

# STT 1000 - STATISTIQUES

ARTHUR CHARPENTIER



## Likelihood / Vraisemblance

Assume that  $\{x_1, x_2, \dots, x_n\}$  are obtained from i.i.d. random variables  $X_1, X_2, \dots, X_n$ , with identical distribution  $F_\theta$ , and density  $f_\theta$ .

$$\mathcal{L}(\theta) = f_\theta(\mathbf{x}) = f_\theta(x_1, \dots, x_n) = \prod_{i=1}^n f_\theta(x_i)$$

$$\log \mathcal{L}(\theta) = \sum_{i=1}^n \log f_\theta(x_i)$$

$\hat{\theta}$  is a **maximum likelihood estimator** of parameter  $\theta$  if

$$\hat{\theta} \in \operatorname{argmax}\{\mathcal{L}(\theta)\} = \operatorname{argmax}\{\log \mathcal{L}(\theta)\}$$

## Likelihood / Vraisemblance

Given some sample  $\{x_1, \dots, x_n\}$   
from a  $\mathcal{N}(\mu, \sigma^2)$  distribution,

FIGS/MLE-lik.png

$$\mathcal{L}(\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2}\right)$$

$$\log \mathcal{L}(\mu, \sigma^2) = -\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2} - \frac{n}{2} \log(2\pi\sigma^2)$$

Here  $\theta = (\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)$ .

FIGS/MLE-log-lik.png

## Likelihood / Vraisemblance

The first order condition (also called likelihood equations) is

$$\left. \frac{\partial \log (\mathcal{L}(\theta; x_1, \dots, x_n))}{\partial \theta} \right|_{\theta=\hat{\theta}} = 0$$

Second order condition is

$$\left. \frac{\partial^2 \log (\mathcal{L}(\theta; x_1, \dots, x_n))}{\partial \theta} \right|_{\theta=\hat{\theta}} < 0$$

**Example:** if  $X \sim \mathcal{P}(\lambda)$ ,

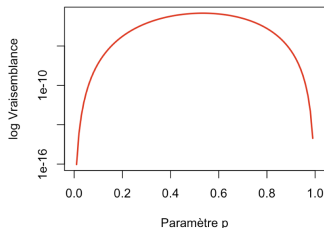
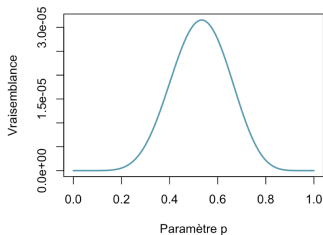
$$\log \mathcal{L}(\lambda; x_1, \dots, x_n) = \sum_{i=1}^n \log \left( e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} \right) = -n\lambda + n\bar{x} \log(\lambda) - \log \left( \prod_{i=1}^n x_i! \right)$$

$$\frac{\partial \log (\mathcal{L}(\lambda; x_1, \dots, x_n))}{\partial \lambda} = -n + \frac{n\bar{x}}{\lambda}, \text{ so } \hat{\lambda} = \bar{x}.$$

## Vraisemblance, cas $\mathcal{B}(p)$

- ce que nous dit la pratique

```
1 > n=15
2 > set.seed(1)
3 > x=sample(0:1,size = n,prob = c(.4,.6),replace=TRUE)
4 > vraisemblance = function(p) prod(dbinom(x,size = 1,
      prob = p))
5 > vect_p = seq(0,1,by=0.01)
6 > plot(vect_p,Vectorize(vraisemblance)(vect_p))
```



## Likelihood / Vraisemblance

### Score

La fonction  $S : x \mapsto \frac{d}{d\theta} \log f_{\theta}(x)$  est appelée **fonction score**.

If  $X \sim f_{\theta}$ ,  $\mathbb{E} \left( \frac{d}{d\theta} \log f_{\theta}(X) \right) = \mathbb{E}(S(X)) = 0$

**Example:** if  $X \sim \mathcal{P}(\lambda)$ ,

$$\frac{d \log f_{\lambda}(X)}{d\lambda} = -1 + \frac{X}{\lambda}, \text{ so } \mathbb{E} \left( \frac{d}{d\lambda} \log f_{\lambda}(X) \right) = -1 + \frac{\mathbb{E}(X)}{\lambda} = 0$$

**Example:** if  $X \sim \mathcal{B}(p)$ ,

$$\frac{d \log f_{\lambda}(X)}{d\lambda} = \frac{X}{p} - \frac{1-X}{1-p}, \mathbb{E} \left( \frac{d}{d\lambda} \log f_{\lambda}(X) \right) = \frac{\mathbb{E}(X)}{p} - \frac{1 - \mathbb{E}(X)}{1-p} = 0$$

## Likelihood / Vraisemblance

Si  $x_1, \dots, x_n$  est tiré suivant  $f_\theta$ ,  $\mathbb{E} \left( \frac{d}{d\theta} \log \mathcal{L}(\mathbf{X}) \right) = 0$

**Example:** if  $X_i \sim \mathcal{P}(\lambda)$ ,

$$\frac{d \log \mathcal{L}}{d\lambda} = -n + \frac{\sum_{i=1}^n x_i}{\lambda}, \mathbb{E} \left( \frac{d}{d\lambda} \log \mathcal{L} \right) = -n + \frac{\sum_{i=1}^n \mathbb{E}(X_i)}{\lambda} = 0$$

**Example:** if  $X_i \sim \mathcal{B}(p)$ ,

$$\frac{d \log \mathcal{L}}{d\lambda} = \frac{\sum_{i=1}^n X_i}{p} - \frac{n - \sum_{i=1}^n X_i}{1 - p},$$

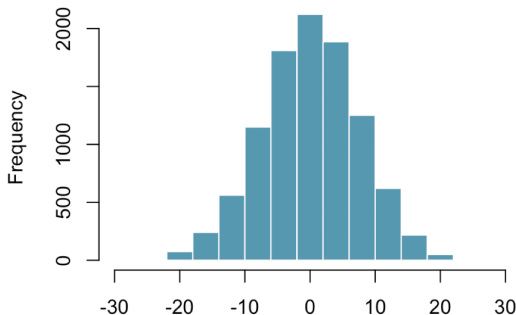
$$\mathbb{E} \left( \frac{d}{d\lambda} \log \mathcal{L} \right) = \frac{\sum_{i=1}^n \mathbb{E}(X_i)}{p} - \frac{n - \sum_{i=1}^n \mathbb{E}(X_i)}{1 - p} = 0$$

## Likelihood / Vraisemblance

```
1 > pente=rep(NA,1e4)
2 > for(s in 1:1e4){
3 +   x=sample(0:1, n,prob = c(.4,.6),replace=TRUE)
4 +   dlogL = function(p) sum(x)/p-(sum(1-x))/(1-p)
5 +   pente[s] = dlogL(0.6) }
6 > hist(pente)
```

La distribution empirique de  $\frac{d}{d\theta} \log \mathcal{L}(\mathbf{X})$  est

```
1 > mean(pente)
2 [1] -0.03458333
3 > var(pente)
4 [1] 63.51731
5 > 1/var(pente)
6 [1] 0.01574374
7 > .4*.6/15
8 [1] 0.016
```





# Information de Fisher

## Information de Fisher

Fisher information associated with a density  $f_\theta$ , with  $\theta \in \mathbb{R}$  is

$$I(\theta) = \mathbb{E} \left( \frac{d}{d\theta} \log f_\theta(X) \right)^2 \text{ where } X \text{ has distribution } f_\theta,$$

$$I(\theta) = \text{Var} \left( \frac{d}{d\theta} \log f_\theta(X) \right) = -\mathbb{E} \left( \frac{d^2}{d\theta^2} \log f_\theta(X) \right).$$

For a sample of size  $n$ ,

$$I_n(\theta) = \mathbb{E} \left( \frac{\partial}{\partial \theta} \log \mathcal{L}(\theta, X_1, \dots, X_n) \right)^2 = nI(\theta)$$

## Information de Fisher

en effet, posons  $U = \frac{\partial}{\partial \theta} \ln f(X; \theta) = \frac{\frac{\partial}{\partial \theta} f(X; \theta)}{f(X; \theta)}$

$$E_{\theta}(U) = \int_{\mathbb{R}} \frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} f(x; \theta) dx = \int_{\mathbb{R}} \frac{\partial}{\partial \theta} f(x; \theta) dx = \frac{\partial}{\partial \theta} \underbrace{\int_{\mathbb{R}} f(x; \theta) dx}_{=1} = 0$$

$$\begin{aligned} \frac{\partial^2}{\partial \theta^2} \ln f(x; \theta) &= \frac{\frac{\partial^2}{\partial \theta^2} f(x; \theta) \cdot f(x; \theta) - \left[ \frac{\partial}{\partial \theta} f(x; \theta) \right]^2}{[f(x; \theta)]^2} \\ &= \frac{\frac{\partial^2}{\partial \theta^2} f(x; \theta)}{f(x; \theta)} - \left[ \frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} \right]^2 \end{aligned}$$

## Information de Fisher

donc

$$E_{\theta} \left[ \frac{\partial^2}{\partial \theta^2} \ln f(X; \theta) \right] = E_{\theta} \left[ \frac{\frac{\partial^2}{\partial \theta^2} f(X; \theta)}{f(X; \theta)} \right] - E_{\theta} [U^2]$$

or

$$E_{\theta} \left[ \frac{\frac{\partial^2}{\partial \theta^2} f(X; \theta)}{f(X; \theta)} \right] = \int_{\mathbb{R}} \frac{\partial^2}{\partial \theta^2} f(x; \theta) dx = \frac{\partial^2}{\partial \theta^2} \int_{\mathbb{R}} f(x; \theta) dx = 0$$

donc

$$\text{Var} \left( \frac{\partial}{\partial \theta} \log f_{\theta}(X) \right) = -\mathbb{E} \left( \frac{\partial^2}{\partial \theta^2} \log f_{\theta}(X) \right).$$

## Likelihood / Vraisemblance

**Example:** if  $X$  has a Poisson distribution  $\mathcal{P}(\theta)$ ,

$$\log f_{\theta}(x) = -\theta + x \log \theta - \log(x!) \text{ and } \frac{d^2}{d\theta^2} \log f_{\theta}(x) = -\frac{x}{\theta^2}$$

$$I(\theta) = -\mathbb{E} \left( \frac{d^2}{d\theta^2} \log f_{\theta}(X) \right) = -\mathbb{E} \left( -\frac{X}{\theta^2} \right) = \frac{1}{\theta}$$

**Example:** if  $X$  has a binomial distribution  $\mathcal{B}(1, \theta)$ ,

$$I(\theta) = -\mathbb{E} \left( \frac{d^2}{d\theta^2} \log f_{\theta}(X) \right) = -\mathbb{E} \left( -\frac{X}{\theta^2} + \frac{1-X}{(1-\theta)^2} \right) = \frac{1}{\theta(1-\theta)}$$

### Borne de Cramér-Rao

If  $\hat{\theta}$  is an **unbiased estimator** of  $\theta$ , then

$$\text{Var}(\hat{\theta}) \geq \frac{1}{nI(\theta)}$$

If that bound is attained, the estimator is said to be **efficient**.  
An unbiased estimator  $\hat{\theta}$  is said to be **optimal** if it has the lowest variance among all unbiased estimators, see **bias**, **minimum variance unbiased estimator**

## Likelihood / Vraisemblance

$$\text{if } \boldsymbol{\theta} \in \mathbb{R}^k, \left. \frac{\partial \log(\mathcal{L}(\boldsymbol{\theta}; x_1, \dots, x_n))}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} = \mathbf{0}$$

Second order condition is

$$\text{if } \boldsymbol{\theta} \in \mathbb{R}^k, \left. \frac{\partial^2 \log(\mathcal{L}(\boldsymbol{\theta}; x_1, \dots, x_n))}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \text{ is definite negative}$$

If  $\boldsymbol{\theta} \in \mathbb{R}^k$ , then Fisher information is the  $k \times k$  matrix  $I = [I_{i,j}]$  with

$$I_{i,j} = \mathbb{E} \left( \frac{\partial}{\partial \theta_i} \log f_{\boldsymbol{\theta}}(X) \frac{\partial}{\partial \theta_j} \log f_{\boldsymbol{\theta}}(X) \right).$$

i.e.

$$I(\boldsymbol{\theta}) = \mathbb{E} \left[ \left( \frac{d}{d\boldsymbol{\theta}} \log f_{\boldsymbol{\theta}}(X) \right) \left( \frac{d}{d\boldsymbol{\theta}} \log f_{\boldsymbol{\theta}}(X) \right)^{\top} \right]$$

$$I(\boldsymbol{\theta}) = -\mathbb{E} \left( \frac{d^2}{d\boldsymbol{\theta} d\boldsymbol{\theta}^{\top}} \log f_{\boldsymbol{\theta}}(X) \right)$$

## Likelihood / Vraisemblance

For a Gaussian distribution  $\mathcal{N}(\theta, \sigma^2)$ ,  $I(\theta) = \frac{1}{\sigma^2}$

For a Gaussian distribution  $\mathcal{N}(\mu, \theta)$ ,  $I(\theta) = \frac{1}{2\theta^2}$

For a Gaussian distribution  $\mathcal{N}(\boldsymbol{\theta})$ ,  $I(\boldsymbol{\theta}) = \begin{pmatrix} 1/\sigma^2 & 0 \\ 0 & 2/\sigma^2 \end{pmatrix}$

Cramér-Rao bound is  $\frac{1}{n}I^{-1} = \frac{1}{n} \begin{pmatrix} \sigma^2 & 0 \\ 0 & \sigma^2/2 \end{pmatrix}$

## Likelihood / Vraisemblance

En effet, la log-densité de la loi normale  $\mathcal{N}(\mu, \nu)$  est

$$\log f(x; \mu, \nu) = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \nu - \frac{1}{2\nu}(x - \mu)^2$$

Les dérivées premières sont

$$\frac{\partial}{\partial \mu} \ln f(x; \mu, \nu) = \frac{1}{\nu}(x - \mu) \text{ et } \frac{\partial}{\partial \nu} \ln f(x; \mu, \nu) = -\frac{1}{2\nu} + \frac{(x - \mu)^2}{2\nu^2}$$

$$E \left[ \frac{1}{\nu^2}(X - \mu)^2 \right] = \frac{1}{\nu^2} \nu = \frac{1}{\nu} = \frac{1}{\sigma^2}$$

$$\begin{aligned} E \left[ \left( -\frac{1}{2\nu} + \frac{(X - \mu)^2}{2\nu^2} \right)^2 \right] &= E \left[ \frac{1}{4\nu^2} - \frac{(X - \mu)^2}{2\nu^3} + \frac{(X - \mu)^4}{4\nu^4} \right] \\ &= \frac{1}{4\nu^2} - \frac{\nu}{2\nu^3} + \frac{3\nu^2}{4\nu^4} = \frac{1}{2\nu^2} = \frac{1}{2\sigma^4} \end{aligned}$$

$$\text{car } E[(X - \mu)^4] = 3\sigma^4.$$



## Likelihood / Vraisemblance

Enfin

$$E \left[ \frac{1}{v} (X - \mu) \left( -\frac{1}{2v} + \frac{(X - \mu)^2}{2v^2} \right) \right] = 0$$

d'où la matrice d'information

$$I(\mu, \sigma^2) = \begin{pmatrix} 1/\sigma^2 & 0 \\ 0 & 1/2\sigma^4 \end{pmatrix}$$

## Likelihood / Vraisemblance

Let  $\{x_1, \dots, x_n\}$  be a sample with distribution  $f_\theta$ , where  $\theta \in \Theta$ .  
The maximum likelihood estimator  $\hat{\theta}_n$  of  $\theta$  is

$$\hat{\theta}_n \in \operatorname{argmax}\{\mathcal{L}(\theta; x_1, \dots, x_n), \theta \in \Theta\}.$$

### Propriétés asymptotiques de l'EMV

Under some technical assumptions  $\hat{\theta}_n$  converges almost surely towards  $\theta$ ,  $\hat{\theta}_n \xrightarrow{a.s.} \theta$ , as  $n \rightarrow \infty$ .

Under some technical assumptions  $\hat{\theta}_n$  is asymptotically efficient,

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{\mathcal{L}} \mathcal{N}(0, I^{-1}(\theta)).$$

See [maximum likelihood estimation](#)

# Optimization

Consider some Poisson model,

```
1 > set.seed(1)
2 > (y=rpois(10,3))
3 [1] 2 2 3 5 2 5 6 4 3 1
4 > mean(y)
5 [1] 3.3
6 > NLogL = function(lambda) -sum(log(dpois(y,lambda)))
7 > optim(fn = NLogL,par = 1)
8 $par
9 [1] 3.3
10
11 $value
12 [1] 18.59581
```

## Calcul numérique du maximum de vraisemblance

Consider a sample  $\mathbf{X} = (X_1, \dots, X_n)$  i.id. from  $F_\theta$ . Let

$$S_{n,\theta}(\mathbf{x}) = \frac{\partial \log \mathcal{L}(\boldsymbol{\theta}; \mathbf{x})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^n S_{1,\theta}(x_i)$$

denote the **score function**. Then  $S_{n,\theta}(\mathbf{X})$  is a random vector. Then

$$\mathbb{E}[S_{n,\theta}(\mathbf{X})] = \mathbf{0}$$

while

$$\text{Var}[S_{n,\theta}(\mathbf{X})] = I_n(\boldsymbol{\theta}) = \mathbb{E} \left( \frac{\partial}{\partial \boldsymbol{\theta}} S_{n,\theta}(\mathbf{X}) \right).$$

$$\frac{S_{n,\theta}(\cdot)}{n} \xrightarrow{a.s.} 0 \quad \text{and} \quad \frac{S_{n,\theta}(\cdot)}{\sqrt{n}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, I(\boldsymbol{\theta})).$$

## Calcul numérique du maximum de vraisemblance

If  $\theta$  is univariate, use Taylor approximation of  $S_n$  in the neighbourhood of  $\theta_0$  (the true value)

$$S_n(x) = S_n(\theta_0) + (x - \theta_0)S'_n(y) \text{ for some } y \in [x, \theta_0]$$

Set  $x = \hat{\theta}_n$ , then

$$S_n(\hat{\theta}_n) = 0 = S_n(\theta_0) + (\hat{\theta}_n - \theta_0)S'_n(y) \text{ for some } y \in [\theta_0, \hat{\theta}_n]$$

$$\text{Hence, } \hat{\theta}_n = \theta_0 - \frac{S_n(\theta_0)}{S'_n(y)} \text{ for } y \in [\theta_0, \hat{\theta}_n].$$

### Algorithme de Newton Raphson

$$\hat{\theta}_n^{(i+1)} = \hat{\theta}_n^{(i)} - \frac{S_n(\hat{\theta}_n^{(i)})}{S'_n(\hat{\theta}_n^{(i)})},$$

from some starting value  $\hat{\theta}_n^{(0)}$  (hopefully well chosen).

# Calcul numérique du maximum de vraisemblance

Newton-Raphson:

$$\hat{\theta}_n^{(i+1)} = \hat{\theta}_n^{(i)} - \frac{S_n(\hat{\theta}_n^{(i)})}{S'_n(\hat{\theta}_n^{(i)})},$$

where

$$S'_n(\hat{\theta}_n^{(i)}) \sim \frac{S_n(\hat{\theta}_n^{(i)} + h) - S_n(\hat{\theta}_n^{(i)} - h)}{2h}$$

from some starting value  $\hat{\theta}_n^{(0)}$  (hopefully well chosen), and some small  $h > 0$ .

Score de Fisher

$$\hat{\theta}_n^{(i+1)} = \hat{\theta}_n^{(i)} + \frac{S_n(\hat{\theta}_n^{(i)})}{nI(\hat{\theta}_n^{(i)})},$$

from some starting value  $\hat{\theta}_n^{(0)}$  (hopefully well chosen).

## Calcul numérique du maximum de vraisemblance

Consider some Poisson model,  $S_1(\theta) = -1 + \frac{x}{\theta}$

```
1 > Sn = function(lambda) sum(-1+y/lambda)
2 > h = 1e-7
3 > dSn = function(lambda) (Sn(lambda+h)-Sn(lambda-h))
   /(2*h)
4 > L = rep(NA,10)
5 > L[1] = 1
6 > for(i in 1:9){
7 +   L[i+1] = L[i] - Sn(L[i])/dSn(L[i])
8 + }
9 > L
10 [1] 1.000 1.697 2.521 3.116 3.290 3.300 3.300 3.300
```

# Calcul numérique du maximum de vraisemblance

Consider some Poisson model, with Fisher information  $I(\theta) = \frac{1}{\theta}$

```
1 > I = function(lambda) 1/lambda
2 > L = rep(NA,10)
3 > L[1] = 1
4 > for(i in 1:9){
5 +   L[i+1] = L[i] - Sn(L[i])/(length(y)*I(L[i]))
6 + }
7 > L
8 [1] 1.0 3.3 3.3 3.3 3.3 3.3 3.3 3.3 3.3 3.3
```