



STT 1000 - STATISTIQUES

ARTHUR CHARPENTIER



Cumulative Distribution Function

Given a random variable X , $F(x) = \mathbb{P}[X \leq x]$

F is an increasing function, taking values in $[0, 1]$.

Fonction de répartition empirique \hat{F}

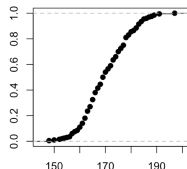
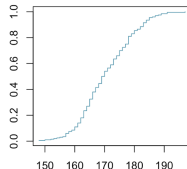
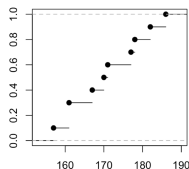
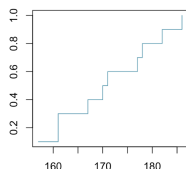
Consider a sample $\mathbf{x} = \{x_1, y_2, \dots, x_n\}$, a natural estimator is the empirical cumulative distribution function \hat{F}

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(x_i \leq x)$$

On peut prouver que $\mathbb{E}[\hat{F}(x)] = F(x)$.

Cumulative Distribution Function

```
1 > x = sort(x)
2 > n = length(x)
3 > y = (1:n)/n
4 > plot(x,y,type="s")
5 > plot(ecdf(x))
```

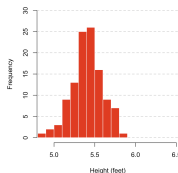
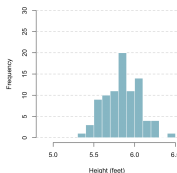
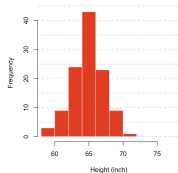
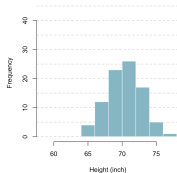
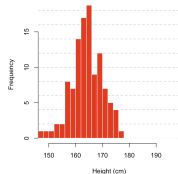
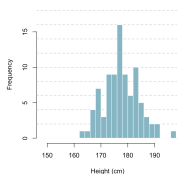
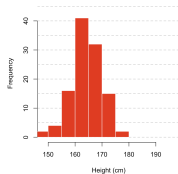
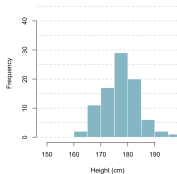


La fonction $x \mapsto \hat{F}(x)$ est une fonction en escalier, qui fait un saut de $1/n$ dès qu'elle croise une observation x_i .

Density & Histogram

Given a random variable X , f is such that $F(x) = \int_{-\infty}^x f(t)dt$
or conversely, $f(x) = F'(x)$.

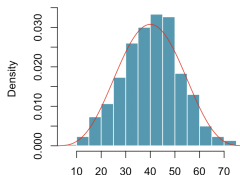
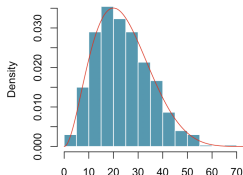
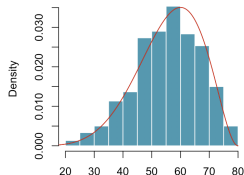
Thus, $\mathbb{P}(X \in [a, b]) = \int_a^b f(t)dt$



Density & Histogram

On dit qu'une distribution est unimodale si elle ne possède qu'un pic majeur.

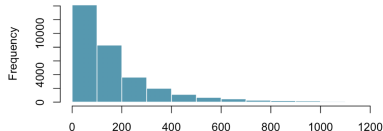
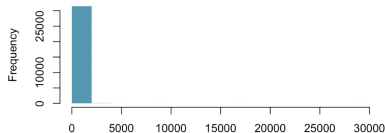
Quand une distribution n'est pas symétrique, elle est dite asymétrique ; on dit qu'une distribution est asymétrique à droite si l'aile (queue) droite de la distribution est plus longue que l'aile gauche.



Histogram

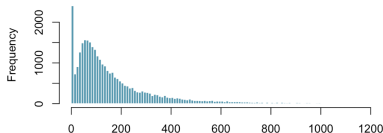
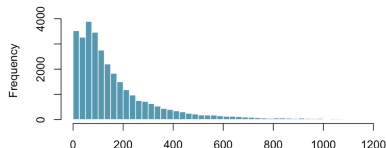
Données de Larry Brown et Haipeng Shen, durée des appels au service à la clientèle d'une banque pendant un mois : 31,492 appels

```
1 > hist(bankcall$Time)
2 > hist(bankcall$Time[bankcall$Time < 1200])
```



On peut se restreindre aux 31,247 appels de moins de 20 minutes

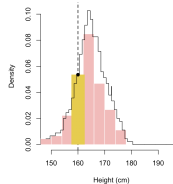
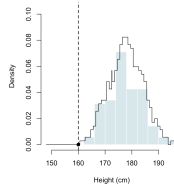
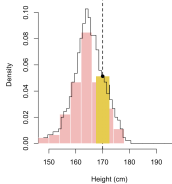
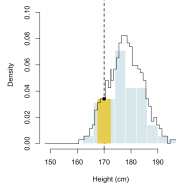
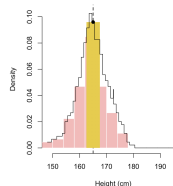
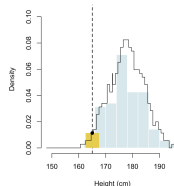
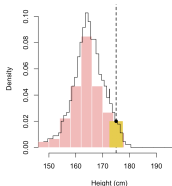
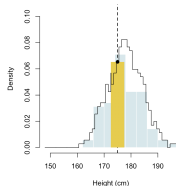
```
1 > hist(bankcall$Time[bankcall$Time < 1200] , breaks=seq(0,1200,by=10))
```



Moving Histogram

Histogramme glissant \hat{f}

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n \mathbf{1}(|x_i - x| \leq h/2)$$



Moving Histogram

\hat{F} cannot be differentiated, but we can consider

$$f_h(x) = \frac{1}{h} \underbrace{F(x + h/2) - F(x - h/2)}_{\mathbb{P}(X \in [x \pm h/2])}$$

i.e.

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n \mathbf{1}(x_i \in [x - h/2, x + h/2])$$

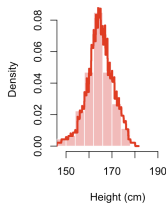
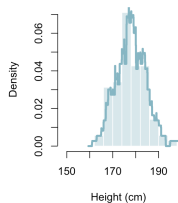
$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n \mathbf{1}(|x_i - x| \leq h/2)$$

One can prove that $\mathbb{E}(\hat{f}_h(x)) = f_h(x) \sim f(x) + \frac{h^2}{24} f''(x)$

i.e. $\text{bias}(\hat{f}_h(x)) \sim \frac{h^2}{24} f''(x)$, while $\text{Var}(\hat{f}_h(x)) \sim \frac{1}{nh} \cdot f_h(x)$

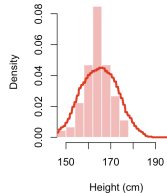
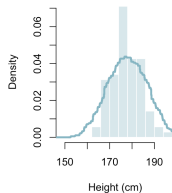
Moving Histogram

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n \mathbf{1}(|x_i - x| \leq h/2)$$



small h

bias $\text{bias}(\hat{f}_h(x))$ small
variance $\text{Var}(\hat{f}_h(x))$ large

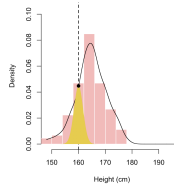
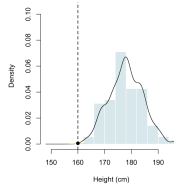
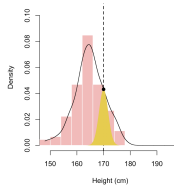
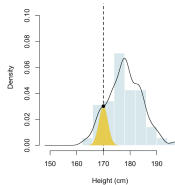
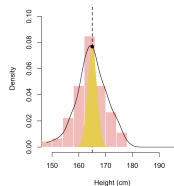
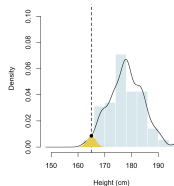
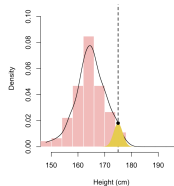
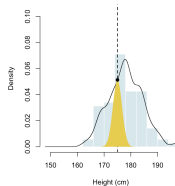


large h

bias $\text{bias}(\hat{f}_h(x))$ large
variance $\text{Var}(\hat{f}_h(x))$ small

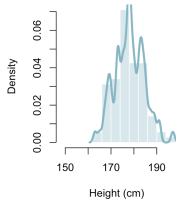
Kernel Density

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{x_i - x}{h}\right)$$



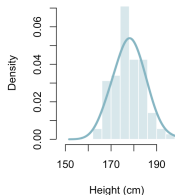
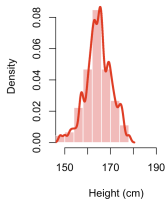
Kernel Density

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{x_i - x}{h}\right)$$



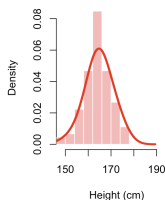
small h

bias $\text{bias}(\hat{f}_h(x))$ small
variance $\text{Var}(\hat{f}_h(x))$ large

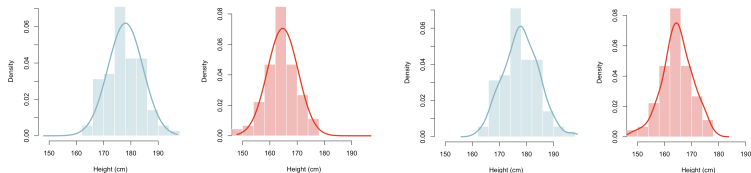


large h

bias $\text{bias}(\hat{f}_h(x))$ large
variance $\text{Var}(\hat{f}_h(x))$ small



Histogram & Density



```
1 > hist(x, probability=TRUE)
2 > plot(density(x))
3 > plot(density(x), kernel="gaussian", bw=1)
```