

# Méthodes statistiques – MST 2

ENSIIE

1<sup>ère</sup> année

2015/16

*Evry* : Nicolas Brunel

*Strasbourg* : Emmanuel Périnel, Nancy Rebout

# Sommaire

0. Introduction
1. Notions de base – Statistique descriptive
2. Bases de l'échantillonnage
3. L'estimation
4. L'estimation par intervalle de confiance
5. Les tests statistiques

## 0. Introduction

# 0. Introduction

## Plan

- Qu'est-ce que la statistique ?
- Domaines d'applications
- La démarche statistique
- Quelques exemples
- L'objet de ce cours

## 0. Introduction

# Qu'est-ce que la statistique ?

### Définition(s)

« Ensemble de données d'observations (une statistique) et activité (la statistique) qui consiste en leur **recueil**, leur **traitement** et leur **interprétation** »

« Compter, dénombrer, **résumer**, synthétiser des **données** afin de mieux comprendre des phénomènes, les expliquer, les modéliser, les prévoir »

- **À l'origine** : ensemble d'informations concernant la population et l'économie
- **Aujourd'hui** : branche des mathématiques appliquées à la frontière de disciplines scientifiques (mathématique et informatique) : théorie des probabilités, algèbre, théorie des graphes, algorithmique, machine learning, datamining (= fouille de données), big data, etc.

## 0. Introduction

# Domaines d'applications

Ils sont extrêmement **variés** !

- médecine
- démographie
- agriculture
- économie
- sociologie
- psychologie
- physique
- contrôle de qualité
- fiabilité
- enquête / sondage
- génomique
- écologie
- astronomie
- analyse sensorielle
- sport
- météorologie
- musique
- analyse de textes,
- web mining, etc.

## 0. Introduction

# La démarche statistique

Elles sont liées aux différentes phases du travail d'un statisticien

- **Recueil des données**  
Plan d'expérience, plan de sondage
- **Statistique descriptive et exploratoire**  
Préparation des données, représentations graphiques, indicateurs numériques, analyse bivariée et multivariée (liaisons entre variables)
- **Statistique inférentielle**  
Raisonnement à partir d'un échantillon, estimation d'une grandeur, qualité de l'estimation, test d'une hypothèse
- **Modélisation et prévision statistique**  
Expliquer / prévoir un phénomène à l'aide de modèles mathématiques

## 0. Introduction

# Quelques exemples

## 1. Construire un plan de sondage dans une enquête marketing

But : évaluer l'appréciation d'un nouveau produit par des consommateurs

- Quels sont les clients potentiels d'un produit (population étudiée) ?
- Quelle technique de sondage (quotas ? aléatoire ? stratifié ? par grappe ?)
- Comment interroger (téléphone ? internet ? auto administré ? En face à face ?)
- Combien de personnes doit-on / peut-on interroger (taille d'échantillon) ?
- Quelle est la précision attendue sur les résultats obtenus selon la taille ?
- Évaluer la représentativité de l'échantillon (faut-il redresser l'échantillon ?)

## 0. Introduction

## 2. Construire un plan d'expérience en agriculture

**But** : Comparer le rendement de variétés de blé

- Quelle est la variable réponse ? *Rendement de chaque variété*
- Quels sont les facteurs contrôlés ? *Les variétés, les doses de fertilisant*
- Quels sont les facteurs aléatoires ? *Hétérogénéité du sol, météo, etc.*
- Quel type de dispositif expérimental ? *Dispositif en blocs, randomisation totale, plans en carrés latins, split plot, criss-cross,  $\alpha$  - plans, etc.*
- Combien de variétés peut-on étudier au maximum ?
- Comment maîtriser les effets de bordure ou de voisinage ?
- etc.



## 0. Introduction

**3. Analyser le lien de dépendance entre deux caractères****But** : étude de la pérennité d'une union selon le type d'habitat*D'après Balakrishnan, 1986**Échantillon de 3864 couples, situation après 5 ans*

Situation / Habitat	Unis	Séparés	Total
rural	287	18	305
Petite ville	1124	89	1213
Grande ville	2081	265	2346
<b>Total</b>	<b>3492</b>	<b>372</b>	<b>3864</b>

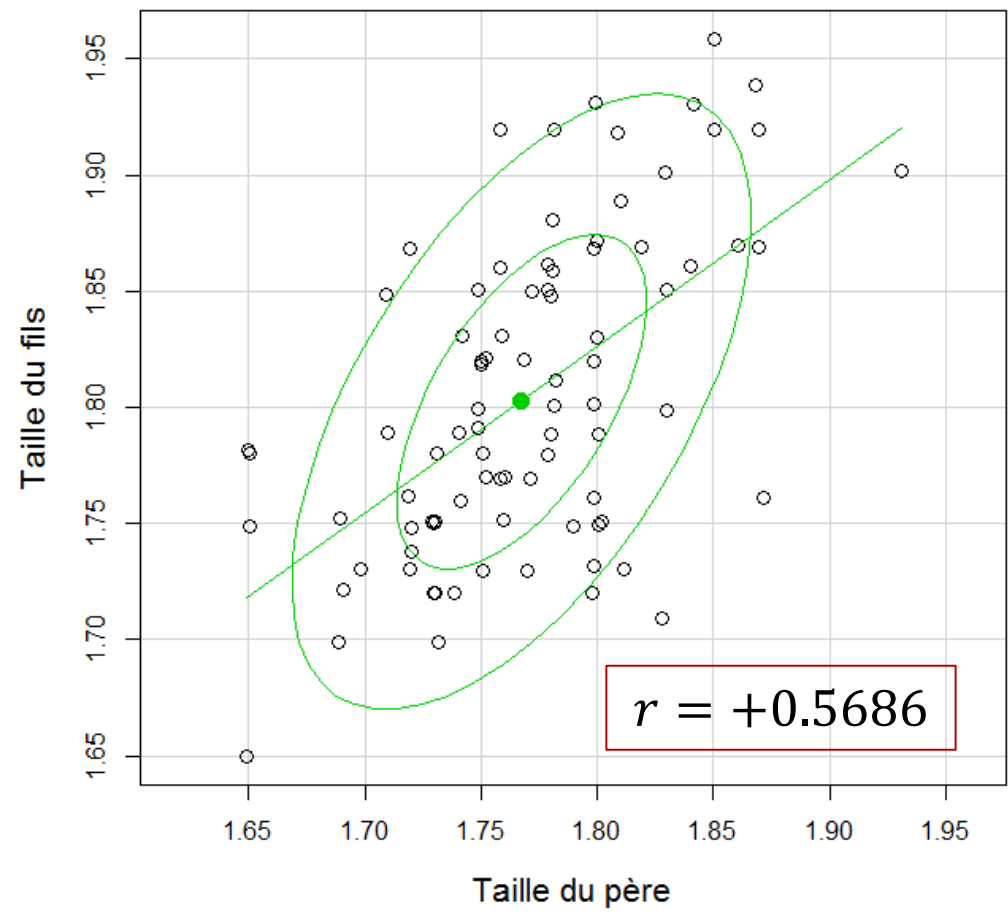
Test d'indépendance (Khi <sup>2</sup> ) :	
Khi <sup>2</sup> (Valeur observée)	19,685
Khi <sup>2</sup> (Valeur critique)	5,991
DDL	2
p-value	< 0,0001
alpha	0,05

*L'hypothèse d'indépendance est rejetée...*

0. Introduction

4. Analyse d'une corrélation linéaire

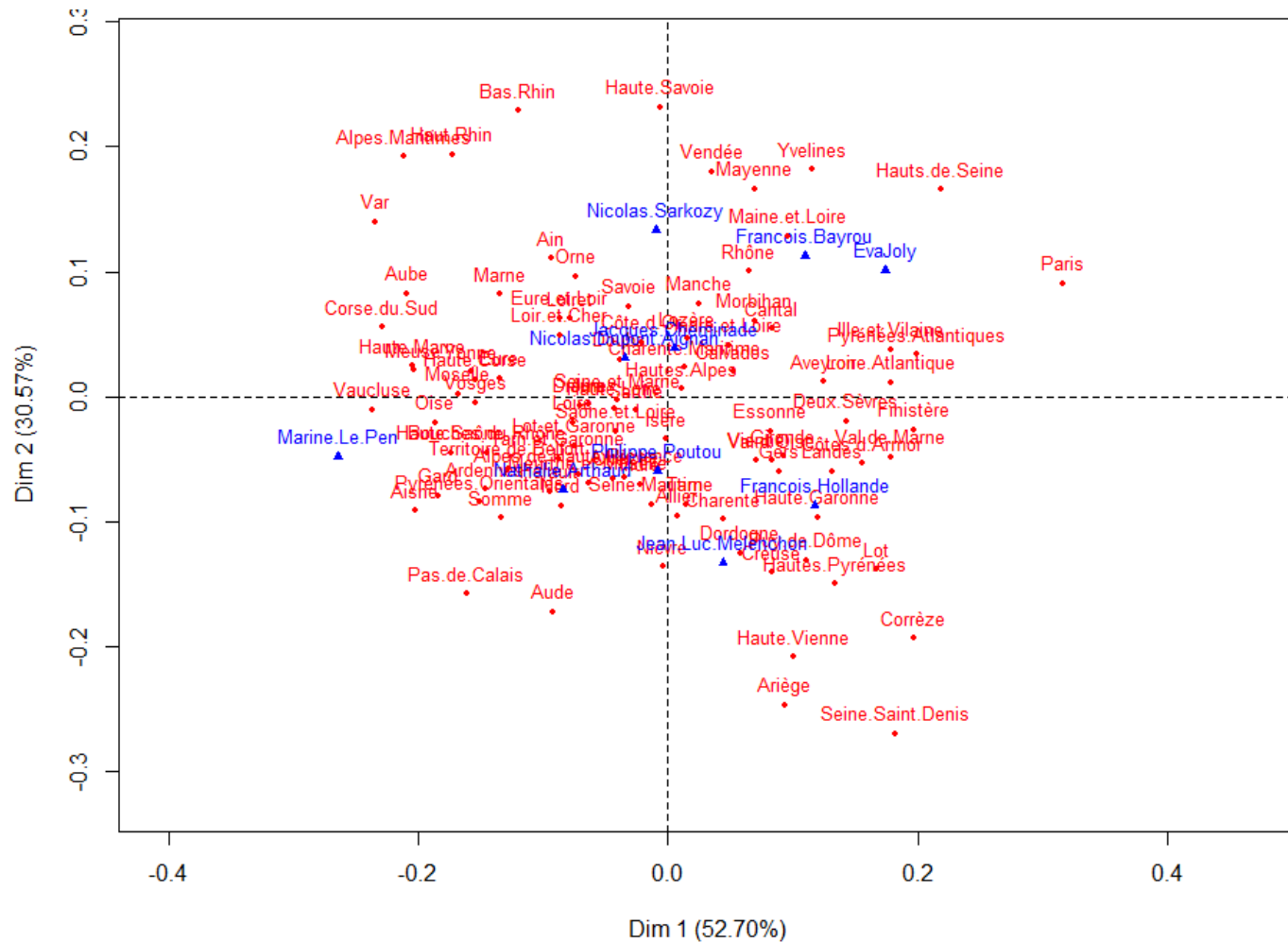
**But :** Dans quelle mesure la taille d'un enfant est-elle liée à celle de ses parents ?



0. Introduction

# 5. Analyse exploratoire multidimensionnelle (AFC)

**But :**  
 Comparer les résultats  
 de la présidentielle 2012  
 selon les départements



## 0. Introduction

# 6. Modélisation par analyse discriminante

**But** : Prévoir l'appartenance d'un individu à un groupe prédéfini

*D'après Saporta (Analyse discriminante, support de cours du CNAM)*

- **Exemple 1. Solvabilité d'emprunteurs auprès de banques – Crédit scoring**
  - Deux groupes : Client à risque (contentieux) ou non
  - Variables explicatives : taux d'endettement, revenu disponible du ménage, statut matrimonial, propriétaire/locataire, profession, ancienneté emploi, âge, nombre d'enfants, etc.
- **Exemple 2. Risque en assurance automobile**
  - Deux groupes : « Bon ou mauvais conducteur »
  - Variables explicatives : CSP, sexe, tranche d'âge, catégorie de véhicule, etc.

## 0. Introduction

### Exemple 1. Solvabilité d'emprunteurs

Principe : établir un score, fonction des variables explicatives, traduisant le niveau de solvabilité du client

$$\text{score} = \text{fonction}(X_1, X_2, \dots, X_p)$$

#### Exemple

Variable	Valeur – catégorie	Score
Ratio d'endettement	15%	+16
Revenu disponible par personne	2300 F	+12
Situation dans le logement	locataire	0
Etat matrimonial, nombre d'enfants	Marié sans enfant	+10
Ancienneté dans l'emploi	6 ans	+22

$$\text{score} = +60$$

## 0. Introduction

### Répartition du nombre de contentieux par tranche de score

<b>Tranche de score</b>	<b>Nbre de deman.</b>	<b>Nbre de conten.</b>	<b>Taux de conten.</b>
<b>90-100</b>	<b>1000</b>	<b>10</b>	<b>1 %</b>
<b>80-90</b>	<b>1500</b>	<b>35</b>	<b>2,3 %</b>
<b>70-80</b>	<b>1500</b>	<b>55</b>	<b>3,6 %</b>
<b>60-70</b>	<b>2000</b>	<b>80</b>	<b>4 %</b>
<b>50-60</b>	<b>2000</b>	<b>100</b>	<b>5 %</b>
<b>40-50</b>	<b>2000</b>	<b>140</b>	<b>7 %</b>
<b>30-40</b>	<b>2000</b>	<b>180</b>	<b>9 %</b>
<b>20-30</b>	<b>1000</b>	<b>110</b>	<b>11 %</b>
<b>10-20</b>	<b>1000</b>	<b>130</b>	<b>13 %</b>
<b>0-10</b>	<b>1000</b>	<b>160</b>	<b>16 %</b>

## 0. Introduction

# 7. Web mining – Text mining

### Web mining

Ensemble des techniques qui visent à explorer, traiter et analyser les grandes masses d'informations consécutives à une activité Internet

But = valoriser un site ; personnaliser un contenu selon le profil de l'utilisateur

Type de données traitées : contenu d'une page (textes, graphiques), sa structure, son usage (adresses IP, date, temps des requêtes), profil de l'utilisateur

### Text mining

Extraction de connaissances dans les textes ; spécialisation de la fouille de données (*data mining*) ; fait partie du domaine de l'intelligence artificielle

Applications : indexation de textes, détection d'anomalies, anti spam, mesure de ressemblance / coïncidence entre textes, etc.

## 0. Introduction

# L'objet de ce cours

- Notions de base, statistique **descriptive**
- Les fondements de la statistique **inférentielle**

### Un cours de statistique inférentielle / statistique mathématique

- **Échantillonnage** : déduire des renseignements sur un échantillon à partir de la connaissance de la population
- **Estimation** : déduire des renseignements sur une population à partir de la connaissance d'un échantillon
- Mettre en place un **test d'hypothèse**



## 1. Notions de base

# 1. Notions de base

## Plan

- Exemples de jeux de données
- Population, échantillon, individus
- Variables statistiques
- Le tableau *individus x variables*
- Représentations graphiques usuelles
- Principaux indicateurs numériques

## 1. Notions de base – Exemple de jeux de données

### 1.1 - Exemples de jeux de données

#### Exemple 1. Chevesne

Longueur (en mm) et poids (en grammes)  
de 20 chevesnes

longueur	poids
105	11
155	36
159	41
165	43
170	46
173	56
181	66
187	70
191	76
195	76
202	101
211	102
220	114
221	118
225	125
232	136
238	139
248	158
252	159
301	274

## 1. Notions de base – Exemple de jeux de données

### Exemple 2. Prison

Enquête par questionnaire auprès de 1500 adultes français sur le thème de la prison

Ident	âge	sexe	âge en classes	diplôme	PCS	orientation politique	détenus					travail souhaitable	respect des droits de l'homme		Poids
							Superficie cellule	par cellule	WC cloisonnés	préservatisés en prison	sexe lors des visites		conditions de détention	plutôt pas/dutout	
1	44	homme	âge (36,50)	DPS	ouvrier	centre droit	10	2	non	oui	dépend	oui	a.bonnes	plutôt	2,16
2	24	femme	âge (18,25)	DPS	employé	autre	6	2	non	oui	non	oui	a.mauvaises	pasdutout	0,39
3	37	femme	âge (36,50)	diplôme	,artisan	centre	6	2	non	non	non	oui	mauvaises	pasdutout	0,66
4	20	homme	âge (18,25)	BAC	inactif	droite	10	3	non	oui	oui	oui	bonnes	plutôt pas	0,66
5	75	femme	âge 66 et +	>BAC+2	profintemé diaires	centre gauche	6	4	non	non	?	oui	mauvaises	pasdutout	0,49
6	45	homme	âge (36,50)	DPS	,artisan	gauche	5	4	non	oui	oui	oui	a.mauvaises	plutôt pas	0,19
7	19	femme	âge (18,25)	BAC+2	inactif	gauche	10	4	non	non	non	oui	mauvaises	pasdutout	0,32
8	46	femme	âge (36,50)	BAC	profintemé diaires	centre gauche	10	4	?	oui	non	oui	bonnes	plutôt	0,28
9	30	homme	âge (26,35)	>BAC+2	cadre, profint.sup.	autre	6	4	non	oui	non	oui	mauvaises	pasdutout	0,11
...															
1497	25	homme	âge (18,25)	>BAC+2	inactif	centre gauche	15	6	non	?	non	oui	a.mauvaises	plutôt pas	0,72
1498	58	homme	âge (51,65)	CEP	,artisan	centre droit	7	8	non	oui	oui	oui	a.bonnes	plutôt	0,6
1499	25	homme	âge (18,25)	BAC	employé	centre gauche	5	6	non	?	oui	oui	a.mauvaises	plutôt pas	0,34
1500	34	femme	âge (26,35)	BAC	commerçant ,artisan	centre	10	4	oui	oui	non	oui	a.mauvaises	plutôt pas	0,6

## 1. Notions de base – Exemple de jeux de données

### Exemple 3. Poussins

Etude de l'effet de trois traitements sur le poids de 24 poussins mâle ou femelle

	Traitement 1	Traitement 2	Traitement 3
Mâles	25	21	23
	30	26	28
	26	22	24
	33	27	29
Femelles	15	16	15
	20	18	19
	18	17	17
	21	20	22

## 1. Notions de base – Exemple de jeux de données

### Exemple 4. Yaourts

Etude de la viscosité (en mPa.s) de quatre yaourts à deux dates différentes

produit	semaine	répétition	viscosité Brookfield (mPa.s)
F1	J+7	1	39400
F1	J+7	2	42000
F2	J+7	1	49200
F2	J+7	2	51000
F3	J+7	1	59400
F3	J+7	2	60800
F4	J+7	1	80000
F4	J+7	2	79600
F1	J+15	1	38800
F1	J+15	2	42400
F2	J+15	1	57800
F2	J+15	2	52600
F3	J+15	1	78200
F3	J+15	2	78000
F4	J+15	1	101400
F4	J+15	2	101600

## 1. Notions de base – Population, échantillon, individus

### 1.2 - Population, échantillon, individus, variables

#### Population et individu statistique

- Très souvent : notion **démographique, biologique** – écologique
- **En statistique** : ensemble des objets ou individus statistiques étudiés
- **Individu** statistique = **unité** statistique : très diverses !

#### Population *versus* échantillon

- **échantillon** : fraction, sous-ensemble de la population étudiée
- **Pourquoi** n'étudier qu'une fraction de la population ?
- Recensement *versus* sondage

## 1. Notions de base – Population, échantillon, individus

# Description et inférence statistique

## Description

- Résumer, synthétiser les résultats l'information contenue dans des données à l'aide de **graphiques** ou d'**indicateurs numériques**
- Les résultats obtenus se limitent aux individus observés

## Inférence

- **Extrapoler, généraliser** les résultats observés sur un échantillon à la population dans sa globalité
- Qualité essentielle attendue pour un échantillon : sa **représentativité**
- Un exemple typique : l'enquête par **sondage**
- Contexte **d'incertitude**. Importance de la théorie des **probabilités**
- Notions importantes : tests d'hypothèses, risque, confiance, prévision, estimation, probabilité critique, etc.

## 1. Notions de base – Variables statistiques

### 1.3 - Variables statistiques

- **Caractéristiques** utilisées pour décrire les individus de la population étudiée  
Informations recueillies sur les unités statistiques
- Terminologie différente selon les domaines d'application :  
*Variables, attributs, paramètres, descripteurs, facteurs, caractères, etc.*

Deux grands types de variables



**Variables quantitatives**

*(taille, salaire, température, etc.)*



**Variables qualitatives**

*(sexe, profession, habitat, etc.)*



## 1. Notions de base – Variables statistiques

### Variables **quantitatives**

Les valeurs prises sont des grandeurs mesurables, numériques

Elles se subdivisent en 2 types :

- **Continues** : observables sur un intervalle continu  
Nombre de valeurs possibles *a priori* infini
- **Discrètes** : prennent un nombre fini de valeurs  
(en général entières et peu nombreuses)

## 1. Notions de base – Variables statistiques

### Variables **qualitatives**

- Les valeurs prises par la variable sont des « qualités », non numériques, appelées **modalités** ou **catégories**
- Les modalités définissent des **sous populations** dans la population étudiée (hommes/femmes, rural/urbain, etc.)
- Variables **ordinales** (modalités ordonnées) ou **nominales** (aucun ordre)

### REMARQUES

- Comment **distinguer qualitative / quantitative** ?  
Opération arithmétique (min, max, somme, moyenne, etc.) possible?
- **Discret – Continu** : distinction parfois **arbitraire**...
- Une variable continue : presque autant de valeurs différentes que d'individus
- Souvent, seule la distinction *quanti / quali* est importante

## 1. Notions de base – Variables statistiques

### Variable qualitative et sous populations associées

produit	semaine	répétition	viscosité
F1	J+7	1	39400
F1	J+7	2	42000
F2	J+7	1	49200
F2	J+7	2	51000
F3	J+7	1	59400
F3	J+7	2	60800
F4	J+7	1	80000
F4	J+7	2	79600
F1	J+15	1	38800
F1	J+15	2	42400
F2	J+15	1	57800
F2	J+15	2	52600
F3	J+15	1	78200
F3	J+15	2	78000
F4	J+15	1	101400
F4	J+15	2	101600

produit	semaine	répétition	viscosité
F1	J+7	1	39400
F1	J+7	2	42000
F1	J+15	1	38800
F1	J+15	2	42400
F2	J+7	1	49200
F2	J+7	2	51000
F2	J+15	1	57800
F2	J+15	2	52600
F3	J+7	1	59400
F3	J+7	2	60800
F3	J+15	1	78200
F3	J+15	2	78000
F4	J+7	1	80000
F4	J+7	2	79600
F4	J+15	1	101400
F4	J+15	2	101600

1. Notions de base – Le tableau *individus x variables*

## 1.4 - Le tableau *individus x variables*

**Principal format de données** sous lequel les informations sont saisies afin d'être traitées par un logiciel statistique

		variables				
		$x_1$	$\dots$	$x_j$	$\dots$	$x_p$
individus	1	$x_{11}$		$x_{1j}$		$x_{1p}$
	$\vdots$			$\vdots$		
	$i$	$x_{i1}$	$\dots$	$x_{ij}$	$\dots$	$x_{ip}$
	$\vdots$			$\vdots$		
	$n$	$x_{n1}$		$x_{nj}$		$x_{np}$

## 1. Notions de base – Le tableau *individus x variables*

### Identificateur des individus

- Identificateur = nom des individus
- Obligatoirement tous différents !
- Appelé aussi « *libellé des observations* », « *nom des cas* » (logiciel R – Package Rcmdr)

produit	semaine	répétition	viscosité Brookfield (mPa.s)
F1	J+7	1	39400
F1	J+7	2	42000
F2	J+7	1	49200
F2	J+7	2	51000
F3	J+7	1	59400
F3	J+7	2	60800
F4	J+7	1	80000
F4	J+7	2	79600
F1	J+15	1	38800
F1	J+15	2	42400
F2	J+15	1	57800
F2	J+15	2	52600
F3	J+15	1	78200
F3	J+15	2	78000
F4	J+15	1	101400
F4	J+15	2	101600

*Individus anonymes*

Athlète	100m	Longueur	Poids	Hauteur	400m
Sebrle	10,85	7,84	16,36	2,12	48,36
Clay	10,44	7,96	15,23	2,06	49,19
Karpov	10,5	7,81	15,93	2,09	46,81
Macey	10,89	7,47	15,73	2,15	48,97
Warners	10,62	7,74	14,48	1,97	47,97
Zsivoczky	10,91	7,14	15,31	2,12	49,4
Hernu	10,97	7,19	14,65	2,03	48,73
Nool	10,8	7,53	14,26	1,88	48,81
Bernard	10,69	7,48	14,8	2,12	49,13
Schwarzl	10,98	7,49	14,01	1,94	49,76
Pogorelov	10,95	7,31	15,1	2,06	50,79
Schoenbeck	10,9	7,3	14,77	1,88	50,3
Barras	11,14	6,99	14,91	1,94	49,41
Smith	10,85	6,81	15,24	1,91	49,27



*Nom des observations*

## 1. Notions de base – Le tableau *individus x variables*

*Quel format pour un tableau individus x variables ?*

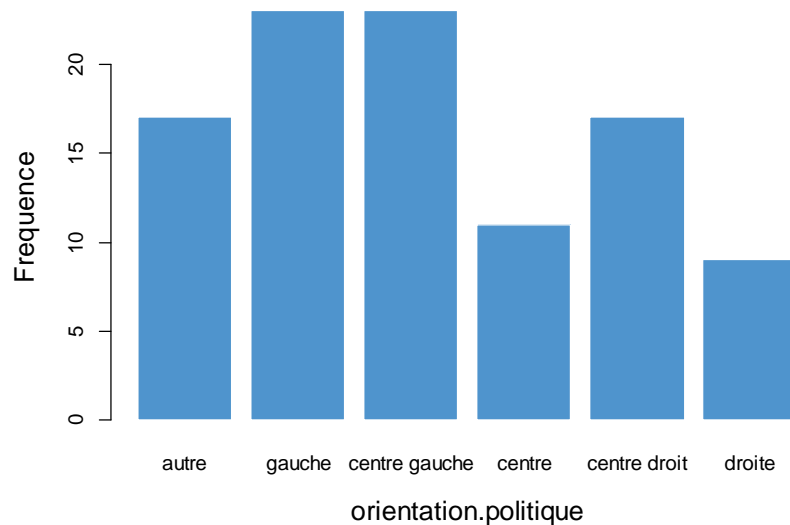
	Traitement 1	Traitement 2	Traitement 3
Mâles	25	21	23
	30	26	28
	26	22	24
	33	27	29
Femelles	15	16	15
	20	18	19
	18	17	17
	21	20	22

## 1. Notions de base – Représentations graphiques

### 1.5 - Représentations graphiques usuelles

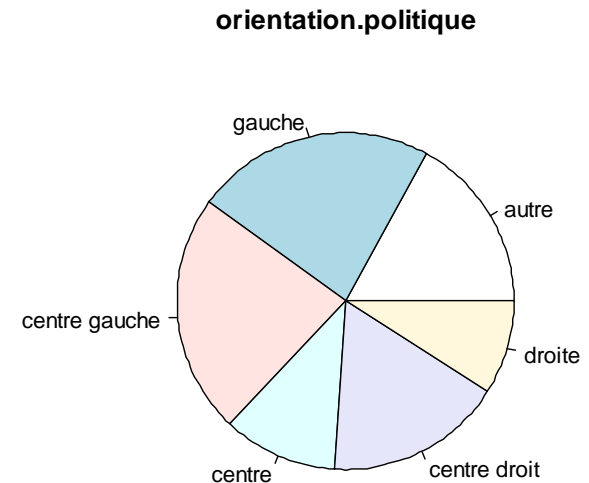
#### Pour une variable qualitative

Diagramme en bâton



```
barplot(table(genepi$orientation.politique),  
xlab="orientation.politique", cex.lab=1.3,  
ylab="Frequence", border="white", col="steelblue3")
```

Diagramme en secteurs - Camembert

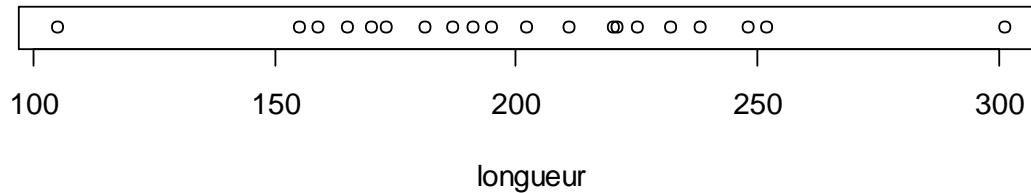


```
pie(table(genepi$orientation.politique),  
main="orientation.politique")
```

## 1. Notions de base – Représentations graphiques

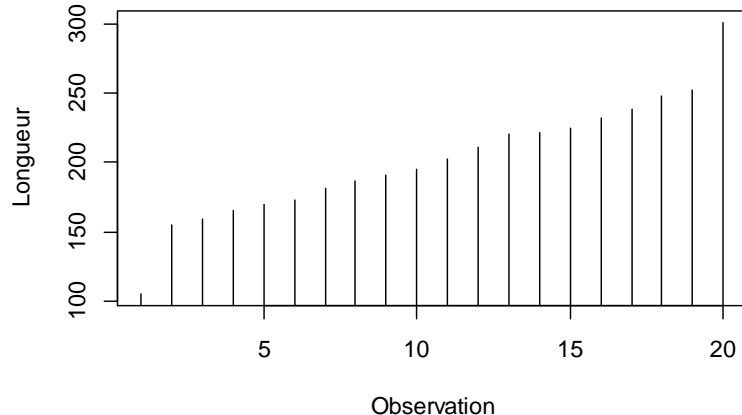
### Pour une variable quantitative – *Représentation des valeurs individuelles*

#### Exemple : longueur d'un chevesne



**Représentation axiale**  
Nuage de points  
Graphique « en bande »

```
stripchart(chevesne$longueur,  
pch=1, xlab="longueur")
```



**Graphique indexé**

```
plot(chevesne$longueur, type="h", ylab="Longueur",  
xlab="Observation")
```



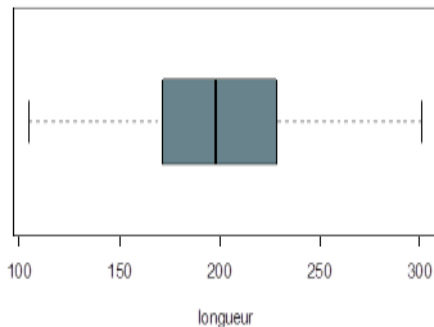
## 1. Notions de base – Représentations graphiques

### Pour une variable quantitative – *Représentation synthétique*



#### Histogramme

```
hist(chevesne$longueur, nclass=6,
     col="lightblue4", border="white", main =
     "Longueur d'un chevesne", ylab="Fréquence",
     xlab="Longueur")
```

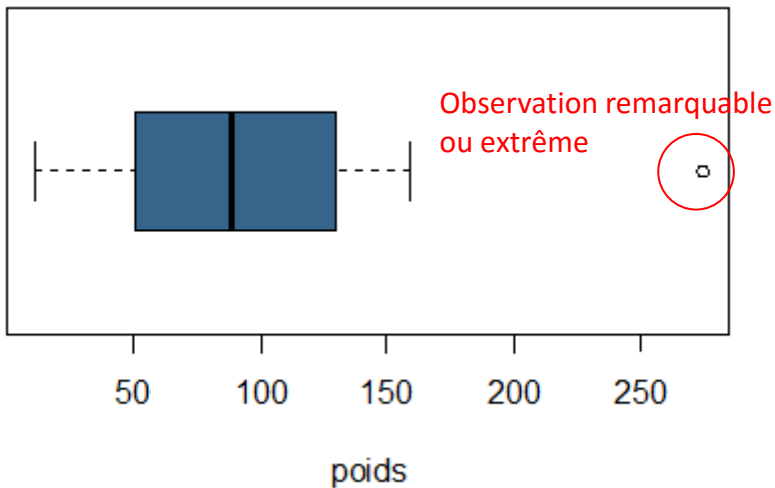
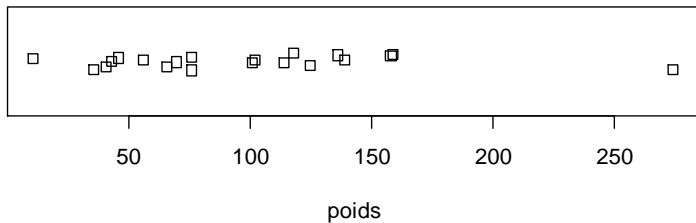


#### Boîte à moustaches - *boxplot*

```
boxplot(chevesne$longueur, xlab="longueur",
        col="lightblue4", horizontal=T)
```

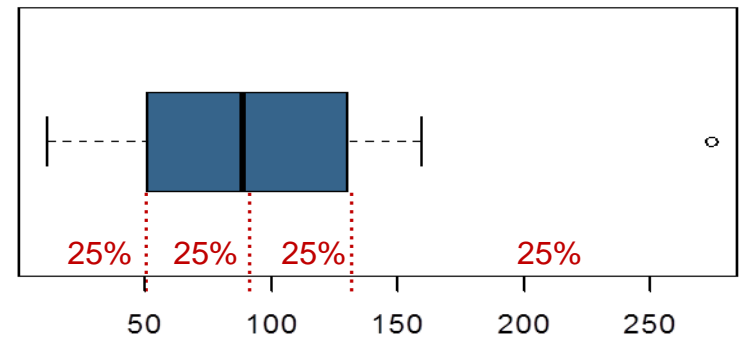
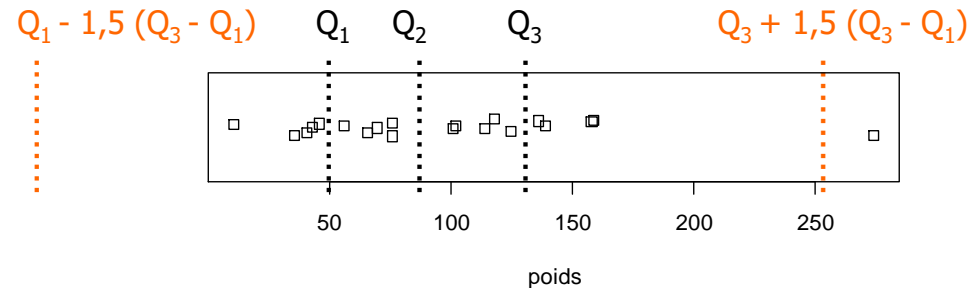
# 1. Notions de base – Représentations graphiques

Pour la variable *Poids*



## Construction d'un box plot

Représentation basée sur les *quartiles*



Les observations qui s'écartent du bord de la boîte de plus d'une fois et demi la longueur de la boîte sont considérées comme « remarquables », « extrêmes »

# 1. Notions de base – Représentations graphiques

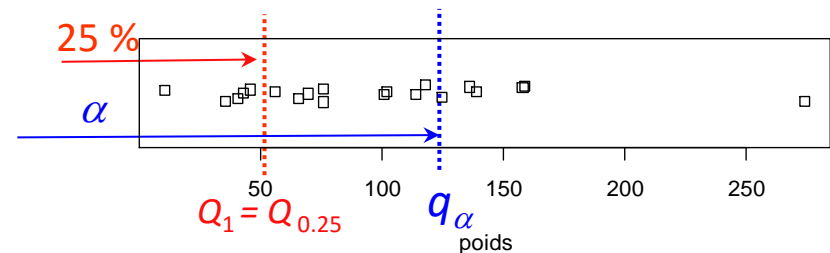
## Quantile d'ordre $\alpha$

- C'est la valeur d'une variable, notée  $q_\alpha$ , associée à une **fréquence cumulée =  $\alpha$**
- « Le pourcentage de valeurs inférieures à  $q_\alpha$  est égal à  $\alpha$  »
- Les quartiles sont des quantiles particuliers

$Q1 = Q_{0.25}$  = quantile d'ordre 25 %

$Q2 = Q_{0.50}$  = quantile d'ordre 50 %

$Q3 = Q_{0.75}$  = quantile d'ordre 75 %



- Définition générale d'un quantile

$$q_\alpha = \inf \{x \mid F_X(x) \geq \alpha\}$$

## 1. Notions de base – Indicateurs numériques

# Indicateurs numériques

### Pour des variables quantitatives

- Indicateurs de **tendance centrale**
- Indicateurs de **dispersion**

### Pour des variables qualitatives

- Tableaux **d'effectifs** ou de **fréquences**
- Mode

## 1. Notions de base – Indicateurs de tendance centrale

### 1.6 – Indicateurs de tendance centrale

Pour des variables quantitatives

#### Illustration

- Population : 25 poudres de lait
- Variable MAT/MST  
*Teneur en protéine / Matière sèche*

N° Poudre lait	MAT/MST
17	82,79
22	82,96
14	83,17
21	83,92
11	84,57
20	84,65
25	85,02
19	85,14
13	85,34
12	85,62
16	85,68
24	85,7
23	85,77
15	86,73
9	87,4
8	87,97
10	88,24
1	88,44
7	89,06
6	89,63
3	89,88
2	90,17
4	91,64
5	92,21
18	97,06

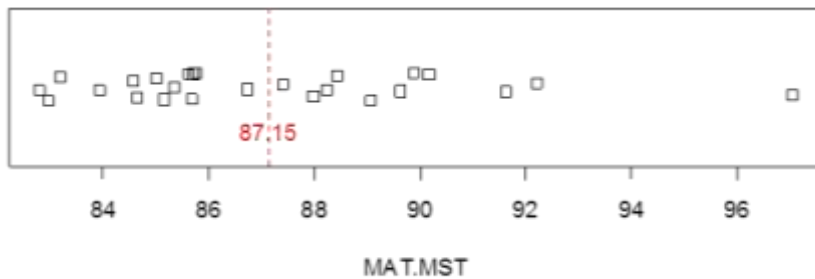
# 1. Notions de base – Indicateurs de tendance centrale

## La moyenne arithmétique

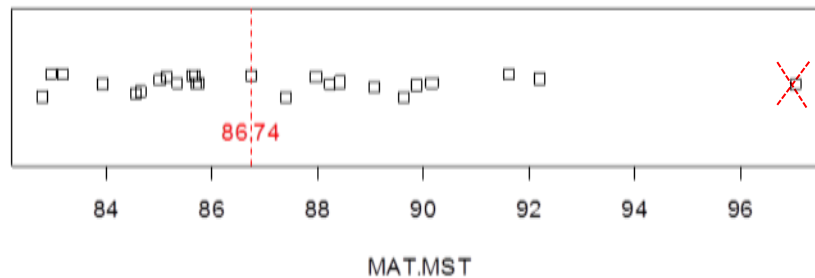
Une série statistique de  $n$  valeurs  $(x_1, x_2, \dots, x_i, \dots, x_n)$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

- Indicateur « universel »
- Sensibilité aux valeurs extrêmes

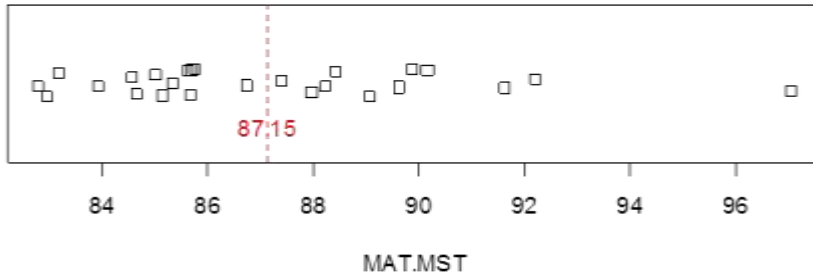


Moyenne (avec) = 87,15

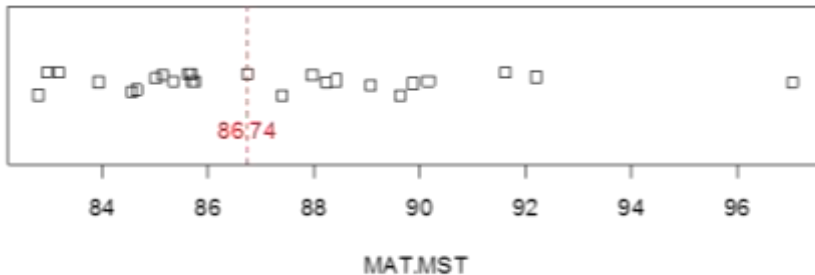


Moyenne (sans) = 86,74

## 1. Notions de base – Indicateurs de tendance centrale



```
stripchart(MAT.MST, method="jitter", xlab="MAT.MST")
text(mean(MAT.MST), 0.7, round(mean(MAT.MST), 2), col="red")
abline(v=mean(MAT.MST), col="red", lty=2)
```

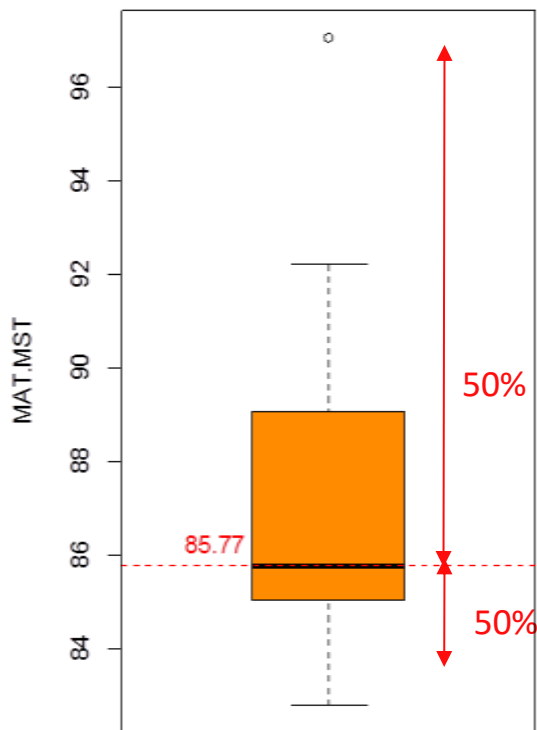


```
stripchart(MAT.MST, method="jitter", xlab="MAT.MST")
text(mean(MAT.MST[1:24]), 0.7, round(mean(MAT.MST[1:24]), 2), col="red")
abline(v=mean(MAT.MST[1:24]), col="red", lty=2)
```

# 1. Notions de base – Indicateurs de tendance centrale

## La médiane

Une série statistique de  $n$  valeurs rangées  $x_{(1)} < x_{(2)} < \dots < x_{(i)} < \dots < x_{(n)}$



$$Me = x_{\left(\frac{n+1}{2}\right)} \quad n \text{ impair}$$

$$Me = \frac{x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n}{2}+1\right)}}{2} \quad n \text{ pair}$$

*La médiane partage la série en deux parties égales*

- Indicateur « robuste »
- peu sensible aux valeurs extrêmes

Médiane (avec) = 85,77

Médiane (sans) = 85,74



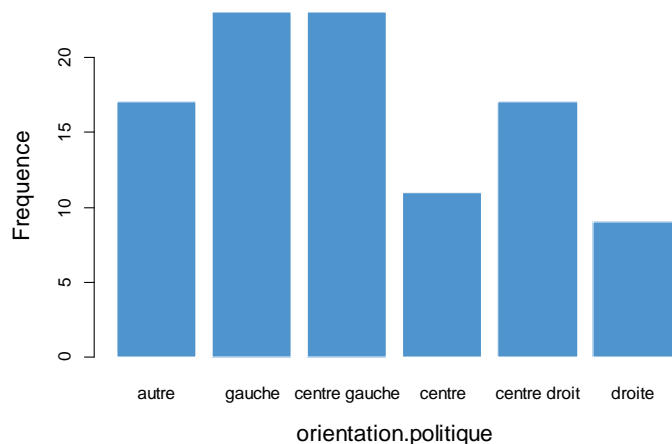
## 1. Notions de base – Tableau de fréquences, Mode

### Pour des variables qualitatives

#### Tableau de fréquences

Tableau associant à chaque modalité d'une variable qualitative, sa fréquence (ou son effectif) observé dans l'échantillon

autre	gauche	centre gauche	centre	centre droit	droite
17	23	23	11	17	9



#### Le Mode

C'est la valeur de la variable **la plus fréquente** (ou d'effectif maximum)  
Déterminé en général pour des *variables qualitatives*

1. Notions de base – Indicateurs de dispersion

# 1.7 – Indicateurs de dispersion

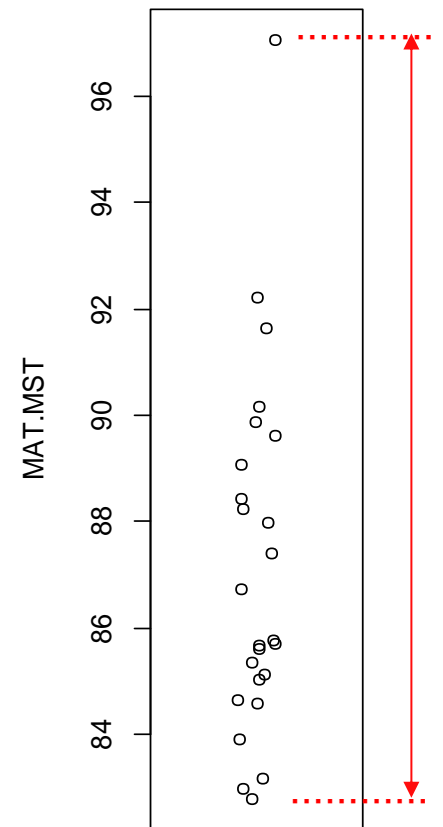
## L'étendue

Une série statistique de  $n$  valeurs  $(x_1, x_2, \dots, x_i, \dots, x_n)$

Étendue = Amplitude =  $(x_{\max} - x_{\min})$

- Le plus simple et le plus intuitif
- Très sensible aux valeurs extrêmes !
- Seules deux valeurs de la série participent au calcul

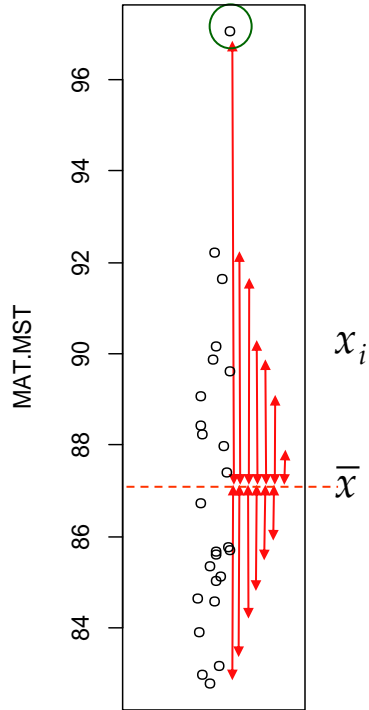
Étendue =  $97,06 - 82,79 = 14,27$



## 1. Notions de base – Indicateurs de dispersion

### L'écart absolu moyen : EAM

Une série statistique de  $n$  valeurs  $(x_1, x_2, \dots, x_i, \dots, x_n)$



$$EAM = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|$$

« Écart moyen à la moyenne »

« Dispersion moyenne autour de la moyenne »

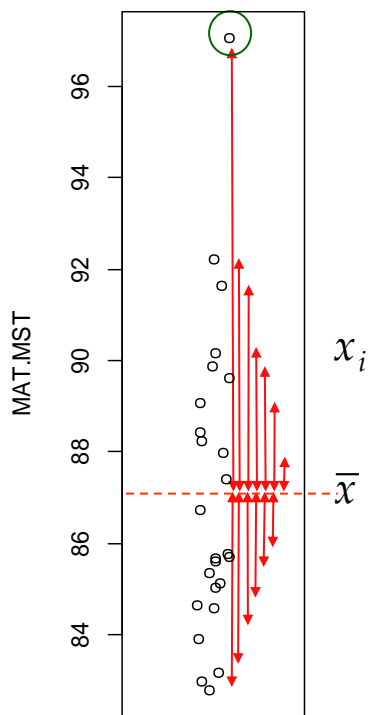
EAM = 2,64 (avec  $\bigcirc$ )

EAM = 2,27 (sans  $\bigcirc$ )

# 1. Notions de base – Indicateurs de dispersion

## La variance : $s^2$

Une série statistique de  $n$  valeurs  $(x_1, x_2, \dots, x_i, \dots, x_n)$



$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

« Moyenne des écarts à la moyenne au carré »

- Pas d'interprétation concrète (unité de mesure au carré)
- Importance fondamentale en statistique (nombreuses propriétés mathématiques)

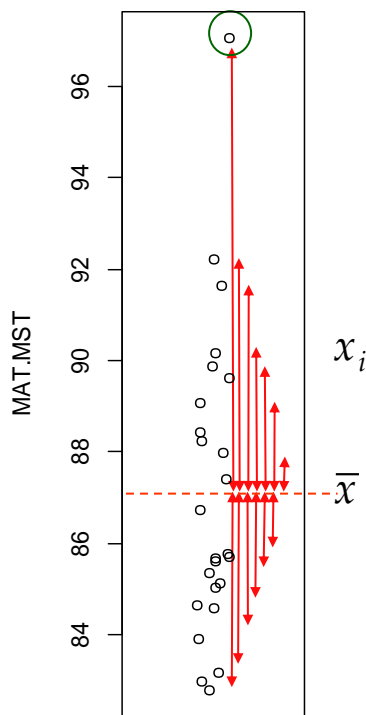
$s^2 = 10,81$  (avec  $\bigcirc$ )

$s^2 = 7,00$  (sans  $\bigcirc$ )

# 1. Notions de base – Indicateurs de dispersion

## L'écart-type : $s$

Une série statistique de  $n$  valeurs  $(x_1, x_2, \dots, x_i, \dots, x_n)$



$$s = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

« Écart moyen à la moyenne » (par abus...)

- Souvent proche de l'EAM (mais supérieur)
- Plus sensible aux valeurs extrêmes
- Interprétable dans l'unité de mesure étudiée

$s = 3,29$  (avec )

$s = 2,65$  (sans )

## 1. Notions de base – Indicateurs de dispersion

### Le coefficient de variation : CV

Une série statistique de  $n$  valeurs  $(x_1, x_2, \dots, x_i, \dots, x_n)$

$$CV = \frac{s}{\bar{x}}$$

« écart-type normalisé / standardisé »

« écart-type en pourcentage de la moyenne »

#### Quel intérêt ?

- **Comparer des dispersions** entre elles
- Le CV permet de comparer la *dispersion de variables* ayant :
  - *des unités de mesure différentes*
  - *des moyennes différentes*

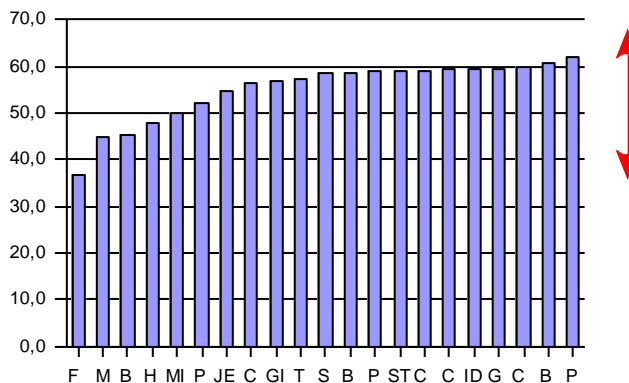
# 1. Notions de base – Indicateurs de dispersion

## Exemple. Mesure cornéométrique à deux dates

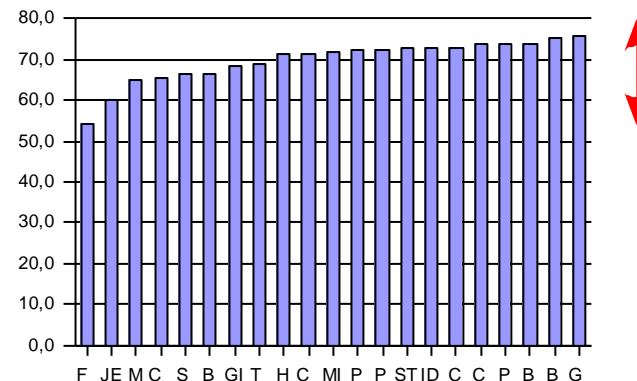
	J0	J30
JEG	54,7	60,1
PUJ	62,1	72,5
CHA	56,3	73,0
STR	58,9	72,7
BER	45,3	66,6
GIR	56,9	68,6
POM	52,2	72,4
CLA	59,8	65,6
IDR	59,5	72,8
MEN	45,0	64,9
CAT	59,0	71,5
SER	58,4	66,4
BERG	58,5	75,2
MIN	50,1	71,9
BUD	60,6	73,7
TRE	57,1	68,9
FEI	36,6	54,2
CEY	59,3	73,6
GUI	59,5	75,9
POU	58,9	73,7
HAB	48,1	71,2

écart-type	6,38	5,19	-18,70%
moyenne	55,09	69,78	
CV	0,12	0,07	-35,80%

Cornéo - Traité - J0



Cornéo - Traité - J30



## 1. Notions de base – Cas de données groupées

### Le cas de données « groupées »

#### Exemple

Examen anticipé d'Analyse des Données  
16 Décembre 2015  
- Documents interdits - Calculatrice autorisée.

**Exercice 1** Une entreprise est constituée de deux usines, appelées *A* et *B*. Le tableau suivant récapitule les salaires en euros par catégorie de personnel et par usine :

Usine A	Salaires	Effectifs	Usine B	Salaires	Effectifs
Ouvriers	700	200	Ouvriers	900	60
Employés	1400	20	Employés	1600	40
Cadres	5300	10	Cadres	7300	20

1. Calculer la moyenne des salaires dans chacune des usines, dans l'entreprise. Vérifier que la moyenne des salaires dans l'entreprise est la moyenne des salaires moyens de chaque usine.
2. Calculer la moyenne des salaires des ouvriers, puis des employés et enfin des cadres dans l'entreprise.
3. Calculer la variance des salaires dans chacune des usines et dans l'entreprise.
4. Vérifier que la variance des salaires dans l'entreprise est égale à la moyenne des variances des usines augmentée de la variance des moyennes calculées dans chaque usine. Quelle est la propriété du cours illustrée ici ?



## 1. Notions de base – Données centrées réduites

### Les données centrées réduites

#### Intérêt ?

- Évaluer le caractère **remarquable** / **extrême** d'une valeur dans une série statistique
- Comparer des données exprimées dans des **unités de mesure différentes**

#### Données initiales

Une série statistique de  $n$  valeurs

$$(x_1, x_2, \dots, x_i, \dots, x_n)$$

#### Données centrées réduites

$$\frac{(x_1 - \bar{x})}{s_x}, \dots, \frac{(x_i - \bar{x})}{s_x}, \dots, \frac{(x_n - \bar{x})}{s_x}$$

centrage

réduction

## 1. Notions de base – Données centrées réduites

**Exemple.** Données météorologique à Strasbourg, le 8 février

Année	minimale (°C)	maximale (°C)	ensoleillement (heures)	précipitations (mm)
2006	2,20	5,30	0,20	1,60
2007	2,10	10,70	2,20	3,00
2008	-2,20	9,40	9,20	0,00
2009	1,30	3,60	0,10	0,00
2010	-1,20	2,40	1,80	0,00
2011	0,40	8,50	0,70	0,20
2012	-9,20	-2,80	7,20	0,00
2013	-0,80	3,60	0,70	0,40
2014	0,60	8,00	0,00	2,20
2015	-4,90	4,10	0,40	0,40
moyenne	-1,17	5,28	2,25	0,78
écart type	3,37	3,80	3,08	1,03

Question : Quelle est l'observation la plus « remarquable », la plus « extrême » ?

## 1. Notions de base – Données centrées réduites

- Une précipitation de 3 mm se situe à 2,15 écart-type au dessus de sa moyenne  
Données « brutes » :  $3,00 = 0,78 + 2,15 \times 1,03$   
Données C-R :  $2,15 = 0 + 2,15 \times 1$
- La température de  $-9,2^{\circ}\text{C}$  se situe à 2,39 écart type sous la moyenne

Année	minimale (°C)	maximale (°C)	ensoleillement (heures)	précipitations (mm)
2006	1,00	0,01	-0,66	0,79
2007	0,97	1,42	-0,02	2,15
2008	-0,31	1,08	2,25	-0,75
2009	0,73	-0,44	-0,70	-0,75
2010	-0,01	-0,76	-0,15	-0,75
2011	0,47	0,85	-0,50	-0,56
2012	-2,39	-2,12	1,60	-0,75
2013	0,11	-0,44	-0,50	-0,37
2014	0,53	0,71	-0,73	1,37
2015	-1,11	-0,31	-0,60	-0,37
moyenne	0,00	0,00	0,00	0,00
écart type	1,00	1,00	1,00	1,00

### Une autre interprétation des données C-R

$(x_i - \bar{x})$  Écart de l'observation (*i*) à sa moyenne

$s_x$  Écart moyen à la moyenne

## 1. Notions de base – Données centrées réduites

### Propriétés des données centrées - réduites

Les données centrées réduites :

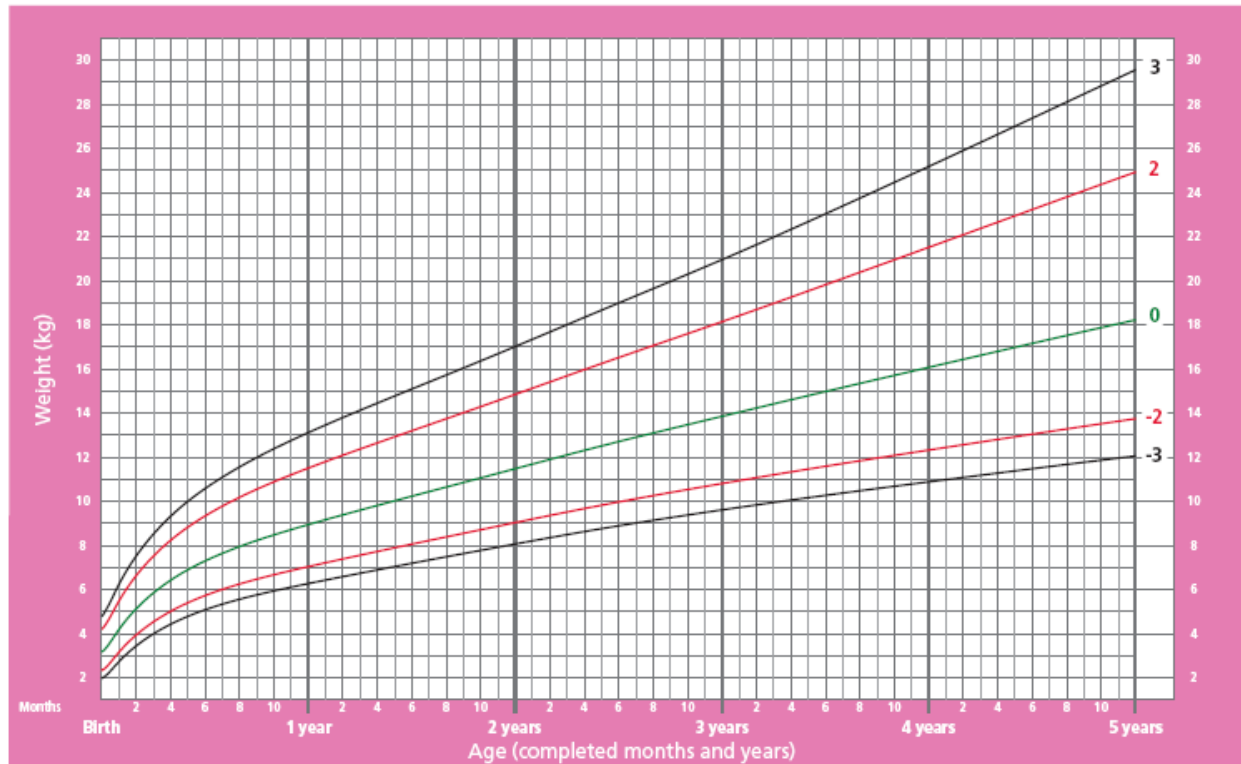
- sont :
  - *de moyenne* = 0 (conséquence du centrage)
  - *d'écart-type* = 1 (conséquence de la réduction)
- s'expriment en **nombre d'écart-type**
- sont **indépendantes de l'unité de mesure** de la variable  
On parle également de données « normalisées » ou « standardisées »

# 1. Notions de base – Données centrées réduites

## Evolution du poids des filles entre 0 et 5 ans

### Weight-for-age GIRLS

Birth to 5 years (z-scores)



WHO Child Growth Standards

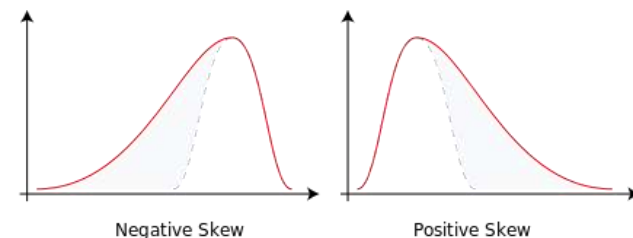
1. Notions de base – Indicateurs de forme

# 1.8 – Indicateurs de forme

## Coefficient d'asymétrie ou *skewness*

- C'est le moment centré d'ordre 3, normalisé
- Référence pour une loi normale :  $\gamma_1 = 0$ 
  - $\gamma_1 > 0$  : asymétrie à droite
  - $\gamma_1 < 0$  : asymétrie à gauche

$$\gamma_1 = \frac{E \left[ (X - E(X))^3 \right]}{\sigma^3}$$



source : Wikipedia

- Estimateur **sans biais**

$$G_1 = \frac{n}{(n-1)(n-2)} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s'_x} \right)^3$$

## 1. Notions de base – Indicateurs de forme

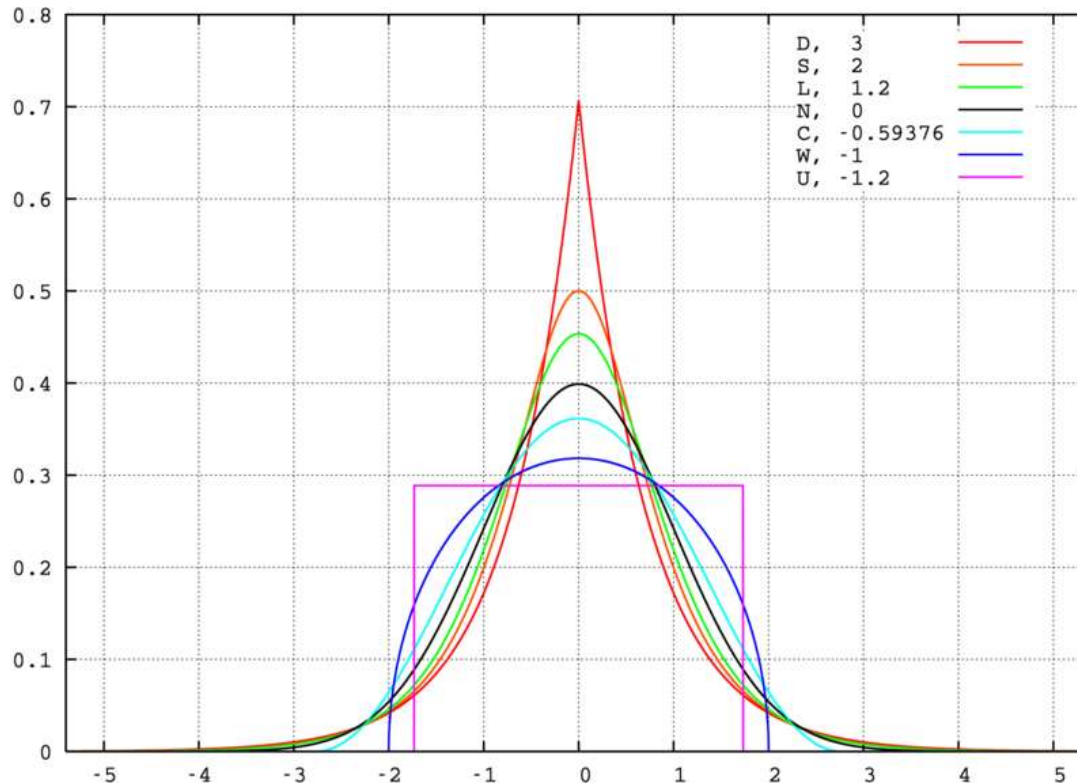
### Coefficient d'aplatissement ou *kurtosis*

- C'est le moment centré d'ordre 4, normalisé  $\gamma_2 = \frac{E[(X - E(X))^4]}{\sigma^4}$
- Référence pour une loi normale :  $\gamma_2 = 3$ 
  - $\gamma_2 > 3$  : queues de distribution plus « épaisses », plus pointue en sa moyenne
  - $\gamma_2 < 3$  : queues de distribution plus « fines »
- Le coefficient est souvent exprimé en termes d'écart à la valeur 3 (*excess kurtosis*)

- Estimateur sans biais  $G_2 = \frac{n(n+1)}{(n-1)(n-2)(n-3)} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s'_x} \right)^4 - 3 \frac{(n-1)^2}{(n-2)(n-3)}$

## 1. Notions de base – Indicateurs de forme

**Illustration** : valeur de :  $(\gamma_2 - 3)$  pour quelques distributions



- Distribution *leptokurtique*  
 $\gamma_2 - 3 > 0$
- Distribution *platikurtique*  
 $\gamma_2 - 3 < 0$
- Loi normale = *mésokurtique*

source : Wikipedia

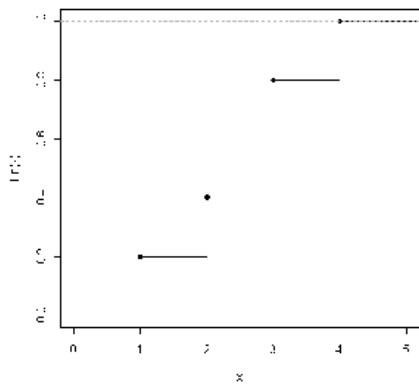


## 1. Notions de base – Fonction de répartition empirique

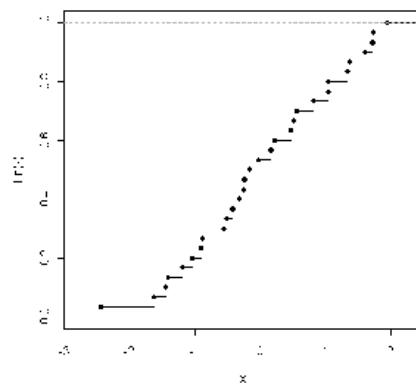
## 1.9 – Fonction de répartition empirique

- Soit  $\{x_1, x_2, \dots, x_n\}$  un échantillon de taille  $n$
- On définit la fonction de répartition empirique pour tout  $x \in \mathbb{R}$

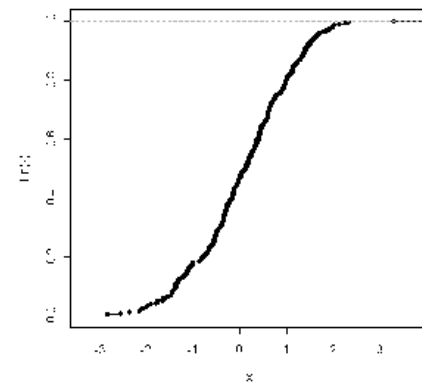
$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n 1_{\{x_i \leq x\}}$$



```
> x=c(1,2,3,3,4)
> plot(ecdf(x))
```



```
> x=rnorm(30)
> plot(ecdf(x))
```



```
> x=rnorm(300)
> plot(ecdf(x))
```

Illustration  
(logiciel R)

## 1. Notions de base – Fonction de répartition empirique

### Inégalité Dvoretzky – Kiefer - Wolfowitz

- Soit  $F_X$  la fonction de répartition de la variable  $X$  de la population
- Soit  $\hat{F}_n$  la fonction de répartition empirique pour un échantillon de taille  $n$

$$\forall \varepsilon > 0, \quad P\left(\sup_{x \in \mathbb{R}} |F_X(x) - \hat{F}_n(x)| > \varepsilon\right) \leq 2e^{-2n\varepsilon^2}$$

## 1. Notions de base – La loi normale

### 1.10 – La loi Normale ou de Laplace-Gauss

#### Pourquoi est-elle si importante ?

1. Loi de probabilité très utilisée pour **modéliser des phénomènes** naturels (phénomènes résultant de l'addition de multiples aléas indépendants)
2. Loi des **erreurs**
3. Loi d'une **moyenne d'échantillon** ; loi **limite** de certaines distributions
4. À l'origine de la définition des nombreuses **autres lois** (ex. : *Student, Fisher, Khi2*)
5. La normalité des données est une **hypothèse nécessaire** pour la réalisation de nombreuses analyses statistiques

## 1. Notions de base – La loi normale



Carl Friedrich Gauss  
(1777 - 1855)

Pierre Simon de Laplace  
(1749 - 1827)



THE  
NORMAL  
LAW OF ERROR  
STANDS OUT IN THE  
EXPERIENCE OF MANKIND  
AS ONE OF THE BROADEST  
GENERALIZATIONS OF NATURAL  
PHILOSOPHY ♦ IT SERVES AS THE  
GUIDING INSTRUMENT IN RESEARCHES  
IN THE PHYSICAL AND SOCIAL SCIENCES AND  
IN MEDICINE AGRICULTURE AND ENGINEERING ♦  
IT IS AN INDISPENSABLE TOOL FOR THE ANALYSIS AND THE  
INTERPRETATION OF THE BASIC DATA OBTAINED BY OBSERVATION AND EXPERIMENT

Citation de  
William Youden  
(*source : Wikipedia*)

« La loi normale des erreurs se distingue dans l'expérience de l'humanité comme une des plus larges généralisations de la philosophie naturelle ♦ Elle sert de guide dans la recherche en sciences physique et sociale, en médecine, en agriculture et en ingénierie ♦ C'est un outil indispensable pour l'analyse et l'interprétation des données de base obtenues par l'observation et l'expérience. »

## 1. Notions de base – La loi normale

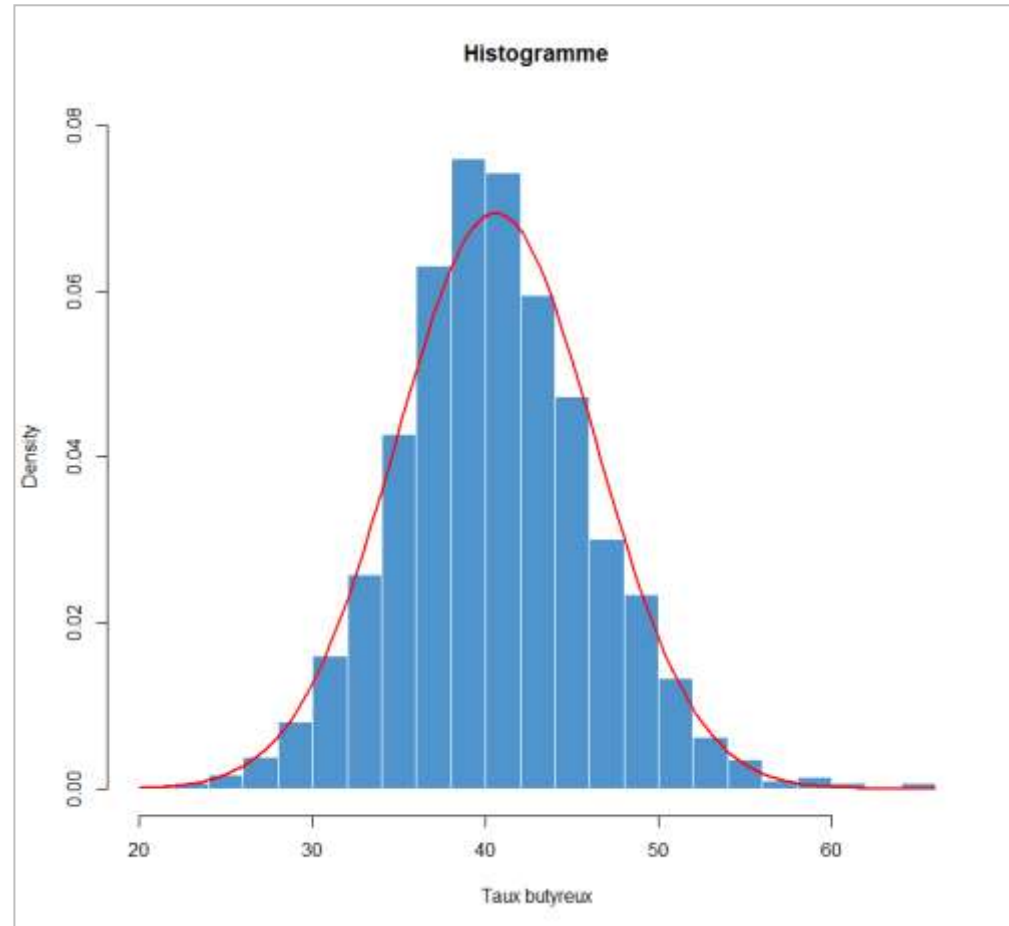
### La loi de LG permet de modéliser un grand nombre de phénomènes

#### Exemple 1.

Taux butyreux du lait de  
1428 vaches montbéliardes

```
hist(vache$TB,  
probability=TRUE, col="steelblue3",  
border="white",  
ylim=c(0,0.08),  
breaks=30,  
main="Histogramme",  
xlab="Taux butyreux")
```

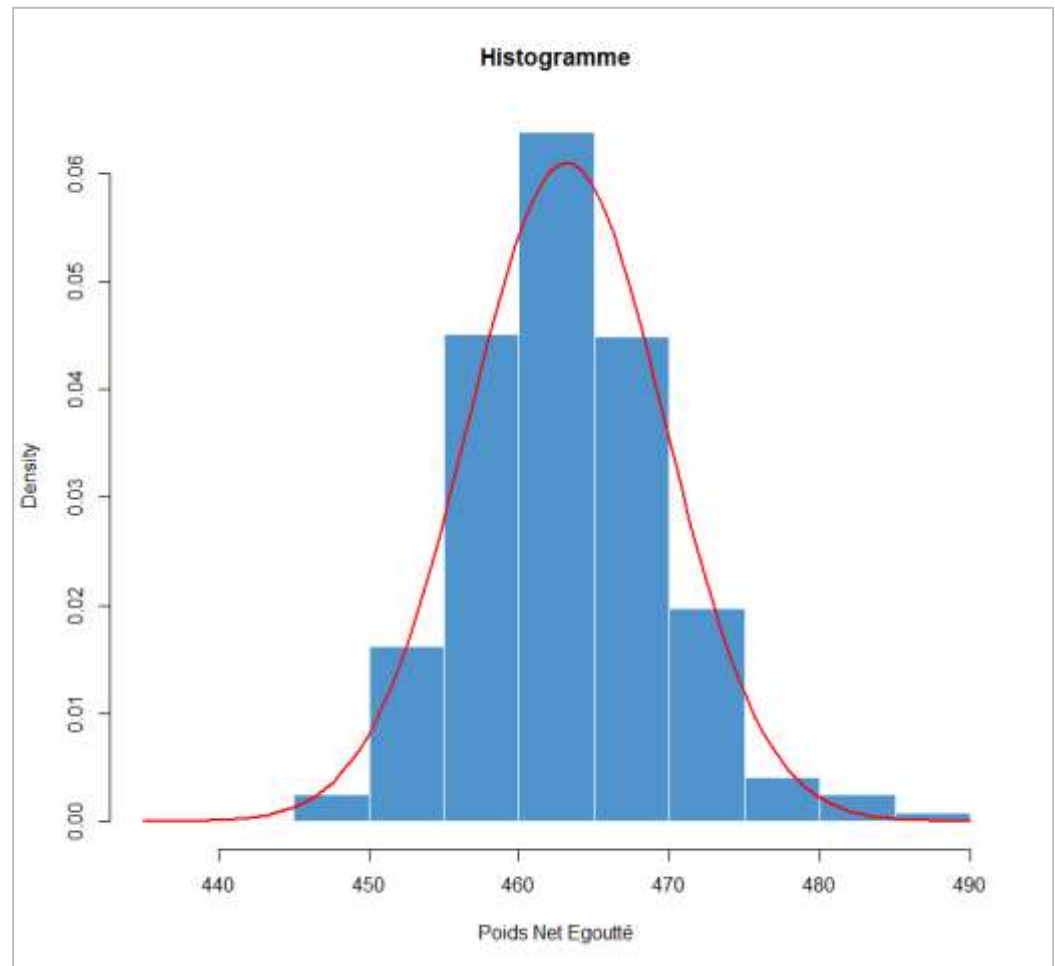
```
curve(dnorm(x, mean=mean(vache$TB),  
sd=sd(vache$TB)), add=T,  
col="red", lwd=2)
```



## 1. Notions de base – La loi normale

### Exemple 2.

Poids net égoutté  
de 2042 boîtes de conserve



## 1. Notions de base – La loi normale

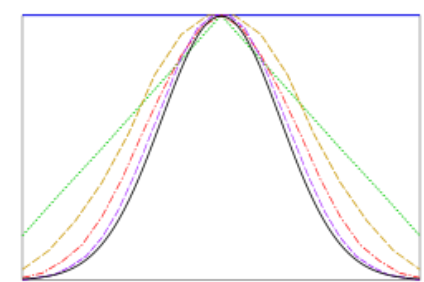
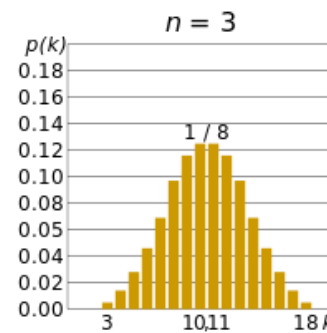
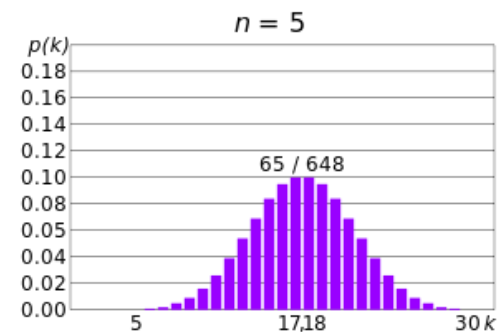
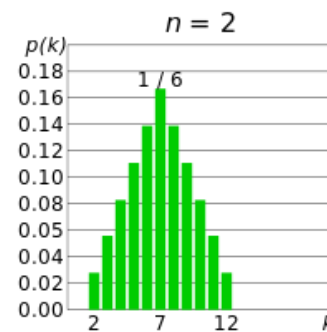
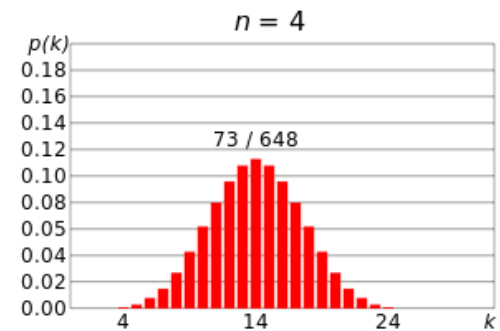
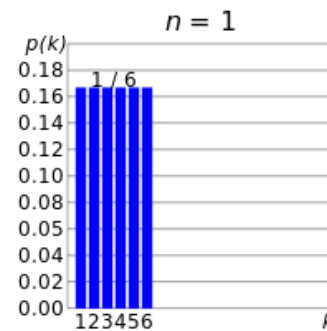
### Autres utilisations ou applications

- Balistique (portée et direction)
- Erreur de mesure en astronomie
- Modélisation du Quotient Intellectuel
- Taille humaine (pour une classe d'âge donnée)
- Courbe de croissance (carnets de santé)
- Un caractère mesurable dans une population peut être modélisé à l'aide d'une loi normale s'il est codé génétiquement par de nombreux allèles ou si le caractère dépend d'un grand nombre d'effets environnementaux
- Accroissement du prix d'une denrée en Bourse (log –  $N$ )

# 1. Notions de base – La loi normale

## La loi de LG comme loi limite

Loi de la somme de  $n$  dés  
(source : Wikipedia)

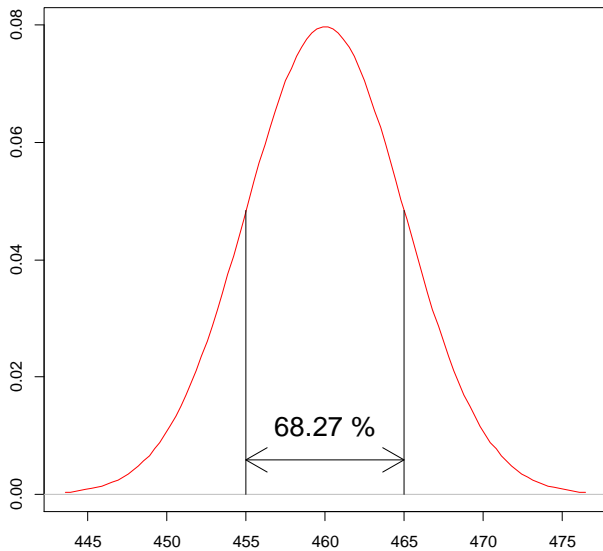




# 1. Notions de base – La loi normale

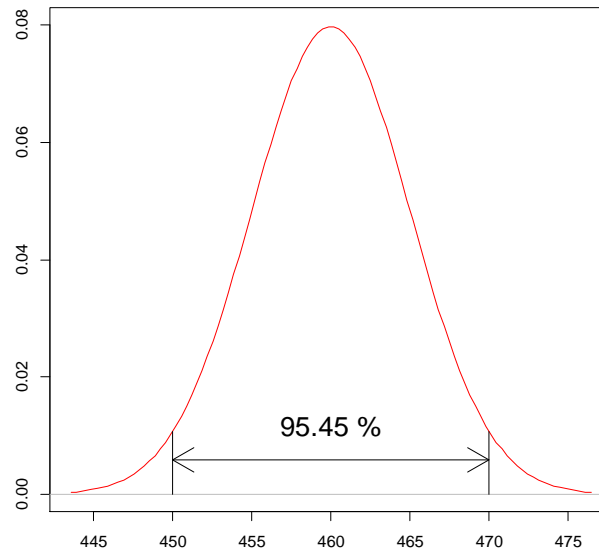
## Quelques valeurs remarquables

Normal Distribution: Mean=460, Standard deviation=5



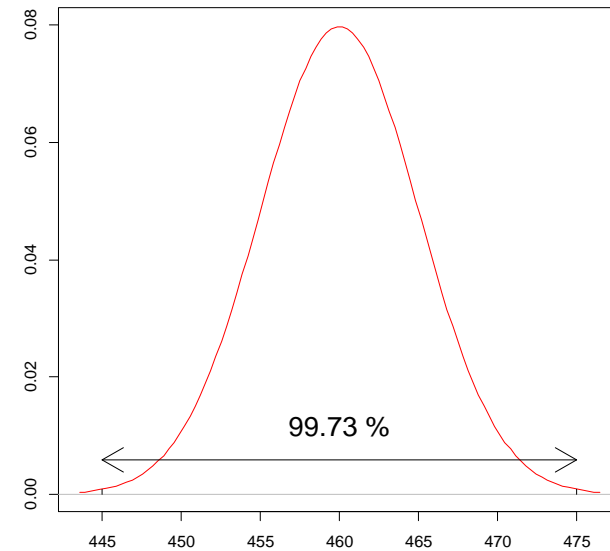
Moyenne +/- 1 ecarts types

Normal Distribution: Mean=460, Standard deviation=5



Moyenne +/- 2 ecarts types

Normal Distribution: Mean=460, Standard deviation=5

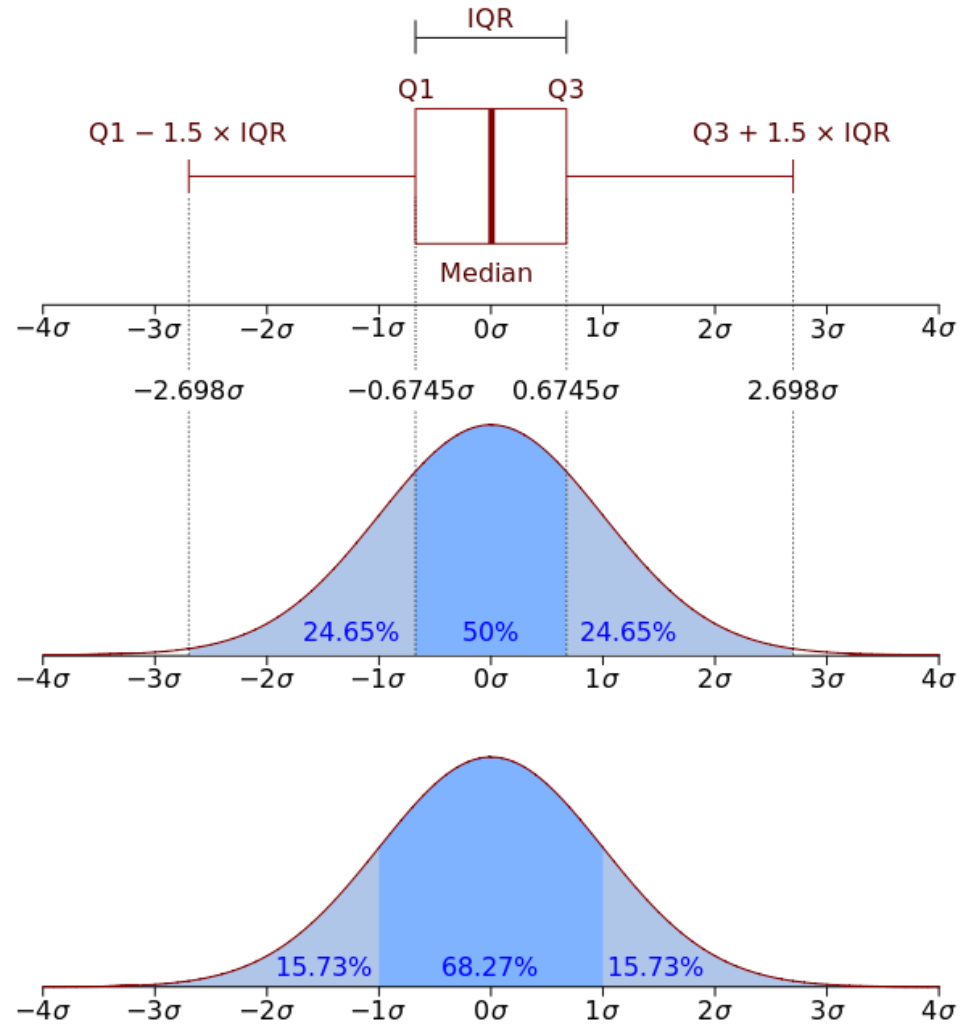


Moyenne +/- 3 ecarts types

```
s=3
plot(.x, dnorm(.x, mean=460, sd=5), xlab=paste("Moyenne +/- ",s,"ecarts types"), ylab="",
     main=paste("Normal Distribution: Mean=460, Standard deviation=5"),type="l", cex.lab=1.5, col="red")
abline(h=0, col="gray")
segments(460-s*5,0,460-s*5,dnorm(460-s*5, mean=460, sd=5))
segments(460+s*5,0,460+s*5,dnorm(460+s*5, mean=460, sd=5))
arrows(460-s*5, 0.006, 460+s*5, 0.006, code=3)
text(460,0.012,paste(round(100*(2*pnorm(c(s), mean=0, sd=1, lower.tail=TRUE)-1),2),"%"),cex=1.7)
```

# 1. Notions de base – La loi normale

## Boîte à moustache et loi normale

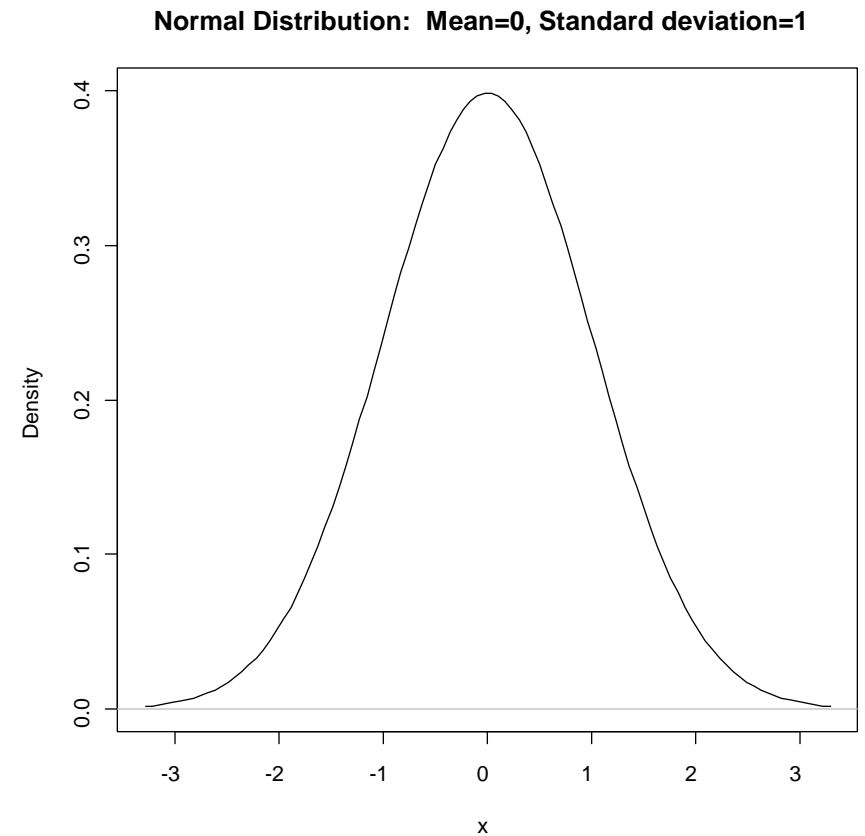


source : Wikipedia

## 1. Notions de base – La loi normale

## La loi normale centrée réduite

$$X \sim \mathcal{N}(\mu, \sigma) \Leftrightarrow Z = \frac{X - \mu}{\sigma} \sim \mathcal{N}(0, 1)$$



## 2. Notions d'échantillonnage

# 2. Notions d'échantillonnage

## Plan

- Etude de la moyenne d'échantillon
- Etude de la variance d'un échantillon
- Etude de la proportion d'un échantillon

## Objectif

Déterminer comment se comporte un indicateur statistique (moyenne, variance, proportion, etc.) au sein d'un échantillon aléatoire issu d'une population donnée

## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

### Contexte

- Soit une variable  $X$  étudiée au sein d'une population (supposée infiniment grande) de moyenne  $E(X) = \mu$  et de variance  $V(X) = \sigma^2$
- On prélève dans cette population un échantillon de  $n$  individus (ayant la même probabilité d'être tirés)
- Les valeurs de  $X$  pour ces  $n$  individus forment ce qu'on appelle un *échantillon aléatoire simple* (EAS) noté  $\{X_1, X_2, \dots, X_n\}$

### Terminologie

- Avant tirage, la valeur de  $X$  n'est pas connue pour l'individu  $i$ .  
 $X_i$  est donc une *variable aléatoire*
- L'échantillon est donc formé de  $n$  variables aléatoires i.i.d : indépendantes et identiquement distribuées (= de même loi que  $X$ )
- Après tirage, on dispose de  $n$  valeurs numériques ou *réalisations* de la variable  $X$  notées  $\{x_1, x_2, \dots, x_n\}$

## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

### 2.1 – Etude de la moyenne d'échantillon

- On calcule au sein d'un échantillon de taille  $n$ , la moyenne d'échantillon notée  $\bar{X}_n$  et définie par

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

- L'échantillon étant aléatoire, la valeur de  $\bar{X}_n$  est également aléatoire

À quelle valeur de  $\bar{X}$  peut-on s'attendre en moyenne ?

- La valeur moyenne attendue pour  $\bar{X}_n$  est calculée par l'espérance de  $\bar{X}_n$  (= moyenne pour un nombre « infiniment grand » d'échantillons)

$$E(\bar{X}_n) = \mu$$

## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

### Illustration

- Population :  $X$  prend les valeurs de 0 à 100.
- $E(X) = 50$
- Prélèvement de  $K$  échantillons aléatoires de taille  $n = 30$
- $K = 3$  échantillons

Moyenne des  $\bar{x}_{30} = 47,47$

- $K = 10$  échantillons

Moyenne des  $\bar{x}_{30} = 49,43$

- $K = 100$  échantillons

Moyenne des  $\bar{x}_{30} = 50,30$

E1	E2	E3	E4	E5	E6	E7	E8	E9	E10	
95	98	57	42	53	48	87	37	8	40	
34	47	76	18	98	51	31	74	88	23	
81	37	95	59	42	16	57	41	28	74	
54	53	71	95	35	39	35	48	24	24	
6	83	3	92	23	73	39	60	1	49	
39	90	74	45	38	53	37	31	87	78	
29	83	71	72	82	31	74	49	34	88	
65	75	70	8	32	72	16	70	66	79	
50	52	81	98	54	6	40	74	2	66	
39	30	88	66	14	10	53	83	3	28	
97	53	55	6	46	29	33	94	54	1	
38	18	1	93	17	11	82	22	29	80	
54	67	19	99	17	51	1	6	66	95	
18	88	93	68	86	38	12	99	86	22	
32	68	31	15	17	43	51	41	31	38	
41	7	33	72	48	48	77	93	18	64	
39	62	39	54	81	33	16	5	19	69	
93	75	84	38	6	44	74	1	34	86	
40	51	41	43	57	71	1	16	100	81	
86	82	31	68	94	100	15	94	10	61	
5	95	23	56	58	43	21	27	17	77	
63	97	2	80	87	3	34	45	37	60	
85	73	19	15	10	2	98	39	58	34	
59	10	31	16	38	39	91	53	49	82	
25	36	22	21	80	46	85	44	4	39	
63	80	68	51	18	69	62	91	66	70	
91	63	86	6	10	4	38	87	21	60	
31	10	65	58	82	12	84	12	64	44	
31	11	55	44	5	89	31	20	50	49	
70	77	55	12	5	93	37	38	57	79	
51,77	59,03	51,30	50,33	44,43	42,23	47,07	49,80	40,37	58,00	49,43

## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

### Quelle est la variabilité attendue pour $\bar{X}$ ?

- La variabilité des moyennes d'échantillons  $\bar{X}_n$  est calculée par la variance de  $\bar{X}_n$

$$V(\bar{X}_n) = \frac{\sigma^2}{n}$$

- *Elle est  $n$  fois plus petite que la variance de  $X$*
- *Elle est d'autant plus faible que  $n$  est élevé*

### Illustration

- Population : valeurs de 0 à 100. On tire 100 échantillons de taille 30
- $V(X) = \sigma^2 = \frac{101^2-1}{12} = 850$  .  $V(\bar{X}) = \frac{\sigma^2}{n} = \frac{850}{30} = 28,33$
- Variance des  $\bar{x}_{30}$  : 27,32



## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

Quelle est la forme de la distribution de  $\bar{X}$  ?

### Le cas d'un échantillon gaussien

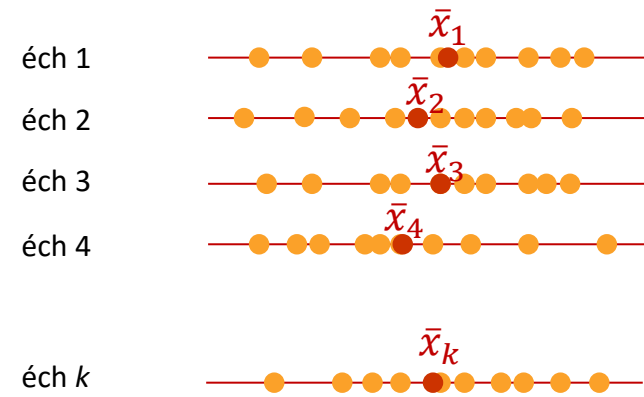
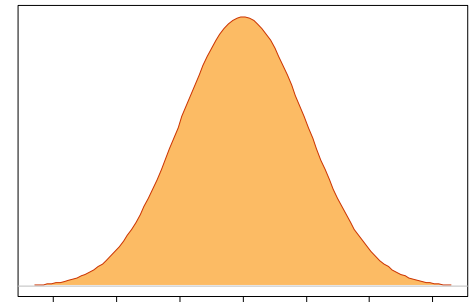
- On suppose que la variable étudiée est distribuée selon une loi Normale :

$$X \sim N(\mu, \sigma)$$

- On prélève des échantillons aléatoires de taille  $n$

$$E(\bar{X}_n) = \mu$$

$$V(\bar{X}_n) = \frac{\sigma^2}{n}$$



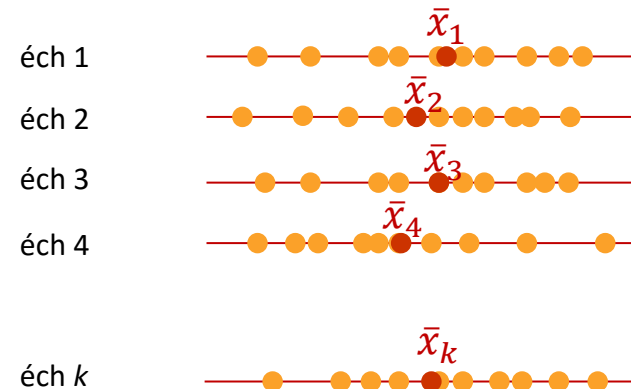
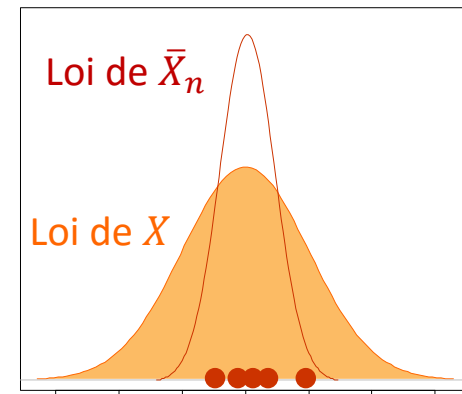
## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

### Le cas d'un échantillon gaussien

- La distribution de la moyenne est elle-même distribuée selon une loi normale

$$\bar{X}_n \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

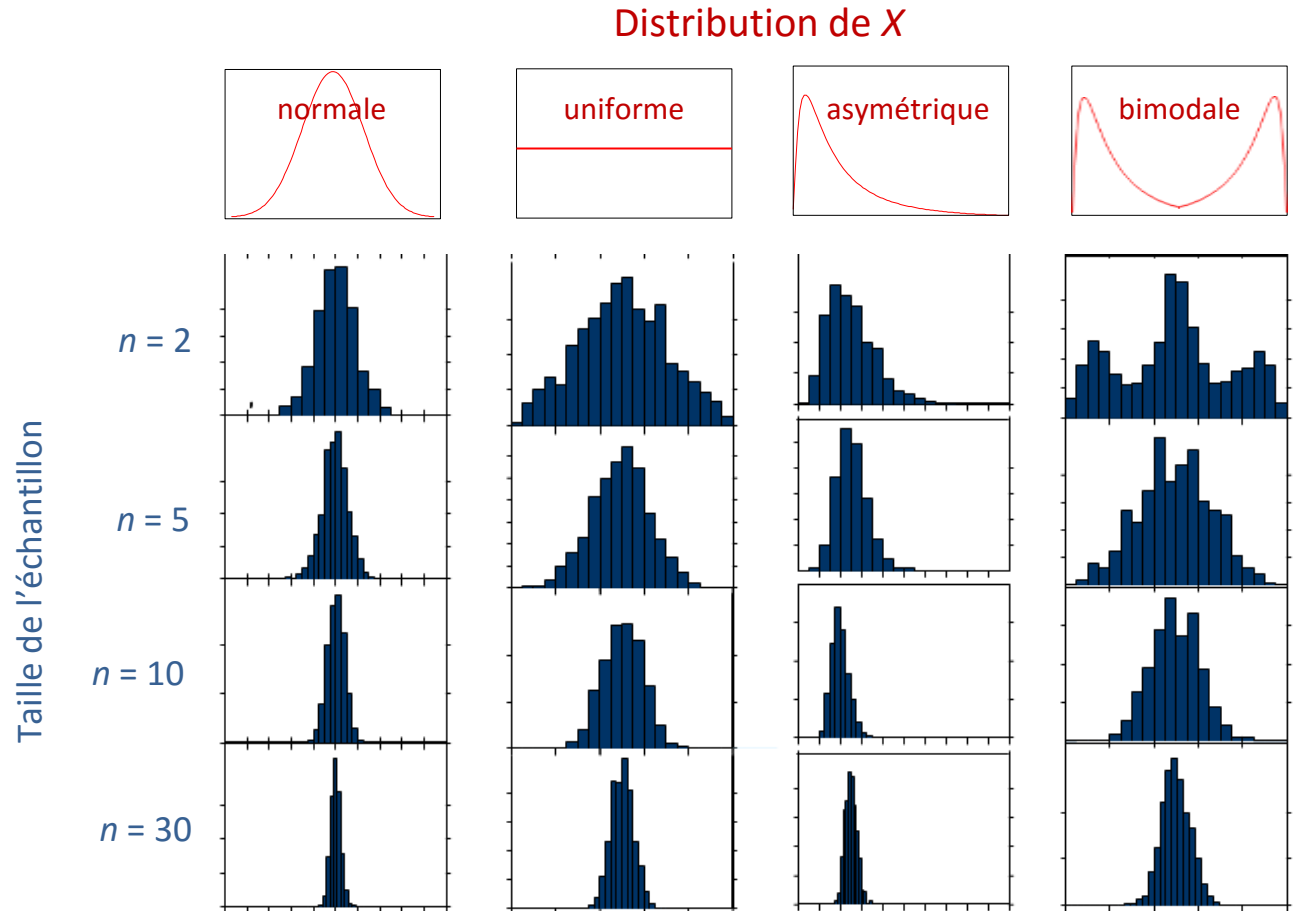
- Justification : toute combinaison linéaire de V.A. distribuées selon une loi normale est également distribuée normalement



## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

### Le cas d'une loi quelconque pour $X$

Distribution de  $\bar{X}_n$   
selon la loi de  $X$  et la  
taille  $n$  de l'échantillon



## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

### Le cas d'une loi quelconque pour $X$

#### Théorème Central Limite (TCL)

Une moyenne d'échantillon, centrée et réduite, converge en loi vers la loi normale  $N(0,1)$

$$\bar{X}_n \xrightarrow{\text{Loi}} N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

$$\frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \xrightarrow{\text{Loi}} N(0,1)$$

#### Interprétation

*Pour  $n$  « assez grand » (en pratique  $n \geq 30$ ) on peut admettre que la loi d'une moyenne d'échantillon est approximativement une loi normale  $N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$*

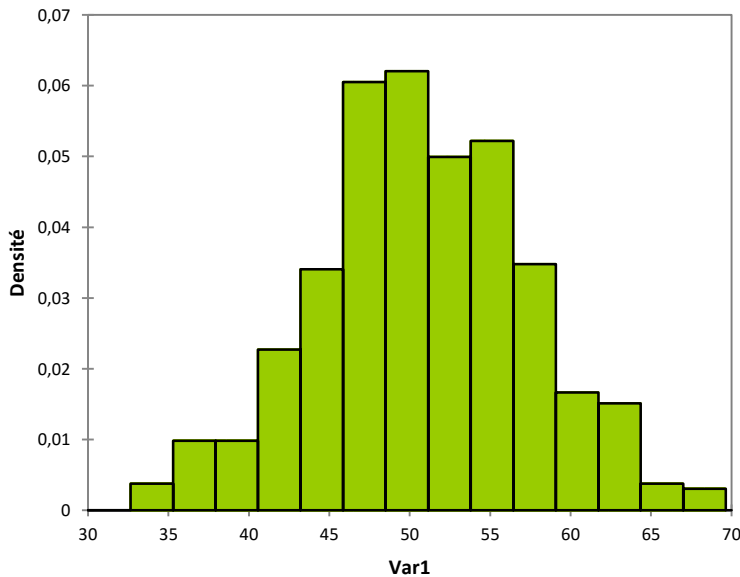
### Résultat capital en statistique !

## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

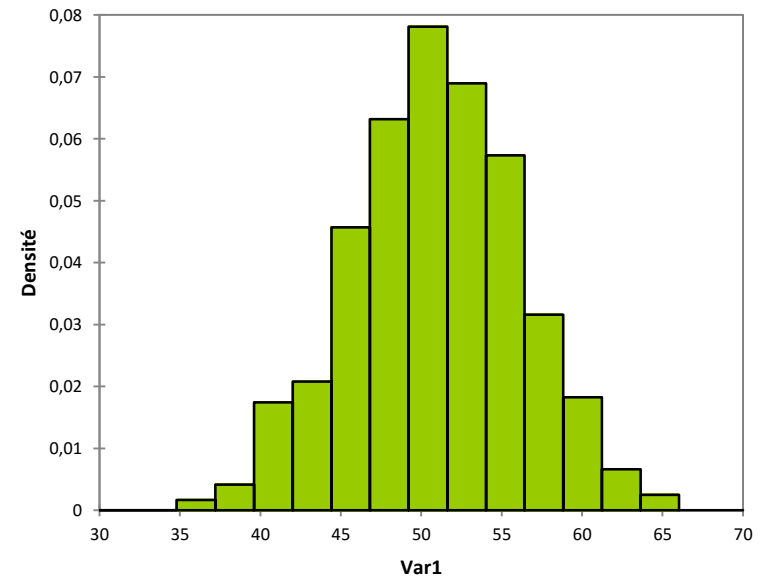
### Illustration

- Population :  $X$  prend les valeurs de 0 à 100.
- Prélèvement de  $K = 500$  échantillons aléatoires de taille  $n = 20$  ou  $n = 30$
- Tracé de l'histogramme des 500 moyennes

Histogramme de 500 moyennes d'échantillon de taille  $n = 20$



Histogramme de 500 moyennes d'échantillon de taille  $n = 30$



## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

### Le cas d'une population de taille **finie** - Tirage SANS remise

- Soit une **population de taille finie**  $N$
- On effectue des tirages **sans remise** d'échantillons de taille  $n$
- À quelle valeur moyenne et quelle variabilité peut-on s'attendre pour  $\bar{X}_n$  ?

$$E(\bar{X}_n) = \mu$$

$$V(\bar{X}_n) = \left( \frac{N - n}{N - 1} \right) \times \frac{\sigma^2}{n}$$

## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

### Le cas d'une population de taille **finie** - Tirage **AVEC** remise

- Soit une **population de taille finie**  $N$
- On effectue des tirages **avec remise** d'échantillons de taille  $n$

Une population finie dans laquelle on effectue un échantillonnage avec remise se comporte comme une population infinie !

$$E(\bar{X}_n) = \mu$$

$$V(\bar{X}_n) = \frac{\sigma^2}{n}$$

**Remarque** : lorsque la population est de **taille infinie**, ou de manière générale lorsque  $n$  est petit devant  $N$ , on ne distingue plus tirage avec ou sans remise

## 2. Notions d'échantillonnage – Loi d'une variance d'échantillon

### 2.1 – Distribution d'une variance d'échantillon

- On calcule au sein d'un échantillon de taille  $n$ , la variance d'échantillon notée  $S^2$  et définie par

$$S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

- L'échantillon étant aléatoire, la valeur de  $S^2$  est aussi aléatoire

À quelle valeur de  $S^2$  peut-on s'attendre en moyenne ?

- La valeur moyenne attendue pour  $S^2$  est calculée par l'espérance de  $S^2$  :

$$E(S^2) = \frac{n-1}{n} \sigma^2$$



## 2. Notions d'échantillonnage – Loi d'une variance d'échantillon

### À quelle variabilité de $S^2$ peut-on s'attendre ?

- La variabilité attendue pour  $S^2$  est calculée par sa variance :

$$V(S^2) = \frac{n-1}{n^3} [(n-1)\mu_4 - (n-3)\sigma^4]$$

$$V(S^2) \xrightarrow{n \rightarrow \infty} \frac{\mu_4 - \sigma^4}{n}$$

Rappel : moment centré d'ordre  $k$  de  $X$  :  $\mu_k = E[(X - \mu)^k]$

### Quelle est la forme de la distribution de $S^2$ ?

- La statistique  $S^2$  centrée réduite converge en loi vers une loi normale  $N(0,1)$
- Pour des échantillons gaussiens, on a :**

$$\frac{nS^2}{\sigma^2} \sim \chi^2_{n-1}$$

## 2. Notions d'échantillonnage – Loi d'une variance d'échantillon

### Le cas de tirages SANS remise

- Soit une **population de taille finie  $N$**
- On effectue des tirages **sans remise** d'échantillons de taille  $n$
- À quelle valeur moyenne peut-on s'attendre pour  $S^2$  ?

$$E(S^2) = \frac{n-1}{n} \frac{N}{N-1} \sigma^2$$

## 2. Notions d'échantillonnage – Loi d'une moyenne d'échantillon

### Exemple

- Soit une population de taille finie  $N = 4 : \{-4, -2, +2, +4\}$
- Calculer moyenne et variance de la population
- Déterminer tous les échantillons possibles de taille  $n = 2$ , sans remise
- Calculer l'espérance et la variance de la moyenne d'échantillon  $\bar{X}_n$
- Calculer l'espérance de la variance  $S^2$

## 2. Notions d'échantillonnage – Loi d'une fréquence d'échantillon

### 2.1 – Distribution d'une fréquence d'échantillon

#### Contexte

- Dans une population, une proportion  $\pi$  d'individus possède un caractère  $A$
- Un échantillon de taille  $n$  est prélevé aléatoirement
- On note  $K$  le nombre d'individus possédant  $A$  dans l'échantillon :  $K \sim B(n, \pi)$
- Soit  $F$ , la fréquence empirique associée :  $F = \frac{K}{n}$

#### Valeur moyenne et variabilité attendues pour $F$ ?

$$E(F) = \pi \qquad V(F) = \frac{\pi(1 - \pi)}{n}$$

## 2. Notions d'échantillonnage – Loi d'une fréquence d'échantillon

### Théorème de De Moivre - Laplace

Lorsque l'échantillon est de grande taille ( $n \geq 30$ )

$$K \approx N \left( n\pi, \sqrt{n\pi(1-\pi)} \right) \quad \text{avec} \quad n\pi \geq 15 \quad \text{et} \quad n\pi(1-\pi) \geq 5$$

Autres conditions « similaires »

$$\begin{cases} n\pi \geq 5 \quad \text{et} \quad n(1-\pi) \geq 5 \\ n\pi(1-\pi) \geq 9 \\ n \geq 30 \quad \text{et} \quad 0,3 \leq \pi \leq 0,7 \end{cases}$$

À quelle distribution peut-on s'attendre pour  $F$  ?

$$\text{Dans les mêmes conditions :} \quad F \approx N \left( \pi, \sqrt{\frac{\pi(1-\pi)}{n}} \right)$$

## 2. Notions d'échantillonnage – Loi d'une fréquence d'échantillon

### Exemple

- On jette  $n = 100$  fois une pièce équilibrée
- On s'intéresse à la fréquence de « face » obtenus
- Quelle est la probabilité que cette fréquence soit comprise entre 0,40 et 0,60 ?

## 3. L'estimation

### Plan

- Introduction
- Modèle, hypothèses, définitions
- Qualités d'un estimateur
- Construction d'estimateurs convergent ( $MM$  et  $MMV$ )
- Exhaustivité, Information de Fisher
- Estimation sans biais de variance minimale
- Normalité asymptotique

## 3. Estimation – Introduction

### 3.1 – Introduction

#### Le contexte : l'inférence statistique

- La population étudiée ne peut être étudiée dans sa globalité : trop volumineuse, inaccessible, potentiellement infinie
- On prélève dans la population un **échantillon aléatoire** (représentatif...)
- On souhaite **généraliser les résultats** obtenus dans l'échantillon au niveau de la population étudiée
- Type de problèmes abordés par l'inférence : **estimation** de paramètres, **test** d'hypothèse, construction de **modèle**



### 3. Estimation – L'estimation

#### Le problème de l'estimation

Il s'agit d'estimer un ou des **paramètres inconnus** (moyenne, variance, etc.) d'une population étudiée à partir d'un échantillon aléatoire issu de cette population

#### **Exemples.** Estimer :

- *Une intention de vote,*
- *La prévalence d'une maladie,*
- *Une ressource maritime,*
- *La précision d'une machine,*
- *L'âge moyen d'une mère à son premier enfant, etc.*

### 3. Estimation – Hypothèses, Définitions, Modèle

## 3.2 – Hypothèses, définitions, modèle statistique

### Population, échantillon, paramètre étudié

- Soit  $U$ , la population étudiée et  $\theta$  un paramètre réel inconnu au sein de cette population (ex : une moyenne)
- On s'intéresse dans cette population à une variable aléatoire  $X$  dont la loi dépend de  $\theta$
- Un échantillon aléatoire de taille  $n$  est prélevé dans la population  $U$ . Cet **échantillon aléatoire** est formé de  $n$  variables aléatoires indépendantes et identiquement distribuées (i.i.d.) :  $\{X_1, X_2, \dots, X_n\}$  dont les valeurs sont notées  $\{x_1, x_2, \dots, x_n\}$

### 3. Estimation – Hypothèses, Définitions, Modèle

#### Modèle statistique

- Les  $n$  variables aléatoires  $\{X_1, X_2, \dots, X_n\}$  suivent une même loi de probabilité  $P_\theta$  dépendant du paramètre  $\theta$  (ex : loi  $N$ , Bin, etc.)  
*On parle de modèle paramétrique*
- Le modèle d'échantillonnage associé au modèle probabiliste  $P_\theta$  constitue ce que l'on appelle un **modèle statistique**

#### But de l'estimation statistique

- Estimer le paramètre  $\theta$  (ou la loi  $P_\theta$ ) ou à partir d'un échantillon de  $n$  observations issues de la population étudiée
- Pour estimer  $\theta$ , on construira ce qu'on appelle un **estimateur** de  $\theta$ , défini à partir des  $n$  variables aléatoires  $\{X_1, X_2, \dots, X_n\}$

### 3. Estimation – Hypothèses, Définitions, Modèle

#### Estimateur

- Soit  $\theta$  le paramètre inconnu (à estimer) d'une population
- Un échantillon aléatoire  $\{X_1, X_2, \dots, X_n\}$  a été prélevé dans cette population
- On appelle estimateur du paramètre  $\theta$  toute fonction  $h$  des  $n$  observations  $\{X_i\}$ , notée  $T_n$  :

$$T_n = h(X_1, X_2, \dots, X_n)$$

**Remarque** : un estimateur est une **variable aléatoire**, car il dépend lui-même de l'échantillon aléatoire  $\{X_1, X_2, \dots, X_n\}$  prélevé

#### Estimation

Une fois l'échantillon observé, on dispose de  $n$  **observations**  $\{x_1, x_2, \dots, x_n\}$  qui nous fournissent une valeur  $h(x_1, x_2, \dots, x_n)$  appelée **estimation** de  $\theta$

### 3. Estimation – Hypothèses, Définitions, Modèle

#### Exemples d'estimateur

- La moyenne  $\bar{X}$  constitue un estimateur « naturel » de la moyenne  $\mu$  d'une population
- La fréquence empirique  $F$  d'un événement constitue un estimateur « naturel » de sa probabilité  $\Pi$  dans la population

#### Questions posées

- Quelles sont les **qualités** d'un « bon » estimateur ?
- Comment **choisir** parmi plusieurs estimateurs ?
- Comment **rechercher** un « bon » estimateur ?

### 3. Estimation – Hypothèses, Définitions, Modèle

#### Exemple

On s'intéresse à la taille  $X$  de personnes adultes (femmes).

On suppose que dans la population :  $N(\mu = 163, \sigma = 10)$

Mais ils sont inconnus et on souhaite les estimer

Un échantillon de taille  $n = 30$  est prélevé au hasard dans la population

Données de l'échantillon (données obtenues par simulation sous R) :

```
[1] 186.47 170.74 178.41 156.42 165.74 177.50 166.02 168.71 149.52 169.82
```

```
[11] 167.41 164.50 162.62 148.24 153.93 156.47 144.86 155.12 171.31 145.03
```

```
[21] 167.47 160.31 162.07 149.82 145.15 173.11 163.42 166.45 167.15 171.83
```

Moyenne de l'échantillon : 162,85

À partir de l'exemple, définir les différentes notions présentées :

*Population, individu, échantillon, variable, paramètres à estimer, modèle probabiliste, modèle statistique, estimateur, estimation*

### 3. Estimation – Qualités d'un estimateur

## 3.3 – Qualités d'un estimateur

### Convergence

Première qualité attendue : lorsque la taille de l'échantillon tend vers  $\infty$ , l'estimateur tend vers  $\theta$  :

$$P(|T_n - \theta| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} 0$$

On dit que l'estimateur  $T_n$  *converge en probabilité* vers  $\theta$

Un estimateur **convergent** est également appelé estimateur **consistant**

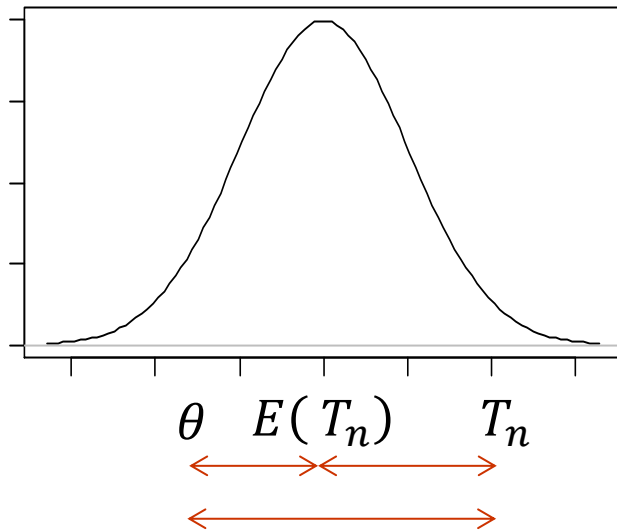
### 3. Estimation – Qualités d'un estimateur

#### Biais d'un estimateur

Lorsqu'on estime le paramètre  $\theta$  à l'aide d'un estimateur  $T_n$  on commet l'**erreur d'estimation** définie par

$$T_n - \theta$$

*Illustration*



Cet écart regroupe deux types d'erreurs :

- **Erreur systématique** : due aux fluctuations aléatoires de  $T_n$  autour de sa moyenne

$$T_n - E(T_n)$$

- **Biais** : écart entre la valeur moyenne de  $T_n$  et la « cible » recherchée

$$\text{Biais}(T_n) = E(T_n) - \theta$$



### 3. Estimation – Qualités d'un estimateur

#### Estimateur sans biais

$T_n$  est un estimateur *sans biais* ou *non biaisé* du paramètre  $\theta$  ssi

$$E(T_n) = \theta \Leftrightarrow \text{Biais}(T_n) = 0$$

$T_n$  est un estimateur *asymptotiquement sans biais* du paramètre  $\theta$  ssi

$$E(T_n) \xrightarrow{n \rightarrow \infty} \theta \Leftrightarrow \text{Biais}(T_n) \xrightarrow{n \rightarrow \infty} 0$$

#### Applications :

- trouver un estimateur sans biais pour  $\mu$
- trouver un estimateur sans biais pour  $\sigma^2$
- trouver un estimateur sans biais pour  $\pi$

### 3. Estimation – Qualités d'un estimateur

#### Bilan

$\bar{X}$  est un estimateur *sans biais* de la moyenne  $\mu$  d'une population

$S_c^2$  est un estimateur *sans biais* de la variance  $\sigma^2$  d'une population

$F$  est un estimateur *sans biais* de la proportion  $\pi$  d'une population

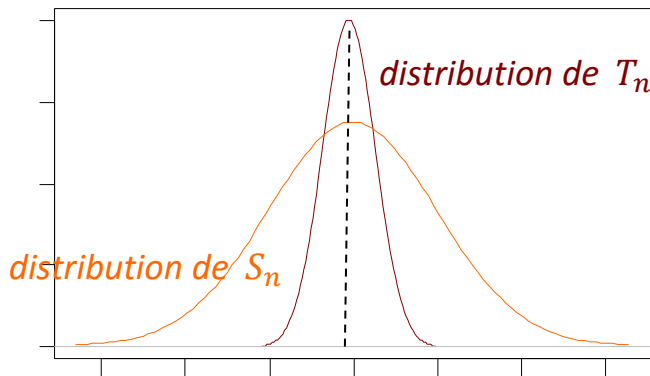
### 3. Estimation – Qualités d'un estimateur

#### Précision d'un estimateur

Dans quelle mesure  $T_n$  permet-il d'estimer avec précision le paramètre  $\theta$  ?

L'estimateur  $T_n$  fournit-il des résultats très variables d'un échantillon aléatoire à l'autre ?

*Illustration : comparaison de la distribution de deux estimateurs*



$$E(T_n) = E(S_n) = \theta$$

- Les deux estimateurs sont **non biaisés**
- L'estimateur  $T_n$  fluctue moins que  $S_n$  autour de la vraie valeur du paramètre
- L'estimateur  $T_n$  est **plus précis** que  $S_n$

### 3. Estimation – Qualités d'un estimateur

#### Précision d'un estimateur

Si  $T_n$  est un estimateur **non biaisé** du paramètre  $\theta$  alors la **précision** de cet estimateur est déterminée par sa **variance**  $V(T_n)$

Face à deux estimateurs non biaisés, il convient donc de choisir celui qui possède **la plus petite variance**

**Exemple.** Comparaison de deux estimateurs de la moyenne  $\mu$  dans le cas d'une loi  $N(\mu, \sigma)$

- La moyenne d'échantillon
- La médiane

### 3. Estimation – Qualités d'un estimateur

#### Précision d'un estimateur : l'erreur quadratique moyenne (EQM)

Si  $T_n$  est un estimateur quelconque du paramètre  $\theta$  alors la **précision** de cet estimateur est déterminée par son **erreur quadratique moyenne** (ou fonction de **risque quadratique**)

$$\text{EQM}(T_n) = E[(T_n - \theta)^2]$$

On peut montrer que l'EQM est égale à :

$$\text{EQM}(T_n) = V(T_n) + [E(T_n) - \theta]^2 = V(T_n) + \text{Biais}(T_n)^2$$

Entre deux estimateurs quelconques du paramètre  $\theta$  on choisira celui dont l'**erreur quadratique moyenne** est la plus faible

### 3. Estimation – Qualités d'un estimateur

#### Exemple 1

*Comparaison de la précision de deux estimateurs de la variance  $\sigma^2$  lorsque la moyenne  $\mu$  est connue*

- La variance d'échantillon corrigée  $S_c^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$
- La variance  $\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$

#### Exemple 2

*Comparaison de la précision de deux estimateurs du paramètre  $\lambda$  d'une loi de Poisson*

### 3. Estimation – Qualités d'un estimateur

#### Propriété importante – Convergence d'un estimateur

Si un estimateur est **non biaisé** et que sa **variance tend vers 0** quand la taille  $n$  de l'échantillon tend vers l'infini,  
Alors cet estimateur est **convergent** (= consistant)

**Exemple.** Convergence de la moyenne d'échantillon

### 3. Estimation – Qualités d'un estimateur

#### Convergence en moyenne quadratique

L'estimateur  $T_n$  est convergent en moyenne quadratique si

$$\lim_{n \rightarrow \infty} \text{EQM}(T_n) = 0$$

L'estimateur  $T_n$  est convergent en moyenne quadratique si les deux propriétés suivantes sont vérifiées :

- $\lim_{n \rightarrow \infty} V(T_n) = 0$
- $T_n$  est asymptotiquement sans biais



### 3. Estimation – Méthode des moments

## 3.4 – Construction d'estimateurs convergents

### La méthode des moments (MM)

Supposons que nous devons estimer  $K$  paramètres  $\theta = (\theta_1, \theta_k, \dots, \theta_K)$

Les estimateurs des moments de  $\theta$ , notés  $\tilde{\theta}$ , sont solutions des équations :

$$E_{\tilde{\theta}}(X^k) = \mu_k = m_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

Moment **théorique**

Moment **empirique**

La méthode consiste donc à égaliser les  $K$  premiers moments *non centrés* théoriques (inconnus) aux moments empiriques correspondants (connus) de l'échantillon

### 3. Estimation – Méthode des moments

#### Remarques

- Méthode intuitive, conceptuellement plus simple que celle du MV (MV = maximum de vraisemblance)
- Chaque  $\tilde{\theta}_i$  est un **estimateur convergent** du paramètre  $\theta_i$  correspondant

#### Exemples

- Estimer par la MM le paramètre  $\lambda$  d'une loi de Poisson,
- Estimer par la MM le paramètre  $\theta$  d'une loi Exponentielle
- Déterminer les estimateurs des moments des paramètres  $a$  et  $b$  d'une loi Uniforme  $U_{[a;b]}$
- Déterminer les estimateurs des moments des paramètres  $\alpha$  et  $\beta$  d'une loi Beta

### 3. Estimation – Méthode des moments

#### Illustration pour la loi Uniforme

- Soit une loi Uniforme  $U_{[2;7]}$
- Les valeurs 2 et 7 des paramètres  $a$  et  $b$  sont supposées inconnues  
On souhaite les estimer
- Un échantillon de taille  $n = 100$  issu aléatoirement de cette loi  
3.769638 6.342496 3.460563 5.484658 4.860614 2.209637 5.638404 ...
- Dans l'échantillon on trouve :  $\bar{x} = 4.25$  et  $s = 1.46$
- Proposer une estimation des paramètres  $a$  et  $b$  par la méthode des moments

### 3. Estimation – Méthode du maximum de vraisemblance

## La méthode du maximum de vraisemblance

### Vraisemblance d'un échantillon

- Soit  $X$  une variable aléatoire dont la densité de probabilité  $f_X(x; \theta)$  dépend d'un paramètre  $\theta$
- Soit  $\{x_1, x_2, \dots, x_n\}$  les réalisations indépendantes de  $X$  issues d'un échantillon aléatoire de taille  $n$

On appelle **vraisemblance de l'échantillon**, la fonction du paramètre  $\theta$ , notée  $L$ , définie par :

$$L(x_1, x_2, \dots, x_n; \theta) = f_X(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f_X(x_i; \theta)$$

### 3. Estimation – Méthode du maximum de vraisemblance

#### Remarque

Si  $X$  est une variable aléatoire **discrète** de loi  $P_X(x; \theta)$

alors la vraisemblance est définie par

$$L(x_1, x_2, \dots, x_n; \theta) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n; \theta) = \prod_{i=1}^n P_X(x_i; \theta)$$

#### Exemples de calcul de vraisemblance

- $X \sim P(\lambda)$
- $X \sim N(\mu, \sigma)$
- $X \sim B(\theta)$

### 3. Estimation – Méthode du maximum de vraisemblance

#### Estimateur du maximum de vraisemblance

On dit que  $\hat{T}_n$  est l'estimateur du maximum de vraisemblance du paramètre  $\theta$  s'il maximise la fonction de vraisemblance  $L(x_1, x_2, \dots, x_n; \theta)$  :

$$\hat{T}_n = \operatorname{argmax}_{\theta} L(x_1, x_2, \dots, x_n; \theta)$$

#### Détermination pratique d'un estimateur du MV

Très souvent, l'expression de la densité conduit à rechercher plus facilement  $\hat{T}_n$  de façon à ce qu'il **maximise la log-vraisemblance**

Ainsi, l'estimateur du maximum de vraisemblance est solution du problème :

$$\begin{cases} \frac{\partial \ln L(\mathbf{x}; \theta)}{\partial \theta} = 0 \\ \frac{\partial^2 \ln L(\mathbf{x}; \theta)}{\partial \theta^2} < 0 \end{cases}$$

### 3. Estimation – Méthode du maximum de vraisemblance

#### Exemple 1

- Estimateur du maximum de vraisemblance de  $\theta = (\mu, \sigma)$  pour  $X \sim N(\mu, \sigma)$  ?
- Estimateur du maximum de vraisemblance de  $\theta = \lambda$  pour  $X \sim P(\lambda)$  ?
- Estimateur du maximum de vraisemblance de  $\theta = (a, b)$  pour  $X \sim U_{[a;b]}$  ?
- Estimateur du maximum de vraisemblance de  $\theta$  pour  $X \sim B(\theta)$  ?

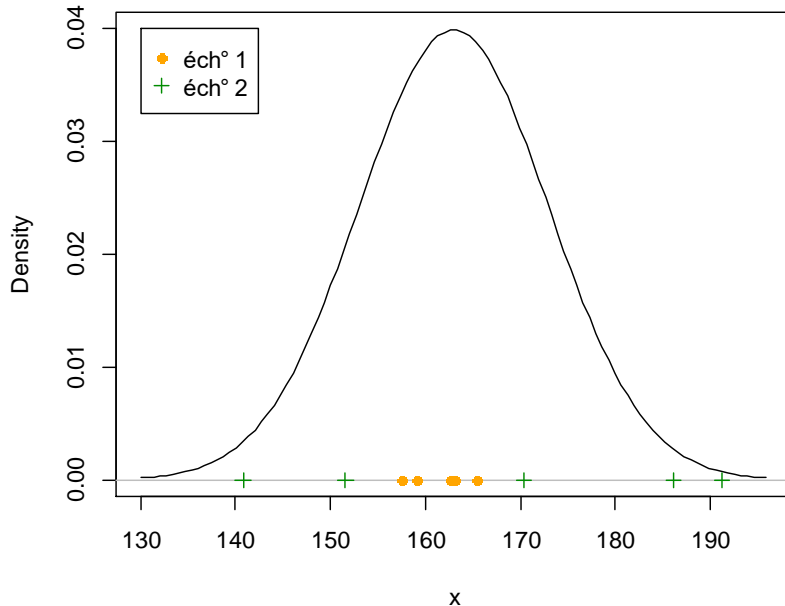
### 3. Estimation – Méthode du maximum de vraisemblance

#### Exemple 2

Soit  $X \sim N(\mu = 163, \sigma = 10)$

- Quelle est la vraisemblance
  - de l'échantillon  $E1 = \{157.5, 159.1, 163.2, 162.7, 165.5\}$  pour la loi  $N(\mu = 163, \sigma = 10)$  ?
  - de l'échantillon  $E2 = \{140.8, 151.6, 170.4, 186.2, 191.3\}$  pour la loi  $N(\mu = 163, \sigma = 10)$  ?
- Quel est l'estimateur du MV de la moyenne  $\mu$  pour  $E1$  ? Pour  $E2$  ?

Normal Distribution: Mean=163, Standard deviation=10



```
.x <- seq(130.095, 195.905, length.out=100)
plot(.x, dnorm(.x, mean=163, sd=10), xlab="x",
     ylab="Density", main=paste("Normal Distribution:
     Mean=163, Standard deviation=10"), type="l")
abline(h=0, col="gray")
points(E1, rep(0, 5), pch=16, col="orange")
points(E2, rep(0, 5), pch=3, col="green4")
legend(130, 0.04, c("éch° 1", "éch° 2"),
      col=c("orange", "green4"), pch=c(16, 3))
```



### 3. Estimation – Méthode du maximum de vraisemblance

```
E1=c(157.5, 159.1, 163.2, 162.7, 165.5)
```

```
L=1
```

```
for (i in 1:length(E1)) {
  print(dnorm(E1[i], 163, 10))
  L=L*dnorm(E1[i], 163, 10)
}
```

```
[1] 0.03429439
[1] 0.03697277
[1] 0.03988625
[1] 0.03987628
[1] 0.03866681
```

densité en  
chaque point de  
l'échantillon

```
print(L)
[1] 7.797964e-08
```

vraisemblable

```
print(log(L))
[1] -16.36682
```

log-vraisemblable

```
E2=c(140.8, 151.6, 170.4, 186.2, 191.3)
```

```
L=1
```

```
for (i in 1:length(E2)) {
  print(dnorm(E2[i], 163, 10))
  L=L*dnorm(E2[i], 163, 10)
}
```

```
[1] 0.003394076
[1] 0.02083078
[1] 0.03033893
[1] 0.00270481
[1] 0.0007274439
```

```
print(L)
[1] 4.220497e-12
```

```
print(log(L))
[1] -26.19107
```

- L'échantillon E1 est « plus vraisemblable » sous la loi  $N(\mu = 163, \sigma = 10)$  que E2 Sa « probabilité » d'être tiré aléatoirement sous  $N(\mu = 163, \sigma = 10)$  est plus élevée

### 3. Estimation – Exhaustivité, information de Fisher

## 3.5 – Exhaustivité, information de Fisher

### Qu'est-ce qu'une statistique exhaustive ?

- Quand on cherche à estimer un paramètre  $\theta$ , l'échantillon  $\{X_1, X_2, \dots, X_n\}$  apporte une certaine quantité d'information sur celui-ci
- En résumant l'information apportée par l'échantillon au travers d'une statistique  $T_n = h(X_1, X_2, \dots, X_n)$ , on souhaite **ne pas perdre cette information**
- Une statistique qui permet de conserver l'information sur le paramètre sera qualifiée de **statistique exhaustive**

### 3. Estimation – Exhaustivité, information de Fisher

#### Définition – Statistique exhaustive

- Soit  $L(x_1, x_2, \dots, x_n; \theta) = L(\mathbf{x}; \theta)$  la densité de l'échantillon  $\{X_1, X_2, \dots, X_n\}$  (également appelée vraisemblance du paramètre  $\theta$ )

$$L(\mathbf{x}; \theta) = \prod_{i=1}^n P_X(X = x_i; \theta) \quad \text{ou} \quad \prod_{i=1}^n f_X(x_i; \theta)$$

$X$  discrète

$X$  continue

- La statistique  $T_n$  est dite exhaustive si la loi conditionnelle de  $X$  sachant  $(T_n = t)$  est indépendante du paramètre  $\theta$

$$L(\mathbf{x}|T_n = t; \theta) = L(\mathbf{x}|T_n = t)$$

**Intuitivement** : une fois connue la statistique  $T_n$ , la connaissance de la totalité des observations  $(x_1, x_2, \dots, x_n)$  n'apporte aucune information supplémentaire sur  $\theta$ ,  $T_n$  apporte toute l'information possible sur  $\theta$  (nous n'avons plus besoin de  $\theta$  dans le calcul de la loi des  $x$ )

### 3. Estimation – Exhaustivité, information de Fisher

#### Principe de factorisation

##### *Une condition nécessaire et suffisante d'exhaustivité*

- Une statistique  $T_n$  est exhaustive si et seulement si la densité de l'échantillon  $\{X_1, X_2, \dots, X_n\}$  peut s'écrire sous la forme suivante :

$$L(\mathbf{x}; \theta) = g(T_n, \theta) h(\mathbf{x})$$

où  $g(T_n, \theta)$  désigne la loi de probabilité de  $T_n$ , dépendant du paramètre  $\theta$

- En d'autres termes, si et seulement si le rapport

$$\frac{L(\mathbf{x}; \theta)}{g(T_n, \theta)}$$

ne dépend plus du paramètre  $\theta$

### 3. Estimation – Exhaustivité, information de Fisher

#### Exemples

- $X \sim P(\lambda)$ . Statistique exhaustive pour le paramètre  $\theta = \lambda$  ?
- $X \sim N(\mu, \sigma)$ . Statistique exhaustive pour le paramètre  $\theta = \sigma^2$  ?
- $X \sim N(\mu, \sigma)$ . Statistique exhaustive pour le paramètre  $\theta = \mu$  ?

#### Limite du principe de factorisation

- Le principe de factorisation = moyen pour reconnaître si une statistique est exhaustive
- Mais ce principe ne nous permet pas en général de construire une statistique exhaustive ni même de savoir s'il en existe une...

### 3. Estimation – Exhaustivité, information de Fisher

## Loi permettant une statistique exhaustive

### Théorème de Darmais

Soit une variable aléatoire  $X$  dont le domaine de définition ne dépend pas du paramètre  $\theta$

Une condition nécessaire est suffisante pour que l'échantillon  $\{X_1, X_2, \dots, X_n\}$  admette une statistique exhaustive est que la forme de la densité soit

$$f(\mathbf{x}; \theta) = \exp(a(\mathbf{x})\alpha(\theta) + b(\mathbf{x}) + \beta(\theta))$$

On dit que la densité est de **forme exponentielle**

Si la densité est de cette forme, alors la statistique est une statistique exhaustive particulière

$$T_n = \sum_{i=1}^n a(X_i)$$

### 3. Estimation – Exhaustivité, information de Fisher

#### Remarques

- La plupart des densités habituelles (*Bernoulli, Poisson, Gauss, gamma, exponentielle* par exemple) appartiennent à la famille exponentielle
- Si le domaine de  $X$  (où la densité est différente de 0) **dépend du paramètre  $\theta$**  (par exemple la distribution uniforme), le théorème de Darmais **ne peut pas s'appliquer**.  
Mais rien n'empêche de trouver une statistique exhaustive...

#### Exemples

$T_n = \sum_{i=1}^n X_i$  est exhaustive pour le paramètre  $p$  d'une loi de Bernoulli ( $p$ )

$T_n = \sum_{i=1}^n X_i$  est exhaustive pour  $\theta$  pour la loi exponentielle  $\mathcal{E}(\theta)$

$T_n = \sum_{i=1}^n X_i$  est exhaustive pour  $\mu$  pour la loi de Gauss  $N(\mu, \sigma)$  ( $\sigma$  connu)

$T_n = \sum_{i=1}^n (X_i - \mu)^2$  est exhaustive pour  $\sigma^2$  pour la loi de Gauss  $N(\mu, \sigma)$  ( $\mu$  connu)

### 3. Estimation – Exhaustivité, information de Fisher

## L'information de Fisher

- L'exhaustivité d'une statistique  $T_n$  renseigne sur son pouvoir à restituer l'information contenue dans un échantillon  $\{X_1, X_2, \dots, X_n\}$  vis-à-vis d'un paramètre  $\theta$  inconnu que l'on souhaite estimer
- **L'information de Fisher** vise elle à estimer la **quantité d'information** qu'un échantillon  $\{X_1, X_2, \dots, X_n\}$  apporte sur le paramètre  $\theta$

### DEFINITION

On appelle **quantité d'information de Fisher**, notée  $I_n(\theta)$ , apportée par un échantillon  $\{X_1, X_2, \dots, X_n\}$  sur le paramètre  $\theta$  la quantité suivante (si elle existe) :

$$I_n(\theta) = E \left[ \left( \frac{\partial \ln L(\mathbf{x}; \theta)}{\partial \theta} \right)^2 \right]$$



### 3. Estimation – Exhaustivité, information de Fisher

**DEFINITION** On appelle **score de Fisher** la quantité définie par :

$$Z(\mathbf{x}; \theta) = \frac{\partial}{\partial \theta} \ln L(\mathbf{x}; \theta)$$

#### PROPRIETES

- Le score est centré :  $E(Z(\mathbf{x}; \theta)) = 0$
- L'estimateur du maximum de vraisemblance  $\hat{\theta}$  est la valeur de  $\theta$  qui annule le score :  $Z(\mathbf{x}; \hat{\theta}) = 0$
- L'information de Fisher est égale à la variance du score :

$$I_n(\theta) = V(Z(\mathbf{x}; \theta))$$

### 3. Estimation – Exhaustivité, information de Fisher

**THEOREME** : si le domaine de définition de  $X$  ne dépend pas de  $\theta$ , alors

$$I_n(\theta) = -E \left( \frac{\partial^2 \ln L(\mathbf{x}; \theta)}{\partial \theta^2} \right)$$

si elle existe...

#### PROPRIETES

- **Additivité.** Si le domaine de définition ne dépend pas de  $\theta$ , alors

$$I_n(\theta) = nI_1(\theta)$$

- **Précision.** Si  $X$  est distribuée selon une loi normale de variance connue, alors

$$I_1(\theta) = \frac{1}{\sigma^2} \quad \text{avec } \mu = \theta$$

- **Dégradation de l'information.** L'information apportée par une statistique  $T$  est inférieure ou égale à celle apportée par l'échantillon tout entier

$$I_T(\theta) \leq I_n(\theta)$$

$$I_T(\theta) = I_n(\theta) \Leftrightarrow T \text{ exhaustive}$$

### 3. Estimation – Exhaustivité, information de Fisher

#### CAS D'UN PARAMÈTRE VECTORIEL

- Soit un échantillon  $\{X_1, X_2, \dots, X_n\}$  dont la loi dépend d'un paramètre vectoriel  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_s)$ . On appelle **matrice d'information de Fisher** associée à l'échantillon, la matrice  $I_n(\boldsymbol{\theta})$  dont le terme général  $(i, j)$  est calculé par

$$I_n(\boldsymbol{\theta})_{(i,j)} = E \left( \frac{\partial \ln L(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_i} \frac{\partial \ln L(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_j} \right)$$

- **Score de Fisher** : c'est un vecteur aléatoire défini par :

$$Z(\mathbf{x}; \boldsymbol{\theta}) = \nabla \ln L(\mathbf{x}; \boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} \ln L(\mathbf{x}; \boldsymbol{\theta}) = Z_1(\mathbf{x}; \boldsymbol{\theta}) \\ \vdots \\ \frac{\partial}{\partial \theta_s} \ln L(\mathbf{x}; \boldsymbol{\theta}) = Z_s(\mathbf{x}; \boldsymbol{\theta}) \end{pmatrix}$$

### 3. Estimation – Exhaustivité, information de Fisher

- Le terme général de la matrice est donc défini par

$$I_n(\boldsymbol{\theta})_{(i,j)} = \text{cov} (Z_i(\mathbf{x}; \boldsymbol{\theta}), Z_j(\mathbf{x}; \boldsymbol{\theta}))$$

On en déduit que l'information de Fisher dans le cas vectoriel est déterminée par la matrice de variance covariance du score de Fisher

Lorsque  $\frac{\partial^2 \ln L(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j}$  existe pour tout  $\boldsymbol{\theta}$  et tout  $\mathbf{x}$ , on montre que le terme général peut se calculer par :

$$I_n(\boldsymbol{\theta})_{(i,j)} = -E \left( \frac{\partial^2 \ln L(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right)$$

### 3. Estimation – Exhaustivité, information de Fisher

#### Exemples

- Pour une loi de Bernoulli ( $\theta$ ) :  $I_n(\theta) = \frac{n}{\theta(1-\theta)}$
- Pour une loi de Gauss  $N(\mu, \sigma)$ , on a :  $I_n(\boldsymbol{\theta} = (\mu, \sigma^2)) = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix}$