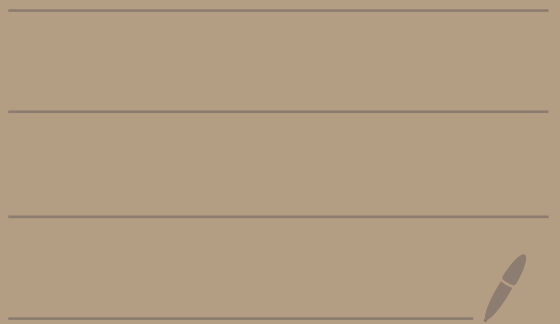


VE

Modélisation

Topic 3:

MLG

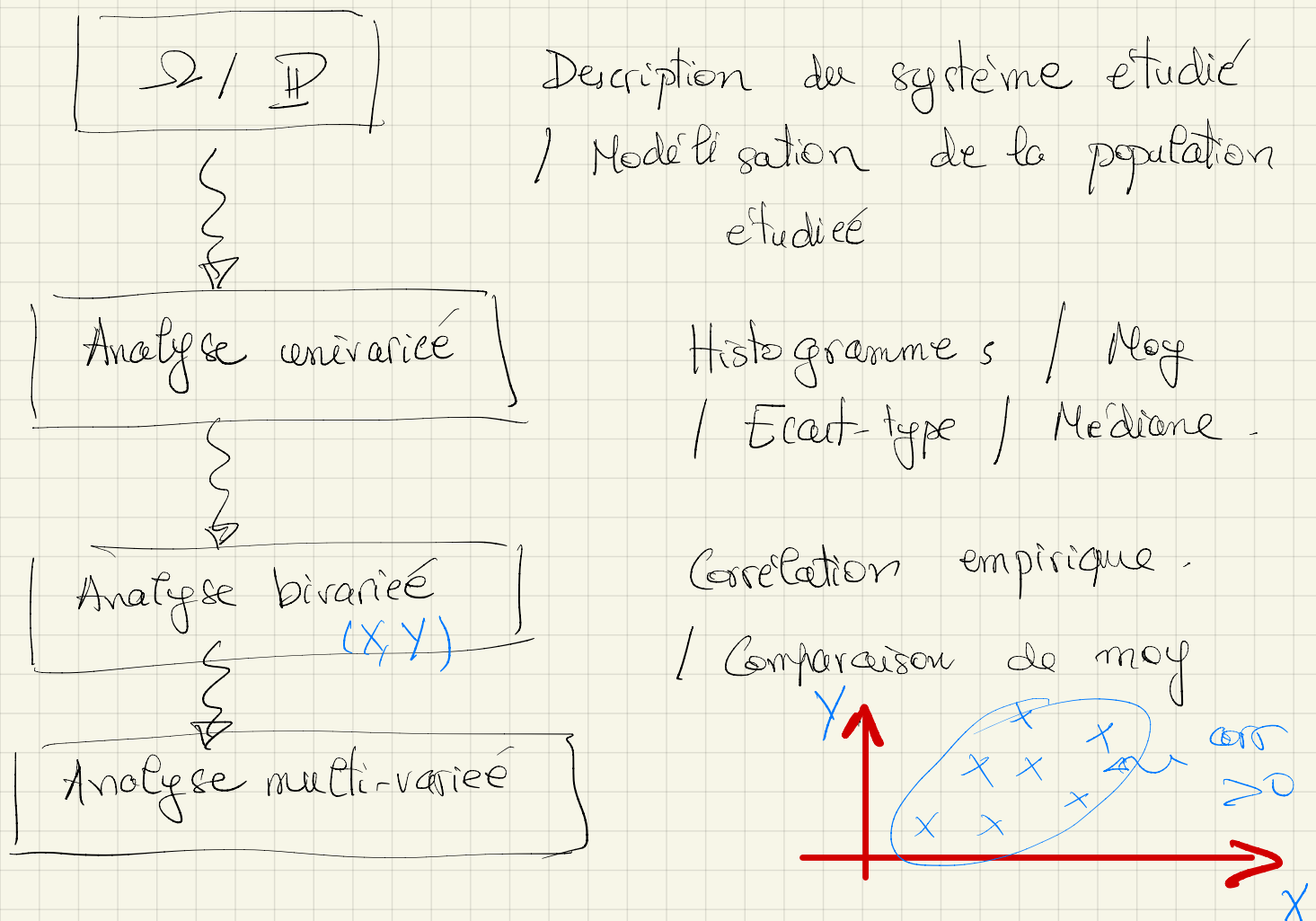


Maintenant commence la partie « statistique »
de l'UE. Plus précisément, le but
est de vous présenter le modèle linéaire
généralisé (MLG).

Uniquement pour fixer les notations et servir de
rappel:

0. Rappel : Le modèle linéaire simple

① Preamble : La démarche statistique.



a pour but de répondre à la question :

« Comment expliquer une variable Y en fonction de plusieurs variables explicatives, X^1, X^2, \dots, X^k ? ? ? \Rightarrow

Dans ce cours, on verra

\hookrightarrow le modèle linéaire

$$Y = \sum_{j=1}^k \theta_j X^j + \varepsilon$$

θ_j paramètres du modèle
 X^j variables explicatives
 ε Bruit / Erreur
 $E(\varepsilon) = 0$
 $\text{Var}(\varepsilon) = \sigma^2$

$f_j(x^j; \theta_j)$
 $= \theta_j x^j$

Écriture matricielle $Y = X \theta + \varepsilon$
avec $X^1 = 1$ (intercept).

Remarque: On peut allégrement modifier le modèle pour inclure des non-linéarités

$$Y = \sum_{j=1}^k f_j(X^j, \theta) + \varepsilon$$

\hookrightarrow fonctions à estimer (ex polynôme)

ou pire encore : (ex: $f_j(z, \theta) = \sum_{k=0}^{d_j} z^k \theta_{j,k}$)

$$Y = F_{\theta}(X^1, X^2, X^3, \dots, X^k) + \varepsilon$$

\hookrightarrow Fonction à estimer

par réseau de neurone
ou plus précisément par descente de gradient

Regression paramétrique

Sondage: Faire un TP de régression paramétrique?

Oui / Non.

Sondage: Faire un TP avec réseaux de neurones.

Oui / Non.

② Estimation MCO (Moindres carrés ordinaires) / OLS (ordinary Least Squares).

$$\begin{cases} Y_1 = X_1 \theta + \varepsilon_1 \\ Y_2 = X_2 \theta + \varepsilon_2 \\ \vdots \\ Y_m = X_m \theta + \varepsilon_m \end{cases} \quad \varepsilon_i \text{ iid} \quad \Leftrightarrow \quad \begin{matrix} Y \\ \downarrow n \end{matrix} = \begin{matrix} X \theta + \varepsilon \\ \downarrow n \end{matrix}$$

L'estimateur MCO / OLS consiste à prendre

$$\hat{\theta}_m^{\text{OLS}} = \underset{\theta}{\text{Argmin}} \underbrace{\| Y - X \theta \|^2}_{\text{Loss / Perte quadratique}}$$

Proposition: $\hat{\theta}_m^{\text{OLS}} = (X^T X)^{-1} X^T Y$



$$\hat{\theta}_m^{\text{OLS}} = \hat{\theta}_m^{\text{MLE}}$$

$$= \underset{\theta}{\text{Argmax}} \log L(\theta)$$

Estimateur de max de vraisemblance

lorsque $\varepsilon = (\varepsilon_1, \dots, \varepsilon_m)$ iid Gaussiens.

En effet $L(\theta) =$ Densité d'observées

$$(y_1, \dots, y_m) = Y$$

$$(x_1, \dots, x_m) = X.$$

= Densité d'observées :

$$\text{iid} \rightsquigarrow \varepsilon = Y - X\theta = \begin{pmatrix} y_1 - x_1\theta \\ \vdots \\ y_m - x_m\theta \end{pmatrix}$$

$$= \prod_{j=1}^m \mathbb{P}_{\varepsilon_j}(y_j - x_j\theta)$$

$$= \prod_{j=1}^m \exp\left(-\frac{\|y_j - x_j\theta\|^2}{2\sigma^2}\right) / \sqrt{2\pi\sigma^2}$$

$$= \exp\left(-\sum_{j=1}^m \frac{\|y_j - x_j\theta\|^2}{2\sigma^2}\right) / (2\pi\sigma^2)^{m/2}$$

$$\rightsquigarrow \log L(\theta) = -\sum_{j=1}^m \frac{\|y_j - x_j\theta\|^2}{2\sigma^2} - \frac{m}{2} \log(2\pi\sigma^2)$$

$$\text{Ainsi } \hat{\theta}_n^{\text{MLE}} = \underset{\theta}{\text{Argmax}} \log L(\theta)$$

$$= \underset{\theta}{\text{Argmax}} -\sum_{j=1}^m \frac{\|y_j - x_j\theta\|^2}{2\sigma^2} - \frac{m}{2} \log(2\pi\sigma^2)$$

$$= \underset{\theta}{\text{Argmax}} -\|Y - X\theta\|^2$$

$$= \underset{\theta}{\text{Argmin}} \|Y - X\theta\|^2$$

$$= \hat{\theta}_n^{\text{OLS}}$$

?

□

Propriétés de $\hat{\Theta}_m^{OLS}$

- $\hat{\Theta}_m^{OLS} = ({}^tXX)^{-1} {}^tXY$
 $= \Theta + ({}^tXX)^{-1} {}^tX \varepsilon$
- $\hat{\Theta}_m^{OLS}$ est consistant ($\hat{\Theta}_m^{OLS} \xrightarrow{P} \Theta$)
et sous hypothèse ε Gaussien
 $\hat{\Theta}_m \stackrel{\mathcal{L}}{=} \mathcal{CP}(\Theta, \sigma^2 ({}^tXX)^{-1})$
- $\hat{\Theta}_m^{OLS}$ sans biais.

③ Résidus et valeurs ajustées

Après estimation de $\hat{\Theta}_m$ on peut former :

↳ « les valeurs ajustées » / les prédictions

$$\hat{Y} = X \hat{\Theta}_m \stackrel{\mathcal{L}}{=} \mathcal{CP}(X\Theta, X({}^tXX)^{-1} {}^tX)$$

↳ si ε Gaussien

↳ « les résidus » / les erreurs de prédiction

$$\hat{\varepsilon} = Y - \hat{Y}$$

Propriétés : Si ε Gaussien

$$\hat{Y} \perp \hat{\varepsilon}$$
$$\hat{\varepsilon} \perp \hat{\Theta}_m$$

④ Estimation de σ^2

Le bon estimateur de σ^2 est

$$\hat{\sigma}_m^2 = \frac{1}{n-k} \|\hat{\epsilon}\|^2 = \frac{1}{n-k} \sum_{j=1}^m \hat{\epsilon}_j^2$$

↑ nombre de facteurs

(l'intercept $X^1 = 1$ inclus)

Sous hypothèse de ϵ Gaussien :

Propriétés :

① $\hat{\epsilon} \perp \hat{\sigma}_m^2$ et $\hat{\sigma}_m^2 \perp \hat{\sigma}_m$

② $\frac{(n-k)\hat{\sigma}_m^2}{\sigma^2} \stackrel{\mathcal{L}}{\sim} \chi_{n-k}^2$ (Somme de $n-k$ carrés de normales)

③ $\hat{\sigma}_m^2$ sans biais et de variance $\frac{2\sigma^4}{n-k} \rightarrow 0$ $n \rightarrow +\infty$

Exercice : Montrez que ② \Rightarrow ③

et ③ \Rightarrow $\hat{\sigma}_m^2$ est un estimateur consistant.

$$\left(\hat{\sigma}_m^2 \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \sigma^2 \right)$$

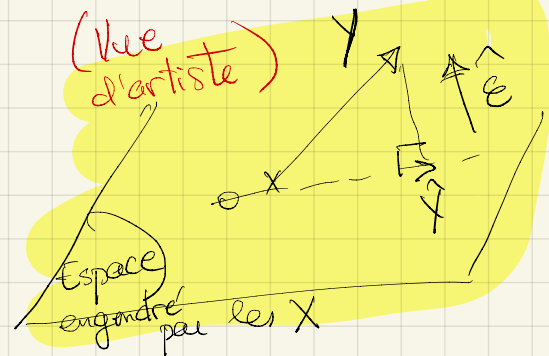
⑤ Analyse et exploitation de la régression.

↳ ANOVA (Analyse de la variance)

Il est question de savoir si la régression est pertinente. Pour cela on observe la décomposition

$$Y = \hat{Y} + \hat{\epsilon}$$

En fait c'est une décomposition



\leadsto en variables indépendantes (Propriétés en ③)

donc
$$\text{Var } Y = \text{Var } \hat{Y} + \text{Var } \hat{\varepsilon}$$

\leadsto en vecteurs orthogonaux (indépendance géométrique)
après centrage

donc
$$\|Y - \bar{Y}\|^2 = \|\hat{Y} - \bar{\hat{Y}}\|^2 + \|\hat{\varepsilon}\|^2$$

La qualité de la regression est mesurée

par le critère
$$R^2 = \frac{\|\hat{Y} - \bar{\hat{Y}}\|^2}{\|Y - \bar{Y}\|^2} \in [0, 1]$$

$\left\{ \begin{array}{l} \text{Si } R^2 \approx 1, \text{ alors } X \text{ explique bien } Y. \\ \text{Si } R^2 \approx 0, \text{ l'inverse} \end{array} \right.$

⑥ Tests de nullité :

$(H_0: \theta_j = 0)$

on a: $\hat{\theta}_j \stackrel{\mathcal{L}}{=} \mathcal{N}(\theta_j, \sigma^2 (tX X)^{-1})$

mais σ^2 inconnue !

Ainsi
$$\left\{ \frac{\hat{\theta}_j - \theta_j}{\sqrt{\sigma^2 (tX X)^{-1}_{jj}}} \stackrel{\mathcal{L}}{=} \mathcal{N}(0, 1) \right.$$

\Downarrow sous hypothèse de ε Gaussien

Par ratio, $T = \frac{\hat{\theta}_j - \theta_j}{\sqrt{\hat{\sigma}_n^2 (X'X)^{-1}_{jj}}} = \frac{dP(0,1)}{\sqrt{\frac{\chi_{n-p}^2}{n-p}}}$

$\stackrel{||}{=} \text{Student } (n-p)$

\uparrow définition de cette loi

Si $t_{1-\frac{\alpha}{2}}$ est la $(1-\frac{\alpha}{2})^e$ quantile de cette loi

Alors $IC_{1-\alpha}(\theta_j) = \hat{\theta}_j + \sqrt{\hat{\sigma}_n^2 (X'X)^{-1}_{jj}} \times \left[-t_{1-\frac{\alpha}{2}} ; t_{1-\frac{\alpha}{2}} \right]$

est l'intervalle de confiance avec certitude $1-\alpha$.

\leadsto On rejette H_0 si $|\hat{\theta}_j| \geq t_{1-\frac{\alpha}{2}} \sqrt{\hat{\sigma}_n^2 (X'X)^{-1}_{jj}}$

Le modèle linéaire généralisé

I. Introduction

Il existe plusieurs variantes de modèle linéaire (polynômes, paramétrique, non-param) :

$$Y = X\beta + \varepsilon$$

Variable d'intérêt $\leftarrow \overline{Y}$ Facteurs $\leftarrow \overline{X}$ $\overline{\beta}$ à estimer $\overline{\varepsilon}$ Bruit

Toutefois, il y a des cas très simples où ceci n'a aucune chance de fonctionner.

Par exemple, si la variable d'intérêt est
 \rightarrow binaire ($\in \{0, 1\}$) et que l'on s'intéresse à la probabilité

$$p = \mathbb{P}(Y = 1)$$

\rightarrow le comptage ($\in \mathbb{N}$).

Ainsi nous verrons dans ce chapitre une généralisation dans deux directions

\rightarrow Utilisation de lois plus générales

dans le cadre des lois de la famille exponentielle.
En effet, le modèle linéaire et son étude sont très « gaussiens ».

↳ Incorporation d'une transformation pour passer de \mathbb{R} à $[0, 1]$; \mathbb{N} ; $\{0, 1\}$

Le modèle linéaire généralisé se décompose en 3 composantes.

① La composante aléatoire Y (\equiv Variable d'intérêt) à laquelle est associée une loi de proba. / de réponse

② La composante déterministe composée d'une combinaison linéaire de X^1, X^2, \dots, X^p (variables explicatives) sous $X\beta$ où $X = (X^1, \dots, X^p)$

③ **Un lien** c'est-à-dire une transformation entre $\left\{ \begin{array}{l} \text{la composante déterministe} \\ \text{un paramètre de la composante aléatoire.} \end{array} \right.$

Ex: $Y \stackrel{\text{lien}}{=} \mathcal{P}(p, 1)$
 $X\beta$ lien

II. Les 3 composantes

II.1. La composante déterministe

C'est juste la combinaison lin. $X\beta = \sum \beta_i X^i$

variables
fixées
par
l'exp.

Autre nom : Prédicteur linéaire

I - 2 - La composante aléatoire.

On dispose de n observations (y_1, y_2, \dots, y_n) de la v.a. Y . On choisit une loi dans la famille exponentielle en fonction de la nature de Y .

$\hookrightarrow Y \in \mathbb{R}$, valeur quantitative.

Ex naturel $\hookrightarrow Y \stackrel{\mathcal{L}}{=} \mathcal{N}(\mu, \sigma^2)$,

et $\mu = \mathbb{E}(Y)$

$\hookrightarrow Y \in \{0, 1\}$ (succès / échec).

Ex naturel \hookrightarrow

$Y \stackrel{\mathcal{L}}{=} \text{Bernoulli}(\pi)$ où $\pi = \mathbb{E}(Y) \in [0, 1]$

$\hookrightarrow Y \in \mathbb{N}$ (Comptage)

Exemple naturel \hookrightarrow

$Y \stackrel{\mathcal{L}}{=} \text{Poisson}(\lambda)$

où $\lambda = \mathbb{E}(Y) \in \mathbb{R}_+$

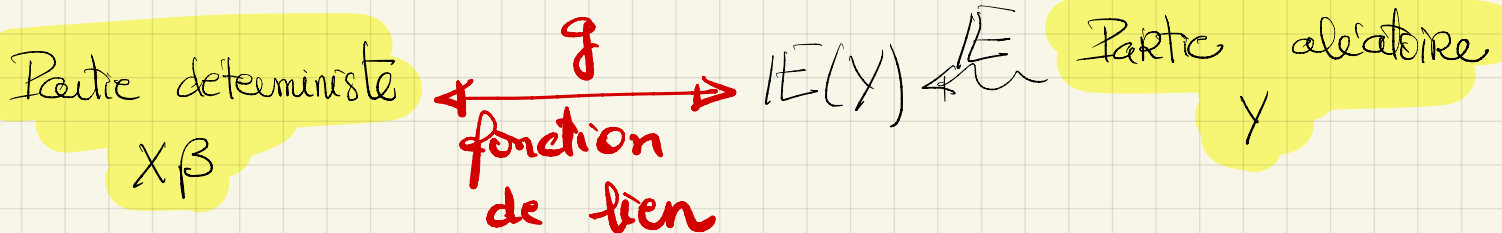
III - 3 - La fonction de lien

Définition: On définit $g: I \text{ intervalle} \rightarrow \mathbb{R}$

donnant $g(\mathbb{E}(Y)) = X\beta$ où $\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_R \end{pmatrix}$

$X = \begin{pmatrix} X^1 & X^2 & \dots & X^R \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix}$

Dit simplement



Exemples: \hookrightarrow Si $Y \stackrel{\mathcal{L}}{=} \mathcal{N}(\mu, \sigma^2) = \mu + \sigma \mathcal{N}(0, 1)$

Alors $g(\mu) = \mu$

$$\Rightarrow Y = X\beta + \underbrace{\sigma \mathcal{N}(0, 1)}_{\varepsilon}$$

C'est le modèle linéaire.

\hookrightarrow Si $Y \stackrel{\mathcal{L}}{=} \mathcal{P}(\lambda)$

Alors $X\beta = g(\mathbb{E}(Y)) = g(\lambda) \stackrel{\text{choix}}{=} \log \lambda$

donne $Y \stackrel{\mathcal{L}}{=} \mathcal{P}(\lambda = e^{X\beta})$

(Modèle log-linéaire).

\hookrightarrow Si $Y \stackrel{\mathcal{L}}{=} \text{Bernoulli}(\pi)$

Alors $X\beta = g(\mathbb{E}(Y)) = g(\pi) \stackrel{\text{choix}}{=} \log\left(\frac{\pi}{1-\pi}\right)$

donne

$Y \stackrel{\mathcal{L}}{=} \text{Bernoulli}\left(\text{logit}^{\leftarrow}(\pi)\right) \stackrel{\text{choix}}{=} \text{logit}(\pi)$

(Modèle de régression logistique)



$$\text{logit}^{\leftarrow}(y) = \frac{1}{1+e^{-y}} = \frac{e^y}{1+e^y} \quad (\text{Fonction sigmoïde})$$

$$y = \text{logit}(x) \iff e^y = \frac{x}{1-x} \iff e^{-y} = \frac{1-x}{x}$$

$$\iff 1+e^{-y} = \frac{1}{x}$$

$$\iff x = \frac{1}{1+e^{-y}}$$

III - Familles exponentielles.

Large famille qui contient

↳ Bernoulli, Binomiale, Poisson (1 param)

↳ Normale, Gamma, Binomiale négative (2 param)

↳ Multinômiale (a + de 2 param).

La densité utilise deux paramètres $\psi = \begin{pmatrix} \theta \\ \phi \end{pmatrix}$

où $\begin{cases} \theta : \text{Paramètre principal lié à } E(Y) \\ \phi : \text{Paramètre de nuisance lié à } \text{Var}(Y) \end{cases}$
(Prendre 0 si loi à 1 param)

Y appartient à la famille exponentielle avec param (θ, ϕ) lorsque

$$P(Y \in dy) = dy f(y; \theta, \phi) \quad (\text{Continue})$$

$$P(Y = y) = f(y; \theta, \phi) \quad (\text{discret})$$

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

Les expressions exactes de $c(y, \phi)$ & $a(\phi)$ ne sont pas importantes et sont identifiées au cas par cas. On a tout de même les propriétés suivantes, que l'on peut établir

en toute généralité :

Thm:

$$\begin{cases} b'(\theta) = \mathbb{E}(Y) \\ \text{Var}(Y) = a(\phi) b''(\theta) \end{cases} \iff \theta = (b')^{\leftarrow 1}(\mathbb{E}(Y))$$

Si un seul paramètre au modèle, on prend $a(\phi) = 1$

Lien avec MLG : Dans toute famille exponentielle

$$\begin{aligned} X\beta &= g(\mathbb{E}(Y)) \quad \text{par def de MLG} \\ &= g(b'(\theta)) \end{aligned} \quad \left. \begin{array}{l} \text{utilisé pour MLG} \\ Y \stackrel{\mathcal{L}}{=} \text{Famille exp} \\ (\theta, \phi) \end{array} \right\}$$

Posez $\theta = X\beta$ comme paramètre de la loi

donne $\theta = g \circ b'(\theta) \iff g = (b')^{\leftarrow 1}$

Preuve : Montrons que $\mathbb{E}(Y) = b'(\theta)$
en toute généralité.

Preuve amusante de statisticien :

Je vais calculer l'estimateur de max de vraisemblance

$\hat{\theta}_n$ pour un tirage iid. Par consistance $\hat{\theta}_n \xrightarrow{\mathbb{P}} \theta$

Calcul : $L(y; \theta, \phi) = \prod_{i=1}^n \exp\left(\frac{\theta y_i - b(\theta)}{a(\phi)} + c(y_i, \phi)\right)$

où $y = (y_1, \dots, y_n)$ échantillon observé.

Donc

$$\hat{\Theta}_m = \underset{\Theta}{\text{Argmax}} \log L(y, \Theta, \phi)$$

$$= \underset{\Theta}{\text{Argmax}} \sum_{i=1}^m \left(\frac{\Theta y_i - b(\Theta)}{a(\phi)} + c(y_i, \phi) \right)$$

$$= \underset{\Theta}{\text{Argmax}} \Theta \left(\sum_{i=1}^m y_i \right) - m b(\Theta)$$

$$= \underset{\Theta}{\text{Argmax}} \Theta \left(\frac{\sum_{i=1}^m y_i}{m} \right) - b(\Theta)$$

d'où $\frac{\partial}{\partial \Theta} = 0 \iff \frac{\sum_{i=1}^m y_i}{m} = b'(\Theta)$

$$\iff \Theta = (b')^{\leftarrow} \left(\frac{\sum_{i=1}^m y_i}{m} \right)$$

Donc

$$\hat{\Theta}_m = (b')^{\leftarrow} \left(\frac{\sum_{i=1}^m y_i}{m} \right) \xrightarrow{m \rightarrow \infty} \Theta$$

$$(b')^{\leftarrow} (E(Y))$$

□

⚠ Ce résultat est difficile à montrer en calculant

$E(Y)$ à partir de la densité!