scalaire Euclidien sur \mathbb{R}^n , $\|\cdot\|_2$ la norme associée à ce produit scalaire et Π_V est la projection orthogonale sur V

On note aussi Π_V la matrice de projection orthogonale sur V

Théorème

On considère le TLLG $Y = m + \varepsilon$
où $m \in V$ et $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$

L'estimateur du maximum de vraisemblance de (m, σ^2) est $(\hat{m}, \hat{\sigma}_n^2)$, où

$$\rightarrow \hat{m} = \Pi_V Y$$

$$\rightarrow \hat{\sigma}_n^2 = \frac{1}{n} \|Y - \Pi_V Y\|_2^2 = \frac{1}{n} \|\Pi_{V^\perp} Y\|_2^2$$

Les variables \hat{m} et $\hat{\sigma}_n^2$ sont indépendantes et de loi $\hat{m} \sim \mathcal{N}(m, \sigma^2 \Pi_V)$

$$\text{et } \frac{n \hat{\sigma}_n^2}{\sigma^2} \sim \chi^2(n - \dim V)$$

Remarque

On a en particulier $E\left(\frac{n \hat{\sigma}_n^2}{\sigma^2}\right) = n - \dim(V)$

Si cet estimateur était non-biaisé, on aurait obtenu n , cet estimateur est donc biaisé

$$\text{On peut utiliser } \hat{\sigma}_n^2 = \frac{1}{n - \dim V} \|Y - \Pi_V Y\|_2^2$$

preuve

Sous la loi \mathbb{P}_{m, σ^2} Y est de loi $\mathcal{N}(m, \sigma^2 I_d)$

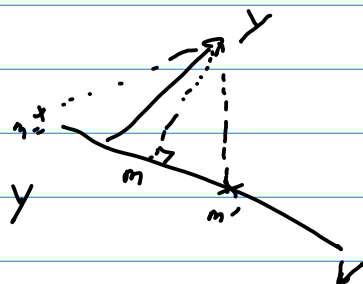
Sa densité est donc

$$f(y; m, \sigma^2) = \frac{1}{(2\pi)^{n/2}} \frac{1}{\sqrt{(\sigma^2)^n}} \exp\left(-\frac{1}{2\sigma^2}(y-m)(y-m)\right)$$
$$= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{\|y-m\|^2}{2\sigma^2}\right)$$

La log-vraisemblance de Y est donc

$$L(y; m, \sigma^2) = -\frac{\|Y-m\|^2}{2\sigma^2} - \frac{n}{2} \log(2\pi\sigma^2)$$

→ Quelle que soit la valeur de σ^2 , cette quantité est maximale par m minimisant $\|Y-m\|$, c'est-à-dire $m = \Pi_V Y$



On fixe $m = \Pi_V(Y)$ et on étudie

$$\sigma^2 \mapsto L(Y; m, \sigma^2) = -\frac{\|Y-m\|^2}{2\sigma^2} - \frac{n}{2} \log(2\pi\sigma^2)$$

En dérivant, cette quantité est maximale par σ^2 tel que

$$\frac{\|Y-m\|^2}{2(\sigma^2)^2} - \frac{n}{2\sigma^2} = 0$$

$$\Rightarrow \sigma^2 = \frac{1}{n} \|Y-m\|^2$$

Par théorème de Cochran, on a
 $\Pi_V Y \perp \Pi_{V^\perp} Y$

$$\text{donc } \hat{m} = \Pi_V Y \perp \hat{S}_n^2 = \frac{1}{n} \|\Pi_{V^\perp} Y\|_2^2$$

$$\begin{aligned} \hat{m} &\sim \mathcal{N}(\Pi_V m, \Pi_V (\sigma^2 I)^k \Pi_V) \\ &\sim \mathcal{N}(m, \sigma^2 \Pi_V) \end{aligned} \quad \begin{aligned} &{}^t \Pi_V = \Pi_V \\ &\text{et } \Pi_V \circ \Pi_V = \Pi_V \end{aligned}$$

$$\text{et } \Pi_{V^\perp} Y \sim \mathcal{N}(0, \sigma^2 \Pi_{V^\perp})$$

$$\dim(V^\perp) = n - \dim(V)$$

$$\text{donc } \frac{\|\Pi_{V^\perp} Y\|_2^2}{\sigma^2} \sim \chi^2(n - \dim V) \quad \square$$

On souhaite reformuler ce théorème sous forme matricielle pour une implémentation plus facile.

Une première étape est de déterminer la projection orthogonale sur $\text{Im}(X)$

Lemme

Soit X une matrice $n \times d$ avec $n \geq d$

$$X \text{ injective} \iff {}^t X X \text{ inversible}$$

preuve

Si ${}^t X X$ est inversible, alors elle est bijective

donc X est injective (et ${}^t X$ est surjective)

i.e. est de rang maximum d

On suppose X injective, et on suppose par l'absurde que ${}^t X X$ n'est pas inversible

Il existe $u \in \mathbb{R}^d$ non nul tel que ${}^t X X u = 0$

$$\text{Donc } {}^t u \cdot ({}^t X X u) = 0$$

$$\text{i.e. } {}^t (X u) (X u) = 0$$

$$\text{d'où } \|X u\|^2 = 0 \text{ ce qui implique } X u = 0$$

Mais $X u = 0$ montre que X n'est pas injective
on a une contradiction. \square

Proposition

Soit X une matrice $n \times d$ avec $n \geq d$ injective

La projection orthogonale de \mathbb{R}^n sur $\text{Im}(X)$ est donnée par $X ({}^t X X)^{-1} {}^t X$

Remarque

Posons $P = X ({}^t X X)^{-1} {}^t X$, alors

$$- P^2 = P$$

$$X ({}^t X X)^{-1} {}^t X X ({}^t X X)^{-1} {}^t X = X ({}^t X X)^{-1} I {}^t X = P$$

$\Rightarrow P$ est une matrice de projection

$$- {}^t P = P$$

$${}^t (X ({}^t X X)^{-1} {}^t X) = ({}^t X) {}^t (({}^t X X)^{-1}) {}^t X$$

$${}^t ({}^t X X) = {}^t X X$$

$$= X ({}^t ({}^t X X))^{-1} {}^t X$$

$$= X ({}^t X X)^{-1} {}^t X$$

$\Rightarrow P$ est une matrice de projection orthogonale

preuve

La projection orthogonale sur $\text{Im}(X)$ est la projection orthogonale sur $W = \text{Vect}(e_1, \dots, e_d)$

$$\text{où } X = (e_1 | e_2 | \dots | e_d)$$

X est injective donc e_1, \dots, e_d sont libres

Soit $u \in \mathbb{R}^n$, on cherche à déterminer $\hat{u} \in W$ tel que $u - \hat{u} \perp W$

Cela revient à imposer $\langle u - \hat{u}, e_j \rangle = 0 \quad \forall j \leq d$

On peut réécrire cette condition

$${}^t e_1 \cdot (u - \hat{u}) = 0$$

$${}^t e_2 \cdot (u - \hat{u}) = 0$$

⋮

$${}^t e_d \cdot (u - \hat{u}) = 0$$

i.e ${}^t X (u - \hat{u}) = 0$

Donc ${}^t X u = {}^t X \hat{u}$

De plus, comme $\hat{u} \in \text{Im}(X)$, on peut écrire $\hat{u} = X \hat{v}$

On a donc ${}^t X u = {}^t X X \hat{v}$

d'où $\hat{v} = ({}^t X X)^{-1} {}^t X u$

On a donc $\hat{u} = X \hat{v} = X ({}^t X X)^{-1} {}^t X u$

On a donc bien montré que $\Pi_{\text{Im}(X)} = X ({}^t X X)^{-1} {}^t X \quad \square$

Théorème Estimation de paramètres (forme matricielle)

X est une matrice injective $n \times d$

L'estimateur du maximum de vraisemblance du
MLG $Y = X \cdot \beta + \varepsilon$, $\beta \in \mathbb{R}^d$ $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$

$$\text{est } \hat{\beta}_n = ({}^t X X)^{-1} X Y$$

$$\hat{S}_n^2 = \frac{1}{n} \|Y - X \hat{\beta}_n\|^2$$

De plus, $\hat{\beta}_n \perp \hat{S}_n^2$ et $\beta_n \sim \mathcal{N}(\beta, \sigma^2 ({}^t X X)^{-1})$
 $\frac{n \hat{S}_n^2}{\sigma^2} \sim \chi^2(n-d)$

preuve

Le modèle vectoriel correspond au modèle matriciel
avec $V = \text{Im}(X)$

$$\text{donc } \hat{m}_n = X \hat{\beta}_n$$

$$\Pi_V Y = X \hat{\beta}_n$$

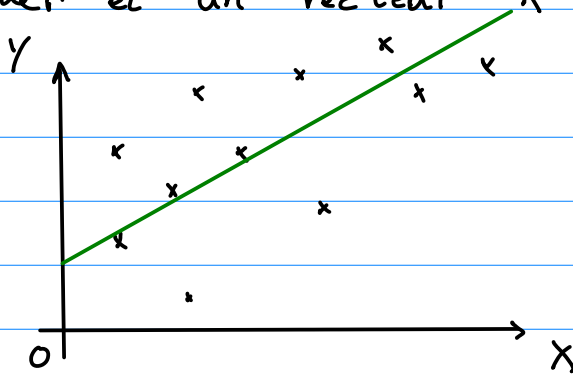
$$\text{(Prop)} \quad X ({}^t X X)^{-1} {}^t X Y = X \hat{\beta}_n$$

$$\text{(injectivité de } X) \quad ({}^t X X)^{-1} {}^t X Y = \hat{\beta}_n$$

$$\text{et } \hat{S}_n^2 = \frac{1}{n} \|Y - \hat{m}_n\|^2 = \frac{1}{n} \|Y - X \hat{\beta}_n\|^2 \quad \square$$

3 - Régression linéaire simple

On suppose qu'on a un vecteur Y de données à expliquer et un vecteur X de données explicatives



On cherche une relation affine entre X et Y ,

c'est-à-dire
$$Y_j = a + b X_j + \varepsilon_j$$

où (ε_j) iid de loi $\mathcal{N}(0, \sigma^2)$

Déterminer P' ETV de (a, b, σ^2)

Observons pour commencer que ce modèle statistique s'écrit $Y = \bar{X} \cdot \beta + \varepsilon$

où $\beta = \begin{pmatrix} a \\ b \end{pmatrix}$, $\bar{X} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$, $\varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} \sim \mathcal{N}(0, \sigma^2 I_n)$

donc, par le théorème précédent, on a

$$\hat{\beta}_n = (\bar{X}^t \bar{X})^{-1} \bar{X}^t Y = \begin{pmatrix} \hat{a}_n \\ \hat{b}_n \end{pmatrix}$$

avec
$$\bar{X}^t \bar{X} = \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix} = n \begin{pmatrix} 1 & \bar{x}_n \\ \bar{x}_n & \overline{x_n^2} \end{pmatrix}$$

$$(\bar{X}^t \bar{X})^{-1} = \frac{1}{n(\overline{x_n^2} - \bar{x}_n^2)} \begin{pmatrix} \overline{x_n^2} & -\bar{x}_n \\ -\bar{x}_n & 1 \end{pmatrix}$$

$$\begin{pmatrix} \bar{X} \\ \bar{Y} \end{pmatrix} = \begin{pmatrix} 1 & \dots & 1 \\ X_1 & \dots & X_n \end{pmatrix} \begin{pmatrix} n\bar{Y}_n \\ X \cdot Y \end{pmatrix} = n \begin{pmatrix} \bar{Y}_n \\ \overline{XY}_n \end{pmatrix}$$

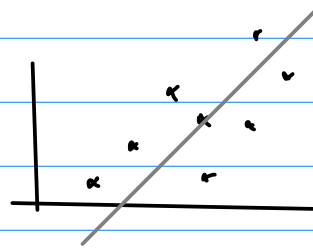
$$\frac{X \cdot Y}{n} = \frac{1}{n} \sum X_i Y_i = \overline{XY}_n$$

donc $\hat{\beta}_n = \frac{1}{\overline{X^2} - \bar{X}_n^2} \begin{pmatrix} \overline{X^2} \bar{Y}_n - \bar{X}_n \overline{XY}_n \\ \overline{XY}_n - \bar{X}_n \bar{Y}_n \end{pmatrix}$

d'où $\hat{a}_n = \frac{\overline{X^2} \bar{Y}_n - \bar{X}_n \overline{XY}_n}{\overline{X^2} - \bar{X}_n^2}$

et $\hat{b}_n = \frac{\overline{XY}_n - \bar{X}_n \bar{Y}_n}{\overline{X^2} - \bar{X}_n^2}$

Remarque Le dénominateur de \hat{a}_n , \hat{b}_n est $\hat{S}^2(X)$, ce dénominateur est d'autant plus grand que les données sont variées



On en déduit l'ETIV de σ^2 comme étant

$$\frac{1}{n} \|Y - \bar{X} \hat{\beta}_n\|^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{a}_n - \hat{b}_n X_i)^2$$

\Rightarrow Cela revient à dire qu'on cherche les valeurs de a et b qui minimisent

$$\sum_{i=1}^n (Y_i - a - b X_i)^2$$

On retombe sur la méthode des moindres carrés

Proposition

L'estimation du maximum de vraisemblance du modèle
 $Y_j = a + b X_j + \varepsilon_j$, $(\varepsilon_j) \text{ iid } \mathcal{N}(0, \sigma^2)$

est donnée par $\hat{b}_n = \frac{\overline{\text{Cov}}(X, Y)}{\overline{\text{Var}}(X)}$, $\hat{a}_n = \bar{Y}_n - \hat{b}_n \bar{X}_n$

$$\text{et } \hat{S}_p^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{a}_n - x_i \hat{b}_n)^2$$

$$\text{où } \overline{\text{Var}}(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2 = \bar{x}_n^2 - \bar{x}_n^2$$

$$\text{et } \overline{\text{Cov}}(X, Y) = \frac{1}{n} \sum (x_i - \bar{x}_n)(y_i - \bar{y}_n) = \bar{x}_n \bar{y}_n - \bar{x}_n \bar{y}_n$$

Remarque

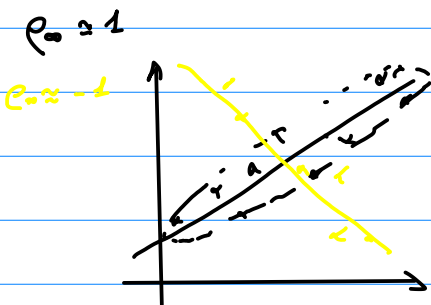
Le coefficient de corrélation de X et Y

$$\rho_n = \frac{\overline{\text{Cov}}(X, Y)}{\sqrt{\overline{\text{Var}}(X) \overline{\text{Var}}(Y)}}$$

$$\rho_n \xrightarrow{n \rightarrow \infty} \rho_0 = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(a + bX + N)}}$$

où $N \sim \mathcal{N}(0, \sigma^2) \perp\!\!\!\perp X$

$$\rho_0 = \frac{b \text{Var}(X)}{\sqrt{\text{Var}(X)(b^2 \text{Var}(X) + \sigma^2)}} = \frac{\text{sgn}(b)}{\sqrt{1 + \frac{\sigma^2}{b^2 \text{Var}(X)}}}$$



⚠ Coefficient de corrélation petit ne correspond pas exactement à relation linéaire inexistante

Test de relation linéaire

$$Y = a + bX + \varepsilon$$

On cherche à tester si l'ordonnée à l'origine a est nulle ou non.

$$H_0 = \{ a = 0 \}$$

$$H_1 = \{ a \neq 0 \}$$

Estimation du MLE

l'hypothèse alternative est

l'existence de données supplémentaires à prendre en compte

L'estimation de a est

$$\hat{a}_n = \frac{\overline{X^2 Y} - \overline{X_n} \times \overline{X Y}_n}{\overline{X^2} - \overline{X_n}^2} \sim \mathcal{N} \left(a, \frac{\sigma^2 \overline{X_n^2}}{n(\overline{X_n^2} - \overline{X_n}^2)} \right)$$

par Théorème d'estimation matriciel

"preuve"

$$\hat{a}_n = \frac{1}{n \overline{\text{Var}}(X)} \sum_{j=1}^n \overline{X_n} Y_j - \overline{X_n} X_j Y_j$$

est combinaison linéaire de coordonnées du vecteur gaussien Y

$$\hat{a}_n = \frac{1}{n \overline{\text{Var}}(X)} \sum_{j=1}^n (a + b X_j + \varepsilon_j) (\overline{X_n^2} - \overline{X_n} X_j)$$

$$= \underbrace{\frac{1}{n \overline{\text{Var}}(X)} \sum_{j=1}^n (a + b X_j) (\overline{X_n^2} - \overline{X_n} X_j)}_c + \underbrace{\frac{1}{n \overline{\text{Var}}(X)} \sum_{j=1}^n \varepsilon_j (\overline{X_n^2} - \overline{X_n} X_j)}_Z$$

$$\text{avec } c = \frac{1}{n \overline{\text{Var}}(X)} \left(a n (\overline{X_n^2} - \overline{X_n}^2) + b n (\overline{X_n^2} \overline{X_n} - \overline{X_n} \overline{X_n^2}) \right) = a$$

$$\text{et } Z \sim \mathcal{N} \left(0, \frac{1}{n^2 \overline{\text{Var}}(X)^2} \sigma^2 \sum_{j=1}^n \overline{X_n^4} - 2 X_j \overline{X_n} \overline{X_n^2} + X_j^2 \overline{X_n^2} \right) \sim \mathcal{N} \left(0, \frac{\sigma^2 \overline{X_n^2}}{n(\overline{X_n^2} - \overline{X_n}^2)} \right)$$

□

Sous H_0 , $\hat{a}_n \sim \mathcal{N}\left(0, \frac{\sigma^2 \bar{X}_n^2}{n \text{Var}(X)}\right)$

\hat{a}_n n'est pas une statistique libre, car sa loi dépend de σ^2

On construit alors une statistique libre avec \hat{S}_n^2

rappelons que $\hat{a}_n \perp \hat{S}_n^2$ et $\frac{n \hat{S}_n^2}{\sigma^2} \sim \chi^2(n-2)$

$$\frac{\chi^2(d)}{\sqrt{\chi^2(d)/d}} \sim \mathcal{O}(d)$$

$$\frac{\hat{a}_n \times \frac{\sqrt{n \text{Var}(X)}}{\sigma^2 \bar{X}_n}}{\sqrt{\frac{\hat{S}_n^2}{\sigma^2} \frac{1}{n-2}}} \sim \mathcal{O}(n-2)$$

$$\frac{\sqrt{\frac{(n-2) \text{Var}(X)}{\bar{X}_n^2}}}{\sqrt{\hat{S}_n^2}} \frac{\hat{a}_n}{\sqrt{\hat{S}_n^2}} \sim \mathcal{O}(n-2) \quad \text{est une statistique libre}$$

On construit donc la statistique de test de niveau α de H_1 contre H_0 par

$$T = \mathbb{1} \left| \frac{\sqrt{\frac{(n-2) \text{Var}(X)}{\bar{X}_n^2}}}{\sqrt{\hat{S}_n^2}} \frac{\hat{a}_n}{\sqrt{\hat{S}_n^2}} \right| \geq \delta$$

en choisissant δ tel que $\mathbb{P}_{H_0}(T=1) = 1-\alpha$

i.e. $\mathbb{P}_{H_0} \left(\left| \frac{\sqrt{\frac{(n-2) \text{Var}(X)}{\bar{X}_n^2}}}{\sqrt{\hat{S}_n^2}} \frac{\hat{a}_n}{\sqrt{\hat{S}_n^2}} \right| \geq \delta \right) = 1-\alpha$

on prend $\delta = t_{1-\frac{1-\alpha}{2}}^{(n-2)}$

et on a $T = \mathbb{1} \left| \frac{\hat{a}_n}{\sqrt{\hat{S}_n^2}} \right| \geq t_{\frac{1-\alpha}{2}}^{(n-2)} \sqrt{\frac{\bar{X}_n^2}{(n-2) \text{Var}(X)}}$

on a toujours intérêt à choisir $\text{Var}(X) \gg 1$

Test de significativité de X

$$Y = a + bX + e$$

$$H_0 = \{ b = 0 \} \text{ contre } H_1 = \{ b \neq 0 \}$$

Ce qu'on teste, c'est si X est effectivement une variable explicative du modèle

$$\hat{b}_n = \frac{\overline{X^2} \times \overline{XY}_n - \overline{X}_n \overline{Y}_n}{\overline{X^2} - \overline{X}_n^2} \sim \mathcal{N}\left(b, \frac{\sigma^2}{\overline{X^2} - \overline{X}_n^2}\right)$$

Construisons le test de niveau $1-\alpha$ de H_1 contre H_0

$$\text{Sous } H_0, \text{ on a } \hat{b}_n \sim \mathcal{N}\left(0, \frac{\sigma^2}{\overline{X^2} - \overline{X}_n^2}\right)$$

$$\text{On observe que } \frac{\hat{b}_n \sqrt{\frac{n \text{Var}(X)}{\sigma^2}}}{\sqrt{\frac{n \hat{S}_n^2}{\sigma^2}} \frac{1}{n-2}} \sim \mathcal{C}(n-2)$$

$\chi^2(n-2)$

La statistique de test est donc

$$\mathbb{1} \left[\sqrt{(n-2) \text{Var}(X)} \left| \frac{\hat{b}_n}{\sqrt{\hat{S}_n^2}} \right| \geq t_{\frac{1+\alpha}{2}}^{(n-2)} \right]$$

$$\text{qu'on peut réécrire } \mathbb{1} \left[\left| \frac{\overline{\text{Cov}}(X, Y)}{\sqrt{\text{Var}(X)} \sqrt{\hat{S}_n^2}} \right| \geq \frac{t_{\frac{1+\alpha}{2}}^{(n-2)}}{\sqrt{n-2}} \right]$$

coefficient de corrélation de X et Y

Prédiction pour une nouvelle donnée $Y = a + bx + \varepsilon$

avec les estimateurs \hat{a}_n , \hat{b}_n et \hat{S}_n^2 calculés

on récupère une nouvelle donnée x , et on souhaite donner un intervalle de confiance sur la quantité y associée

un estimateur naturel de y est

$$\hat{y}_n = \hat{a}_n + \hat{b}_n x$$

de plus $y - \hat{y}_n = a + bx + \varepsilon - \hat{a}_n - \hat{b}_n x$

$$= \underbrace{a - \hat{a}_n}_{\text{erreur sur } a} + \underbrace{(b - \hat{b}_n)}_{\text{erreur sur } b} x + \varepsilon$$

$$\sim \mathcal{N}\left(0, \sigma^2 \left(1 + \frac{(x - \bar{x}_n)^2}{n \text{Var}(x)}\right)\right)$$

On montre que

$$\frac{\hat{y}_n - y}{\sqrt{\hat{S}_n^2 \left(1 + \frac{(x - \bar{x}_n)^2}{\text{Var}(x)}\right)}} \sim \mathcal{L}(n-2)$$

$$\text{donc un IC}_{1-\alpha}(y) = \left[\hat{y}_n \pm t_{1-\frac{\alpha}{2}}^{(n-2)} \sqrt{\hat{S}_n^2 \left(1 + \frac{(x - \bar{x}_n)^2}{\text{Var}(x)}\right)} \right]$$

Exercice Vérifier les calculs ci-dessus

4 - Régions de confiance et tests fondamentaux

a. Intervalle de confiance et tests pour la variance

Sous le modèle linéaire gaussien, l'estimateur du maximum de vraisemblance de σ^2 est

$$\hat{\sigma}_n^2 = \frac{1}{n} \|Y - \Pi_V Y\|_2^2 = \frac{1}{n} \|Y - X\hat{\beta}_n\|_2^2$$

et une version débiaisée de cet estimateur est

$$\hat{\sigma}_n^2 = \frac{1}{n - \dim(V)} \|Y - \Pi_V Y\|_2^2 = \frac{1}{n - \text{Rg}(X)} \|Y - X\hat{\beta}_n\|_2^2$$

On sait que $\hat{\sigma}_n^2 \perp (\hat{m}_n, \hat{\beta}_n)$

$$\text{et que } \frac{n-d}{\sigma^2} \hat{\sigma}_n^2 \sim \chi^2(n-d)$$

$$\text{où } d = \dim(V) = \text{Rg}(X)$$

Par conséquent, l'intervalle de confiance unilatère à droite de niveau $1-\alpha$ pour σ^2 est

$$\text{IC}_{1-\alpha}(\sigma^2) = \left[0, \frac{(n-d)}{\chi_{\alpha}^2(n-d)} \hat{\sigma}_n^2 \right]$$

Une statistique de test pour $H_0: \{\sigma > \sigma_{\text{ref}}\}$ contre $H_1: \{\sigma \leq \sigma_{\text{ref}}\}$ est donnée par

$$T = \mathbb{1}_{\hat{\sigma}_n^2 < \delta} \quad \text{avec } \delta \text{ tel que}$$
$$\sup_{m, \sigma \in H_0} \mathbb{P}_{m, \sigma} (T=1) < \alpha$$
$$= \mathbb{P}_{m, \sigma_{\text{ref}}} (T=1) < \alpha$$

d'où on obtient le test de niveau α

$$T = \mathbb{1} \left\{ \hat{\sigma}_n^2 \leq \frac{\sigma_{ref}^2}{n-p} \chi_{\alpha}^{2(n-p)} \right\}$$

$$\text{En effet } \mathbb{P}_{n, \sigma_{ref}^2} \left(\hat{\sigma}_n^2 \leq \lambda \right) = \mathbb{P}_{n, \sigma_{ref}^2} \left(\frac{\hat{\sigma}_n^2 (n-d)}{\sigma_{ref}^2} \leq \frac{\lambda (n-d)}{\sigma_{ref}^2} \right)$$

$$= \mathbb{P} \left(\chi_{\alpha}^{2(n-d)} \leq \frac{\lambda (n-d)}{\sigma_{ref}^2} \right) = \alpha$$

$$\Rightarrow \frac{\lambda (n-d)}{\sigma_{ref}^2} = \chi_{\alpha}^{2(n-d)} \Rightarrow \lambda = \frac{\sigma_{ref}^2}{n-d} \chi_{\alpha}^{2(n-d)}$$

b. Test d'une relation affine

On se place dans le cadre du modèle matriciel

Pour $c \in \mathbb{R}^d$ et $a \in \mathbb{R}$, on cherche à tester l'hypothèse $H_0 : \{ c \cdot \beta = a \}$ contre $H_1 : \{ c \cdot \beta \neq a \}$

$$\text{où } c \cdot \beta = c_1 \beta_1 + c_2 \beta_2 + \dots + c_d \beta_d$$

Exemple

On peut souhaiter tester $\beta_1 = \beta_2$

On peut souhaiter tester l'hypothèse $\beta_j = 0$, c'est à dire la significativité de la donnée X_j .

La statistique de test naturelle à considérer

$$\text{est la variable } T_n = \frac{c \cdot \hat{\beta}_n - a}{\sqrt{\hat{\sigma}_n^2 c' (X'X)^{-1} c}}$$

$$\text{En effet, sous } H_0, \quad c \cdot \hat{\beta}_n - a \sim \mathcal{N} \left(\underbrace{c \cdot \beta - a}_0, \sigma^2 c' (X'X)^{-1} c \right)$$

En utilisant que $\frac{(n-d)\hat{\sigma}_n^2}{\sigma^2} \sim \chi^2(n-d) \perp\!\!\!\perp c\hat{\beta}_{n-a}$

$$\text{donc } \frac{c \cdot \hat{\beta}_{n-a}}{\sqrt{\hat{\sigma}_n^2 t_c (X'X)^{-1} c}} = \frac{c \cdot \hat{\beta}_{n-a}}{\sqrt{\sigma^2 t_c (X'X)^{-1} c}} \frac{1}{\sqrt{\frac{(n-d)\hat{\sigma}_n^2}{\sigma^2} / (n-d)}}$$
$$\sim \mathcal{N}(0, 1) \quad \chi^2(n-d)$$

$$\sim \mathcal{C}(n-d)$$

Donc T_n est une statistique libre du modèle

On propose donc le test de niveau α
de H_0 contre H_1

$$\parallel |T_n| > t_{1-\frac{\alpha}{2}}^{(n-d)}$$

$$\begin{aligned} \text{car } \mathbb{P}_{H_0}(|T_n| > t_{1-\frac{\alpha}{2}}^{(n-d)}) &= \mathbb{P}(T_n > t_{1-\frac{\alpha}{2}}^{(n-d)}) \\ &\quad + \mathbb{P}(T_n < -t_{1-\frac{\alpha}{2}}^{(n-d)}) \\ &= \frac{\alpha}{2} + \frac{\alpha}{2} = \alpha \end{aligned}$$

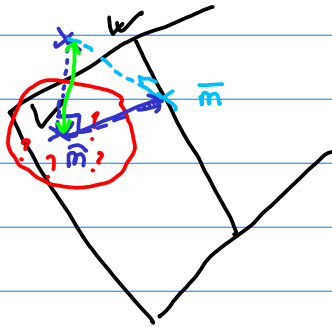
Avec les mêmes méthodes, on construit

$$IC_{1-\alpha}(c, \beta) = \left[c \cdot \hat{\beta}_n \pm t_{1-\frac{\alpha}{2}}^{(n-d)} \hat{\sigma}_n \sqrt{t_c (X'X)^{-1} c} \right]$$

c. Test de Fisher d'un sous-modèle

On se place dans le modèle vectoriel $Y = m + \varepsilon$
 où $m \in V \subset \mathbb{R}^n$ et $\varepsilon \sim N(0, \sigma^2 I_n)$

On se donne un sous-espace vectoriel W de V
 et on souhaite tester $H_0 : \{m \in W\}$ contre $H_1 : \{m \notin W\}$



$$\frac{\|\text{green vector}\|^2}{\sigma^2} \sim \chi^2(n - \dim V)$$

$$\frac{\|\text{red vector}\|^2}{\sigma^2} \sim \chi^2(\dim V - \dim W)$$

Définition

La loi de Fisher à k et p degrés de liberté notée $F(k, p)$ est la loi de $Z = \frac{U/k}{V/p}$, où

$$U \sim \chi^2(k) \quad V \sim \chi^2(p) \quad \text{et} \quad U \perp V$$

Remarque

$$\text{On a } \lim_{p \rightarrow \infty} Z_{k,p} = U/k$$

$$\text{donc pour } p \gg 1, \quad F(k, p) \approx \frac{\chi^2(k)}{k}$$

Théorème

$$d = \dim(V) \quad d' = \dim(W)$$

Lorsque $m \in W$, la statistique

$$F = \frac{\|\pi_W Y - \pi_V Y\|^2 / (d - d')}{\|Y - \pi_V Y\|^2 / (n - d)} \sim F(d - d', n - d) \quad \text{et est indépendante de } \pi_W Y$$

preuve Par théorème de Cochran
on a $\pi_W Y$, $\pi_{W^\perp} Y$ et $\pi_{V^\perp} Y$ qui

sont des v.a. gaussiens indépendants

de plus $\pi_W Y \sim \mathcal{N}(m, \sigma^2 I_W)$

$$\pi_{W^\perp} Y \sim \mathcal{N}(0, \sigma^2 I_{W^\perp})$$

$$\pi_{V^\perp} Y \sim \mathcal{N}(0, \sigma^2 I_{V^\perp})$$

$$\text{Or } \pi_V Y - \pi_W Y = \pi_{W^\perp} Y$$

$$Y - \pi_V Y = \pi_{V^\perp} Y$$

$$\text{Donc } F = \frac{\|\pi_{W^\perp} Y\|^2 / \dim(W^\perp)}{\|\pi_{V^\perp} Y\|^2 / \dim(V^\perp)} \sim F(\dim(W^\perp), \dim(V^\perp))$$

et est \perp de $\pi_W Y$

On en déduit le test suivant de niveau α
de $H_0 = \{m \in W\}$ contre $H_1 = \{m \in V \setminus W\}$

$$\|F > \delta_{1-\alpha}^{(d-d', n-d)}$$

Remarque

Si $d-d' = 1$, i.e. W est un hyperplan de V
on retrouve le test de Student bilatéral

$$T = \frac{N}{\sqrt{k_0/k_1}} \quad \text{donc } T^2 = \frac{N^2}{k_0/k_1} \sim F(1, k)$$

d. Test de Wald de plusieurs relations affines

Considérons un TLO $Y = X\beta + \varepsilon$

On teste l'hypothèse $H_0 = \{c\beta = a\}$
contre H_1

où $c \in M_{d', d}$ injective et $a \in \mathbb{R}^{d'}$

Exemple $c = \begin{bmatrix} -2 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ et $a = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

dans ce cas on teste sous H_0 si $2\beta_1 = \beta_2$
et $\beta_3 = 0$

Rappelons que ${}^t X X$ est symétrique définie positive
on va noter $\Gamma^2 = ({}^t X X)^{-1}$

donc $c ({}^t X X)^{-1} c^t = ({}^t \Gamma^2 c) ({}^t \Gamma^2 c)$ est symétrique

où $W = \{ \beta : c\beta = a \}$ ← définie positive
 $\| \Pi_W Y - \underbrace{\Pi_V Y}_{X\hat{\beta}_n} \|^2$

Théorème

Sous H_0 ,

$$W = \frac{[(c\hat{\beta}_n - a) c ({}^t X X)^{-1} c^t (c\hat{\beta}_n - a)] / d'}{\| Y - X\hat{\beta}_n \|^2 / (n-d)}$$

suit une loi $F(d', n-d)$ de degrés de liberté

preuve
Sous H_0 , $\hat{\beta}_n \sim \mathcal{N}(\beta, \sigma^2 (X'X)^{-1})$

donc $c\hat{\beta}_n - a \sim \mathcal{N}(\underbrace{c\beta - a}_0, \sigma^2 c'(X'X)^{-1}c)$

or $c'(X'X)^{-1}c$ est symétrique définie positive

il existe Δ telle que $\Delta^2 = c'(X'X)^{-1}c$

on obtient donc $\Delta^{-1}(c\hat{\beta}_n - a) \sim \mathcal{N}(0, \sigma^2 I_{d'})$

$\frac{1}{\sigma^2} (c\hat{\beta}_n - a)' c'(X'X)^{-1} c (c\hat{\beta}_n - a) = \frac{1}{\sigma^2} \|\Delta^{-1}(c\hat{\beta}_n - a)\|^2$
 $\sim \chi^2(k)$

$\hat{\beta}_n$ étant indépendant de $\sum_n \hat{\epsilon}_n^2 = \|Y - X\hat{\beta}_n\|^2/n$

on obtient $\frac{\frac{1}{\sigma^2} (c\hat{\beta}_n - a)' c'(X'X)^{-1} c (c\hat{\beta}_n - a) / d'}{\frac{1}{\sigma^2} \|Y - X\hat{\beta}_n\|^2 / (n-d)}$
 $\sim F(d', n-d)$

Le test de $H_0: \{c\beta = a\}$ contre H_1 s'écrit alors

$$\{W \geq f_{1-\alpha}^{(d', n-d)}\}$$

c'est un test de niveau α

Le test de Wald permet la création de régions de confiance pour $\hat{\beta}_n$

Pour tout c matrice injective, on peut définir

$$\mathcal{E}_c = \left\{ a \in \mathbb{R}^d : \frac{(c\hat{\beta}_n - a) c (X'X)^{-1} c' (c\hat{\beta}_n - a) / d}{\|Y - X\hat{\beta}_n\|^2 / (n-d)} \leq F_{d, n-d} \right\}$$

est une région de confiance de niveau exactement $1 - \alpha$ pour $c \cdot \beta$

\mathcal{E}_c est un ellipsoïde centré en $c \hat{\beta}_n$

En particulier $c = I_d$

$$\mathcal{E} = \left\{ a \in \mathbb{R}^d : \frac{(\hat{\beta}_n - a) (X'X)^{-1} (\hat{\beta}_n - a) / d}{\|Y - X\hat{\beta}_n\|^2 / (n-d)} \leq F_{d, n-d} \right\}$$

est une région de confiance de β

5 - Variables quantitatives: La régression linéaire multiple

On suppose qu'on souhaite expliquer une variable Y à partir de plusieurs variables quantitatives X_1, \dots, X_d .

Définition

Une variable quantitative est une variable dont la valeur numérique a un sens: taille, aire, volume, température,

Dans cette situation, on pose

$$Y_u = \mu + \beta_1 X_1 + \dots + \beta_d X_d$$

donc les paramètres sont $(\mu, \beta_1, \dots, \beta_d) \in \mathbb{R}^{d+1}$

La matrice des données $X = \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & x_{n,1} & \dots & x_{n,d} \end{pmatrix}$

On appelle $(x_{1,1}, \dots, x_{1,d})$ les régresseurs et $(\mu, \beta_1, \dots, \beta_d)$ les coefficients de régression.

Par exemple, on peut utiliser un test de Student pour déterminer la significativité de la j^e variable X_j

$$H_0 = \{ \beta_j = 0 \} \quad \text{contre} \quad H_1$$

On se fixe $d' \leq d$, et on teste

$H_0 = \{ \beta_1 = \dots = \beta_{d'} = 0 \}$ contre H_1 : il existe un coefficient non nul

$$c = \begin{pmatrix} 1 & 1 & 0 & \dots & 0 \\ 0 & 1 & 1 & \dots & 0 \end{pmatrix} \quad a = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$X_{\text{ref}} = \begin{pmatrix} 1 & X_{1,d'} & \dots & X_{1,d} \\ \vdots & \vdots & \dots & \vdots \\ 1 & X_{n,d'} & \dots & X_{n,d} \end{pmatrix}$$

La matrice des données significatives sous H_0

On pose $V = \{ X \cdot \beta, \beta \in \mathbb{R}^d \} = \text{Im}(X)$

$W = \{ X_{\text{ref}} \cdot \beta, \beta \in \mathbb{R}^{d-d'} \} = \text{Im}(X_{\text{ref}})$

$H_0 = \{ X \cdot \beta \in W \}$

$H_1 = \{ X \cdot \beta \notin W \}$

donc la statistique de test sera

$$F = \frac{\| \pi_V Y - \pi_W Y \|^2 / d'}{\| Y - \pi_V Y \|^2 / (n-d)}$$

$$= \frac{\| X (X^T X)^{-1} X^T Y - X_{\text{ref}} (X_{\text{ref}}^T X_{\text{ref}})^{-1} X_{\text{ref}}^T Y \|^2 / d'}{\| Y - X (X^T X)^{-1} X^T Y \|^2 / (n-d)} \sim F(d', n-d)$$

On construit le test $\mathbb{1}_{\{ F \geq f_{1-\alpha}^{(d', n-d)} \}}$

6. Application au cas des variables qualitatives : l'analyse de variance ANOVA

On s'intéresse à un modèle linéaire gaussien dont les données sont des variables qualitatives avec un nombre fini de modalités

Définition

Une variable qualitative est une variable ne permettant pas de traduction chiffrée immédiate. On appelle modalités l'ensemble des valeurs possibles de cette variable

a. ANOVA à un facteur

On considère un modèle linéaire gaussien où X prend des valeurs dans un ensemble I avec $|I| \geq 2$

Le modèle d'analyse de variance s'écrit

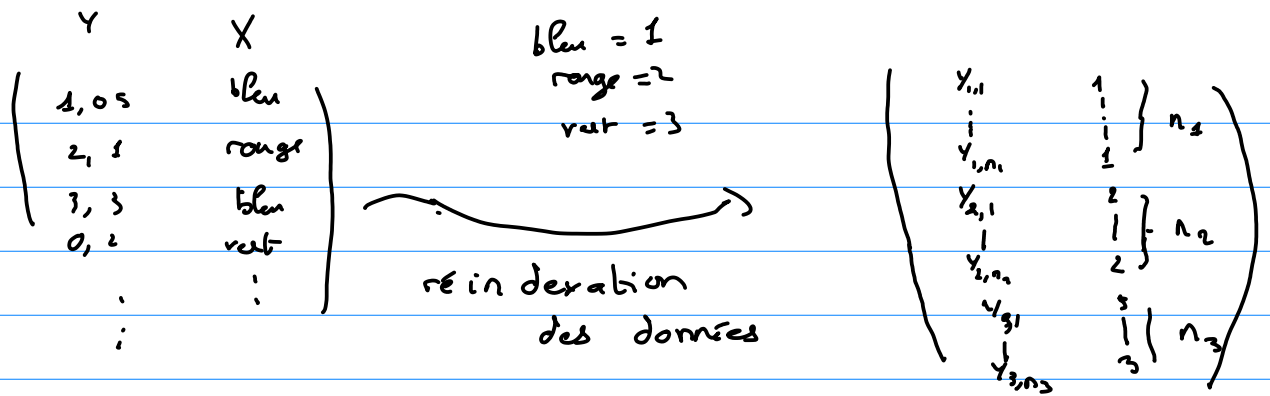
$$Y_k = \mu + a_i + \varepsilon_k \quad \text{si } X_k = i$$

où $\varepsilon_k \sim \mathcal{N}(0, \sigma^2)$ iid

Quitte à renumérotter, on exprime le modèle sous la forme

$$\left\{ \begin{array}{l} Y_{i,k} = \mu + a_i + \varepsilon_{i,k} \\ i \in \{1, 2, \dots, |I|\} \quad k \in \{1, \dots, n_i\} \end{array} \right. \quad \varepsilon_{i,k} \sim \mathcal{N}(0, \sigma^2) \text{ iid}$$

n_i : nombre de données avec la modalité i et $n = n_1 + \dots + n_{|I|}$



La matrice des données de ce ALG s'écrit

$$X = \begin{pmatrix} 1 & 1 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 0 \end{pmatrix} \quad \text{et} \quad \beta = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_I \end{pmatrix}$$

cette matrice n'est pas injective, le modèle est surparamétrisé. On imposera donc $\mu = 0$, pour retrouver une formulation matricielle

$$\begin{cases} Y_{i,k} = a_i + \varepsilon_{i,k} & \varepsilon_{i,k} \text{ iid } \mathcal{N}(0, \sigma^2) \\ i \leq I, \quad k \leq n_i \end{cases}$$

Les estimateurs de $\beta = (a_1, \dots, a_I)$ et σ^2

$${}^t X X = \begin{pmatrix} n_1 & 0 & \dots & 0 \\ 0 & n_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & n_I \end{pmatrix}$$

donc $\hat{\beta}_n = \left(\frac{\bar{Y}_i}{n_i}, 1 \leq i \leq I \right)$

où $\bar{Y}_i = \frac{Y_{i,1} + \dots + Y_{i,n_i}}{n_i}$

$$\hat{\beta}_n \sim \mathcal{N} \left(\beta, \sigma^2 \begin{pmatrix} \frac{1}{n_1} & & & 0 \\ & \frac{1}{n_2} & & 0 \\ & & \ddots & 0 \\ 0 & & & \frac{1}{n_I} \end{pmatrix} \right)$$

$$\hat{\sigma}_n^2 = \frac{1}{n-I} \sum_{i=1}^I \sum_{k=1}^{n_i} (Y_{i,k} - \bar{Y}_i)^2$$

$$\frac{\hat{\sigma}_n^2 (n-I)}{\sigma^2} \sim \chi^2 (n-I)$$

Le test statistique classique est le test de pertinence de la factorisation

$$H_0 = \{ a_1 = a_2 = \dots = a_I \} \quad \text{contre } H_1 := \{ \exists i \neq j \text{ avec } a_i \neq a_j \}$$

On pose c la projection orthogonale sur $\begin{pmatrix} 1 \\ 1 \\ \dots \\ 1 \end{pmatrix}$

$$H_0 = \{ c\beta = 0 \} \quad \text{contre } H_1$$

Sous H_0 , la projection de Y sur $\begin{pmatrix} 1 \\ 1 \\ \dots \\ 1 \end{pmatrix}$ est $\frac{1}{n} \sum_{i=1}^I \sum_{k=1}^{n_i} Y_{ik} = \bar{Y}$.

Dans ce cas, on a

$$F = \frac{\| X(X'X)^{-1} X'Y - X_0(X_0'X_0)^{-1} X_0'Y \|^2 / (I-1)}{\| Y - X(X'X)^{-1} X'Y \|^2 / (n-I)}$$

$$= \frac{n-I}{I-1} \frac{\sum_{i=1}^I n_i (\bar{Y}_i - \bar{Y}_0)^2}{\sum_{i=1}^I \sum_{k=1}^{n_i} (Y_{ik} - \bar{Y}_{i\cdot})^2}$$

qui suit sous H_0 la loi $F(I-1, n-I)$

Le test est alors $\| F \geq \lambda$

b. Extension à l'ANOVA à deux facteurs

Si on a plusieurs modalités, on écrit le modèle sous la forme $I \times J$

$$\begin{cases} Y_{i,j,k} = a_{i,j} + \varepsilon_{i,j,k} \\ 1 \leq i \leq I, 1 \leq j \leq J \text{ et } k \leq n_{i,j} \end{cases}$$

On peut se poser la question

$$H_0 = \{ a_{i,1} = a_{i,2} = \dots = a_{i,I} \text{ pour tout } i \in I \}$$

VLG à 2 paramètres $I \times J$

$$Y_{ij,u} \sim N(c_{ij} + \epsilon)$$

anova

anova X

VLG à 1 paramètre I

$$Y_{i,u} \sim N(a_i + \epsilon)$$

anova

VLG à 1 paramètre J

$$Y_{j,u} \sim N(b_j + \epsilon)$$

anova

modèle gaussien
simple

$$Y \sim N(\mu + \epsilon)$$

7 - Analyse des résidus

Quand on étudie un modèle linéaire gaussien, on fait l'hypothèse que $Y \sim N(m, \sigma^2 I_n)$

Il est souvent important de tester cette hypothèse par l'analyse des résidus

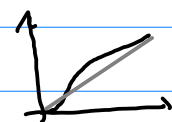
$$\hat{\epsilon} = Y - \hat{m}_n = Y - X\hat{\beta}_n$$

• $E(\hat{\epsilon}) = 0$? $\overline{\hat{\epsilon}}$ comparé à 0

• $\text{Var}(\hat{\epsilon}) = \sigma^2$ $\overline{\text{Var}(\hat{\epsilon})}$ comparé à $\hat{\sigma}_n^2$

• Variables qualitatives tester l'indépendance de ϵ

• Loi de ϵ est normale ?



q-q plot

III Estimation de densité

On a un n -échantillon (x_1, \dots, x_n) d'une variable aléatoire et on cherche à déterminer la loi de X

1. Consistance de la fonction de répartition empirique

On suppose que la loi de (x_1, \dots, x_n) admet une densité par rapport à la mesure de Lebesgue

On pose $F: x \mapsto P(X_1 \leq x)$ leur fonction de répartition

Rappelons que F est continue, croissante et

$$\lim_{x \rightarrow \infty} F(x) = 1 \quad \lim_{x \rightarrow -\infty} F(x) = 0$$

Définition

La fonction de répartition empirique de $X = (x_1, \dots, x_n)$ est la v.a. définie par

$$F_n: x \mapsto \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq x} \quad : \quad \begin{array}{l} \text{La proportion de données} \\ \text{inférieures ou égales} \\ \text{à } x \end{array}$$

Remarque

F_n est une fonction en escaliers continue à droite avec des limites à gauche en tout point, dont les discontinuités sont aux points

$$x^{(1)} \leq x^{(2)} \leq \dots \leq x^{(n)}$$

où $(x^{(1)} \dots x^{(n)})$ est la réorganisation de (x_1, \dots, x_n) dans l'ordre croissant

$$F_n(x) = \frac{j}{n} \quad \text{par} \quad x^{(j)} \leq x \leq x^{(j+1)}$$

Proposition

On a $\lim_{n \rightarrow \infty} F_n(x) = F(x)$ p.s. pour tout $x \in \mathbb{R}$

Autrement dit, $F_n(x)$ est un estimateur fortement consistant de $F(x)$ qui est asymptotiquement normal

$$\lim_{n \rightarrow \infty} \sqrt{n}(F_n(x) - F(x)) = \mathcal{N}\left(0, F(x)(1-F(x))\right) \text{ en } \mathbb{P}_i$$

preuve

$$\text{Soit } x \in \mathbb{R}, F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq x}$$

et $(\mathbb{1}_{x_i \leq x}, i \in \mathbb{N})$ est une suite de v.a. iid

de loi Bernoulli $(F(x))$

Par LLN , on a $F_n(x) \rightarrow F(x)$ p.s.

Par TCL on a $\sqrt{n}(F_n(x) - F(x)) \rightsquigarrow \mathcal{N}(0, F(x)(1-F(x)))$
en \mathbb{P}_i

Corollaire **Théorème de Glivenko-Cantelli**

F_n converge vers F uniformément

preuve Théorème de Dini

Définition

Les quantiles empiriques de X sont les quantités définies par $F_n^{(-1)}(q) = \inf \{x \in \mathbb{R} : F_n(x) \geq q\}$

$$= X^{(Lq_n)}$$

Propriétés

Les quantiles empiriques sont des estimateurs fortement consistants des quantiles de la loi de X

2- Test d'ajustement à une loi ou à une famille de lois

a. Test d'ajustement à une loi unique

On teste l'hypothèse $H_0 : \{X \text{ a pour fdr } F\}$
contre H_1

On a un estimateur fortement consistant de la fdr de X , un test naturel est de la forme

$$h(X) = \mathbb{1} \left\{ \|F_n - F\|_\infty > \delta \right\}$$

On remarque que F_n est une fonction en escaliers

$$\|F_n - F\|_\infty = \sup_{x \in \mathbb{R}} |F_n(x) - F(x)|$$

$$= \max_{1 \leq i \leq n} \max \left(\left| \frac{i-1}{n} - F(x^{(i)}) \right|, \left| \frac{i}{n} - F(x^{(i)}) \right| \right)$$

$$\begin{array}{ll} \text{Sous } H_1 & \liminf \left\| F_n - F \right\|_\infty > 0 \\ \text{sous } H_0 & \lim_{n \rightarrow \infty} \left\| F_n - F \right\|_\infty = 0 \end{array}$$

On souhaite fixer le seuil δ tel que

$$\mathbb{P}_{H_0} \left(\|F_n - F\|_\infty > \delta \right) = 1 - \alpha$$

Théorème

La variable aléatoire $\|F_n - F\|_\infty$ ne dépend pas de la valeur de F (c'est une statistique libre)

preuve

Soit (U_1, \dots, U_n) n var. iid de la loi Uniforme sur $[0, 1]$

$(F^{-1}(U_1), F^{-1}(U_2), \dots, F^{-1}(U_n))$ est un n -échantillon

de fonction de répartition F

$$\mathbb{P}(F^{-1}(U) \leq x) = \mathbb{P}(U \leq F(x)) = F(x)$$

Des lors

$$\|F_n - F\|_\infty = \max_{1 \leq i \leq n} \max \left(\left| \frac{i-1}{n} - F(F^{-1}(U^{(i)})) \right|, \left| \frac{i}{n} - F(F^{-1}(U^{(i)})) \right| \right)$$

$$= \max_{1 \leq i \leq n} \max \left(\left| \frac{i-1}{n} - U^{(i)} \right|, \left| \frac{i}{n} - U^{(i)} \right| \right)$$

ne dépend pas de F

□

Cela permet donc de définir des tables de la

Lemme Soit F_n la fdr empirique d'un échantillon de la loi uniforme sur $[0, 1]$

On a $\lim_{n \rightarrow \infty} (\sqrt{n} (F_n(x) - x), x \in [0, 1]) = (B_x, x \in [0, 1])$

où B est un pont brownien (mouvement brownien conditionné à $B_1 = 0$)

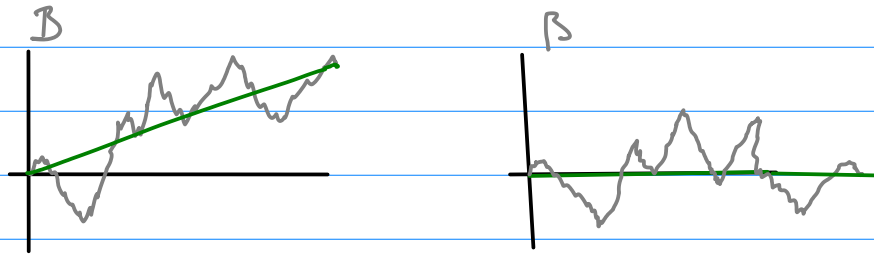
Corollaire $\lim_{n \rightarrow \infty} \mathbb{P}(\sqrt{n} \|F_n - F\|_\infty \geq \lambda) = \mathbb{P}(\sup_{x \in [0, 1]} |B_x| \geq \lambda)$

$$= 2 \sum_{k=1}^{\infty} (-1)^{k+1} e^{-2k^2 \lambda^2}$$

Remarque

Soit $(B_t, t \geq 0)$ un mouvement Brownien

On pose $\beta_s = B_s - sB_1$, pour tout $s \in [0,1]$



Calculons $\text{Cov}(B_s - sB_1, B_1)$

$$\text{Cov}(B_s, B_1) - s \text{Cov}(B_1, B_1) = s - s = 0$$

$$\text{Cov}(B_t, B_s) = \min(s, t)$$

$$s \leq t$$

$$B_t - B_s \perp B_s$$

$$B_t = \underline{B_s} + B_t - B_s$$

Donc comme $\text{Cov}(B_s, B_1) = 0$, $\beta_s \perp B_1$

Calculons $\text{Cov}(\beta_s, \beta_u) = \text{Cov}(B_s - sB_1, B_u - uB_1)$

$$= \text{Cov}(B_s, B_u) - s \text{Cov}(B_1, B_u) - u \text{Cov}(B_s, B_1) + us \text{Cov}(B_1, B_1)$$

$$= \min(s, u) - su.$$

Le pont brownien a \hat{m} la que $(\beta_s, s \in [0,1])$

et est un processus gaussien de matrice de covariances
 $(\min(s, u) - su, s, u \in [0,1])$

preuve de Lemme

Soit (U_1, \dots, U_n) un n -échantillon de \mathcal{U} uniforme sur $[0, 1]$
et F_n la fonction de répartition empirique de cet échantillon
 $(\sqrt{n} (F_n(x) - x), x \in [0, 1]) \xrightarrow[n \rightarrow \infty]{\text{en loi}} (\beta_x, x \in [0, 1])$

Pour tout x fixe, $\sqrt{n} (F_n(x) - x) \xrightarrow[n \rightarrow \infty]{\text{en loi}} \mathcal{N}(0, x(1-x))$

et $x(1-x)$ est bien la variance de β_x

Soit $(x_1, \dots, x_k) \in [0, 1]^k$, supposons $x_1 < x_2 < \dots < x_k$

On montre $(\sqrt{n} (F_n(x_j) - x_j), j \leq k) \xrightarrow[n \rightarrow \infty]{\text{en loi}} (\beta_{x_j}, j \leq k)$

par Théorème Central limite multidimensionnel

$$\sqrt{n} \times \frac{1}{n} \sum_{p=1}^n \left(\begin{pmatrix} \mathbb{1}_{U_p \leq x_1} \\ \mathbb{1}_{U_p \leq x_2} \\ \vdots \\ \mathbb{1}_{U_p \leq x_k} \end{pmatrix} - \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{pmatrix} \right)$$

$$= \frac{1}{\sqrt{n}} \sum_{p=1}^n \left(\begin{pmatrix} Y_1(p) \\ \vdots \\ Y_k(p) \end{pmatrix} - \begin{pmatrix} \bar{Y}_1(p) \\ \vdots \\ \bar{Y}_k(p) \end{pmatrix} \right)$$

$$\xrightarrow{\text{si}} \mathcal{N}(0, \Sigma)$$

$$\text{or } \Sigma_{ij} = \text{Cov}(Y_i(l), Y_j(l))$$

$$= \text{Cov}(\mathbb{1}_{U_i \leq x_i}, \mathbb{1}_{U_i \leq x_j})$$

$$= \mathbb{E}(\mathbb{1}_{U_i \leq x_i, U_i \leq x_j}) - \mathbb{E}(\mathbb{1}_{U_i \leq x_i}) \mathbb{E}(\mathbb{1}_{U_i \leq x_j})$$

$$\Sigma_{i,j} = \min(x_i, x_j) - x_i x_j$$

Or la matrice de covariance de β est

$$\begin{pmatrix} \beta_{x_1} \\ \vdots \\ \beta_{x_n} \end{pmatrix} \sim \mathcal{N}(0, \Sigma)$$

On a donc bien $(\sqrt{n} (F_n(x_j) - x_j), j \leq k) \rightarrow (\beta_{x_j}, j \leq k)$

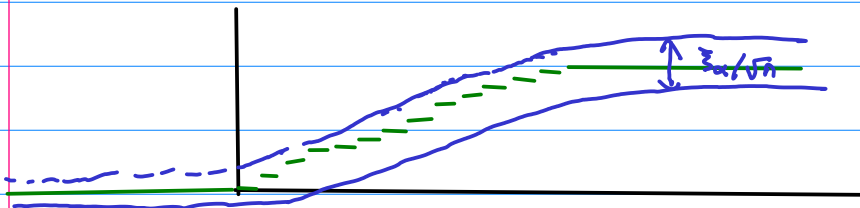
En prenant le $\gg 1$, $x_j = \frac{j}{k}$

Théorème Test de Kolmogorov

Soit $\alpha \in]0,1[$ on pose ξ_α le quantile d'ordre α de la variable $\sup_{t \in [0,1]} |B_t|$, le test $\|F_n - F\|_\infty > \frac{\xi_\alpha}{\sqrt{n}}$

est un test asymptotique de niveau α de $H_0 = \{F \text{ est la f.d.r.}\}$ contre H_1

De plus, on peut construire des intervalles de confiance pour la fonction de répartition F



Remarque Si X prend un nombre fini K de valeurs, on utilise un test du χ^2 pour tester l'adéquation à une loi

$$\text{Toutefois } (\sqrt{n} (F_n(j) - \sum_{i=1}^j p_i), j \leq K) \rightarrow (\beta_{\sum_{i=1}^j p_i}, j \leq K)$$

donc le test de Kolmogorov reste applicable

b. Adéquation à une famille de lois

$X = (X_1, \dots, X_n)$ un n -échantillon

On souhaite tester l'hypothèse, par exemple
 $H_0 = \{ X \text{ est de loi exponentielle} \}$
contre H_1

De façon générale, on teste $H_0: X \text{ est de } P_\theta, \theta \in \Theta$
contre H_1

Une façon de réaliser ce test est de déterminer
un estimateur $\hat{\theta}_n$ de θ , et d'utiliser la statistique

$$\mathbb{1}_{\|F_n - F_{\hat{\theta}_n}\| \geq \lambda}$$

Mais cette statistique n'est pas en général libre, on
utilise des estimations asymptotiques ad-hoc.

Toutefois, il existe deux familles de lois pour lesquelles
il est possible d'obtenir des statistiques libres

→ Ajustement à une famille de lois exponentielles

$$\mathcal{E} = \{ F_\lambda : x \mapsto (1 - e^{-\lambda x}) \mathbb{1}_{x \geq 0}, \lambda > 0 \}$$

L'estimateur de maximum de vraisemblance de λ

$$f(x_1, \dots, x_n, \lambda) = \lambda^n e^{-\lambda(x_1 + \dots + x_n)}$$

$$\log f(x_1, \dots, x_n; \lambda) = n \log \lambda - \lambda(x_1 + \dots + x_n)$$

$$\text{est } \hat{\lambda}_n = \frac{1}{\bar{x}_n}$$

Lemme

La variable $\sup_{x>0} |F_n(x) - F_{\frac{1}{x_n}}(x)|$ est une statistique libre du paramètre λ

preuve Si (x_1, \dots, x_n) est un n. échantillon de Poi Exp (λ), on peut écrire $X_j = \frac{e_j}{\lambda}$ où (e_1, \dots, e_n) n-échantillon de Poi Exp (1)

$$\bar{X}_n = \frac{\bar{e}_n}{\lambda}, \text{ et}$$

$$\sup_{x>0} |F_n(x) - F_{\frac{1}{\bar{X}_n}}(x)| = \max_{j \leq n} \max \left| \frac{j-1}{n} - F_{\frac{1}{\bar{X}_n}}(x^{(j)}) \right|, \left| F_{\frac{1}{\bar{X}_n}}(x^{(j)}) - \frac{j}{n} \right|$$

$$= \max_{j \leq n} \max \left| \frac{j-1}{n} - F_{\frac{\lambda}{\bar{e}_n}}\left(\frac{e^{(j)}}{\lambda}\right) \right|, \left| F_{\frac{\lambda}{\bar{e}_n}}\left(\frac{e^{(j)}}{\lambda}\right) - \frac{j}{n} \right|$$

$$\text{or } F_{\frac{\lambda}{\bar{e}_n}}\left(\frac{e^{(j)}}{\lambda}\right) = 1 - e^{-\frac{\lambda}{\bar{e}_n} \left(\frac{e^{(j)}}{\lambda}\right)} = F_{\frac{1}{\bar{e}_n}}(e^{(j)})$$

$$= \max_{j \leq n} \max \left| \frac{j-1}{n} - F_{\frac{1}{\bar{e}_n}}(e^{(j)}) \right|, \left| F_{\frac{1}{\bar{e}_n}}(e^{(j)}) - \frac{j}{n} \right|$$

qui ne dépend pas de λ . \square

On peut donc construire des test exact d'a déquation à la famille de Poi exponentielles

→ Test d'ajustement à la loi normale

$H_0 = \{ (X_1, \dots, X_n) \text{ n-échantillon de loi normale} \}$ contre H_1

On a également la statistique libre donnée par

$$\sup_{x \in \mathbb{R}} \left| F_n(x) - \Phi \left(\frac{x - \bar{X}_n}{\sqrt{S_n^2}} \right) \right|$$

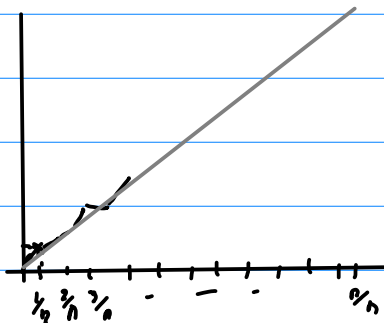
||

$$\Phi_{\bar{X}_n, \sqrt{S_n^2}}(x)$$

Ce test est très utilisé pour l'analyse des résidus d'un modèle linéaire gaussien

$(Y_j - (X \hat{\beta}_n)_j, j \leq n)$ suit un n-échantillon de loi normale si l'analyse de TLG est justifiée

On représente ce test sous la forme d'un "q-q plot"



3 - Test d'homogénéité de Kolmogorov - Smirnov

Ce test permet de vérifier que deux jeux de données suivent la même loi.

(X_1, \dots, X_n) de loi F et (Y_1, \dots, Y_m) de loi G

On souhaite tester $H_0 = \{F = G\}$ contre H_1

On définit la statistique $\|F_n - G_m\|_\infty$

Cette statistique dépend uniquement de $I_m(F)$ et $I_m(G)$ donc si X et Y sont à densité par rapport à la mesure de Lebesgue, c'est une statistique libre.