

# Méthode statistique en médecine

## 1) Introduction

- **Biostatistiques** : statistiques appliquées au domaine de la santé publique

### 3 objectifs :

- Description d'une maladie par rapport à une population
- Évaluation des traitements, des techniques et des coûts
- Mise en place des observations épidémiologiques et en tirer des conclusions

## 2) Définitions

- **Statistique** : art de collecter, analyser et interpréter des données.

Il en existe 2 types en biostatistiques :

- **Descriptives** : description de données à l'aide de paramètres.

*Ex : on collecte des données sur les étudiants en LAS : QI, âge, taille, note en biostat ...*

- **Déductives** : l'observation est-elle due au hasard, y a-t-il une autre explication

*Ex : on constate que les personnes qui aiment la biostat ont une meilleure espérance de vie : est-ce dû au hasard ?*

- **Données** : c'est le résultat de l'observation d'un individu, grâce à un instrument de mesure, ou par le sens d'un observateur (signes cliniques, biologiques, ...)

> Une donnée n'est intéressante que si on l'observe ou la compare à d'autres individus.

> On parle alors de variable car elle est différente selon les individus.

*Ex : taille, âge, poids, groupe sanguin...*

La variabilité peut être :

*inter sujet* (=entre 2 sujets)  
comparaison de 2 sujets

*intra sujet* (= pour un même sujet)  
comparaison du sujet à lui-même

*On revient aux variables après tkt*

- **Paramètre** : grandeur apportant une information résumée sur la variable étudiée.

*Ex : moyenne, médiane, ...*

- **Série statistique** : collection d'objets de même nature avec des caractéristiques différentes d'un objet à l'autre.

*Ex : Les étudiants de LAS de Nice (même nature, caractéristiques différentes)*

- **Population** : série exhaustive de tous les individus étudiés, sur lesquels on peut appliquer (inférer) des décisions.

*Ex : La population française, une école*

- **Échantillon** : sous-ensemble fini et d'effectif limité, extrait de la population. Il doit être représentatif de la population, d'où la nécessité de tirage au sort = randomisation

*Ex : 100 LAS tirés au sort*

La population est *inaccessible* dans son entièreté pour des raisons d'organisation et de moyens limités. Du coup on réalise l'étude sur l'échantillon puis on fait un « pari » sur l'application des résultats à la population.

Logique, on va pas prendre tous les étudiants en médecine de France pour tester leur niveau de stress, alors qu'on peut prendre un échantillon à Nice, du coup faut faire gaffe à la fiabilité du résultat !

**L'ÉCHANTILLON EST CONNU,  
ALORS QUE LA POPULATION  
EST INCONNUE.**

*Ca c'est très important please*

## 3a) Variables

Variable <b>qualitative</b>	Variable <b>quantitative</b>
Non mesurable, <i>Couleur des yeux, prénom...</i>	Mesurable, <i>Taille, poids...</i>
Binaires <i>Femme/homme</i>	Discrètes <i>Age (sans virgule)</i>
Nominales <i>Couleur des cheveux</i>	Continues <i>Poids, glycémie (avec virgule)</i>
Ordinales <i>Echelle de douleur de 1 à 10</i>	

*Tips : Une variable discrète est trop timide pour utiliser des virgules !*

Une variable qualitative ordinale peut être approximée en une variable pseudo quantitative : la variable est qualitative mais ressemble à une quantitative !

En gros, si on a une échelle par exemple de douleur, de satisfaction... on va les classer de 1 à 10 (donc ça c'est pseudo quantitatif), mais ça reste des scores subjectifs qui dépendent de la personne donc qu'on peut pas vraiment mesurer (donc qualitatif)

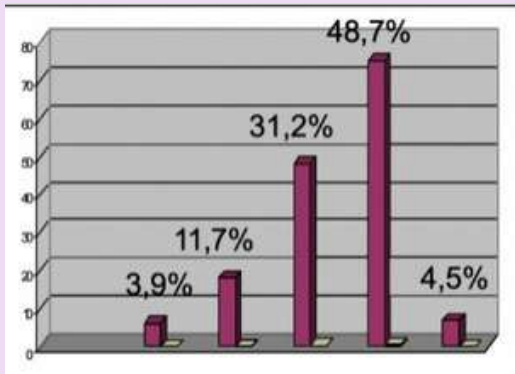
**UNE VARIABLE  
PSEUDO  
QUANTITATIVE RESTE  
QUALITATIVE !!**

Ca aussi grrrr

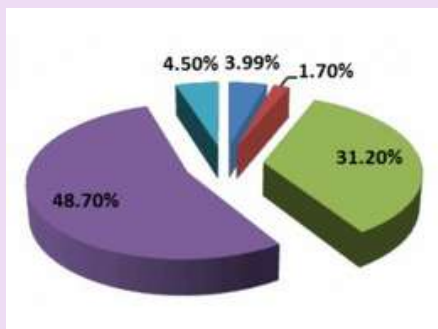
## 3b) Représentation des variables

### Qualitatives :

On les représente sous forme de tableau de % , d'histogramme, de secteurs ... *un pourcentage est une variable qualitative*



Degré de satisfaction	Nb mères	%
Très insatisfait	6	3,9%
Plutôt insatisfait	18	11,7%
Plutôt satisfait	48	31,2%
Très satisfait	75	48,7%
Pas d'opinion	7	4,5%

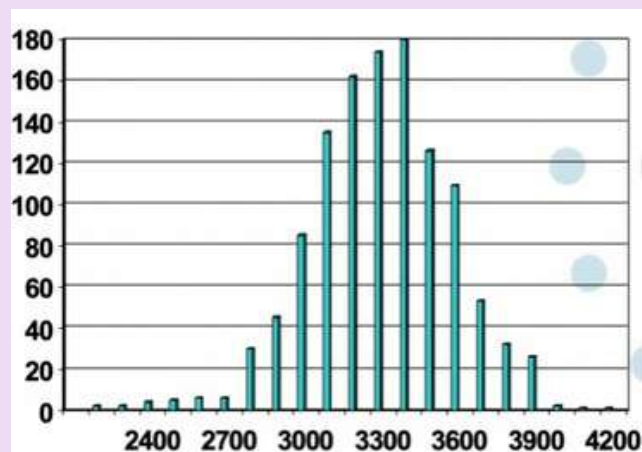


(N'apprenez pas ces images, c'est pour illustrer <3)

### Quantitatives :

On les représente sous forme de tableau, de diagramme en bâton ou d'histogramme

Poids (g)	Nb bébés
2200	2
2300	2
2400	4
2500	5
...	
3100	121
3200	150
3300	162
3400	170



## 4) Paramètres

### MOYENNE

Variable  
quantitative  
**discrète**

$$m = \frac{\sum x_i}{n}$$

Variable  
quantitative  
**continue**

$$m = \frac{\sum n_i x_i}{n}$$

### MÉDIANE

**N pair** : moyenne des  $\frac{n}{2}$  et  $\frac{n+1}{2}$  valeurs

Valeur centrale qui **sépare** la  
série d'effectif n en 2 sous  
séries de même effectif

**N impair** :  $\frac{n}{2}$  eme valeur

### VARIANCE

Indique la dispersion des  
valeurs autour de la moyenne.

### QUARTILES

Valeurs de la variable qui partagent la série  
d'effectif n en 4 sous séries de même effectif

Pas de formules yayy



Soyez comme réré <3

# Exemples :

**Enoncé :** Les notes des LAS en biostat au premier EB (tintintinnn) : 10, 7, 15, 20, 2

**Calculer la :** moyenne, médiane, quartiles :

MOYENNE :  $(10 + 7 + 15 + 20 + 2) / 5 = 10,8$       ça c'est easy

MEDIANE :  
 1) remettre dans l'ordre la suite : 2, 7, 10, 15, 20  
 2) parité de la suite : ici impair car 5 valeurs  
 3) application : on prend la valeur du milieu : 10

QUARTILES :  
 1) premier quartile on fait  $1/4 \times 5 = 1,25$  avec 5 le nombre de valeurs  
 2) donc Q1 se trouve entre la **1e et la 2e note**  
 3)  $Q1 = (2+7)/2 = 4,5$   
 4) 25% des LAS seulement ont une note inférieure à 4,5

Perso j'ai mis du temps à comprendre les quartiles donc si vous arrivez pas go fofo ;)

	♥ Avantages	♥ Inconvénients
♥ Moyenne	<ul style="list-style-type: none"> <li>-<b>Simple</b> à calculer</li> <li>-Facile à manipuler dans des tests stats donc <b>adaptée</b> aux calculs statistiques</li> <li>-Très <b>significative</b> si la répartition des données est assez <b>symétrique</b> avec une <b>faible</b> dispersion</li> </ul>	<ul style="list-style-type: none"> <li>-<b>Sensible</b> aux valeurs anormales (max et min)</li> </ul>
♥ Médiane	<ul style="list-style-type: none"> <li>-Calcul <b>facile</b></li> <li>-<b>Peu sensible</b> aux valeurs anormales</li> <li>-Utilisable pour des valeurs <b>ordinales</b>, des classes</li> </ul>	<ul style="list-style-type: none"> <li>-Se prête <b>moins</b> aux calculs statistiques</li> </ul>

# Statistiques descriptives

## 1) Variabilité

**Toutes les données biologiques possèdent une variabilité.**

Il faut la connaître pour pouvoir classer nos données comme « normales » ou « anormales » :

- Une variabilité **maîtrisée** permet une *estimation*
- Une variabilité **non maîtrisée** conduit à des *biais*

*Exemple : les valeurs normales de la glycémie sont comprises entre 0,75 et 1,25 g/L. Si on est en dessous de 0,75 g/L on a une valeur anormale, on est en hypoglycémie*

## 2) Estimation statistique

- Les études en biostatistique sont réalisées sur un **échantillon** représentatif de la population après « **échantillonnage** »
- Après l'étude on réfléchit à la légitimité des résultats et à leur **extrapolation** à la population.
  - > On réalise donc une **estimation** du résultat vrai à partir des données de l'échantillon.



On détermine des paramètres au niveau d'une **population** à partir d'**observations** réalisées sur un **échantillon** de cette population.



Ce charabia résumé

**ECHANTILLON** —————> **ESTIMATION** —————> **POPULATION CIBLE**

## On retrouve deux types d'estimations :

♥ **L'estimation ponctuelle** : valeur unique jugée la meilleure à l'instant t (PEU FIABLE)

♥ **L'estimation par intervalle** : un intervalle de valeurs comprenant la valeur recherchée, c'est **l'Intervalle de Confiance ou IC** (BEAUCOUP + FIABLE)

♥ **2 estimations ponctuelles** réalisées sur 2 échantillons donneront des résultats **proches mais différents**

♥ **2 estimations par intervalles** réalisées sur 2 échantillons donneront 2 IC se **recouvrant** mais pas nécessairement le même IC.

> Cependant, si on refait la même estimation sur un autre échantillon, elle recouvrira la première, ce qui ne serait sûrement pas le cas avec des valeurs ponctuelles

**L'ESTIMATION PAR INTERVALLE EST  
MOINS PRÉCISE MAIS PLUS JUSTE**

## 3) Estimation des données quantitatives


### - Méthodologie :

1. Détermination précise de la population étudiée (=population cible)
2. Tirage au sort (TAS) d'un échantillon représentatif (n sujets)
3. Calcul de l'intervalle de confiance

**Pour les données quantitatives, on va estimer la moyenne.**

> L'estimation assure la correspondance entre ce qu'il se passe au niveau de l'échantillon et ce qu'il se passe au niveau de la population.



 <b>ECART-TYPE</b>	Mesure la <b>dispersion</b> d'un ensemble de données autour de la moyenne. C'est la variabilité des mesures entre elles et par rapport à la moyenne.
--	---

> Plus l'écart type est faible plus le caractère étudié est homogène ( les valeurs sont proches de la moyenne ).


Ex : A l'épreuve de biostar 3 étudiants ont eu 0, 10 et 20, la moyenne est de 10

> La médiane et de 10.

Ici c'est l'écart-type qui permettra le mieux de résumer la dispersion de la série.

Si les étudiants avaient eu 9, 10 et 11 la moyenne et la médiane seraient les mêmes, l'écart-type serait plus petit.

En gros plus les valeurs sont éloignées plus l'écart-type est grand, et inversement.

 <b>DEGRÉ DE LIBERTÉ DDL</b>	Le nombre de valeurs nécessaires à <b>connaître</b> pour pouvoir résoudre l'équation et connaître toutes les valeurs de la série.
--	---

On définit « m » la moyenne, «  $x_i$  » les valeurs dont on veut faire la moyenne, « n » l'effectif, «  $x_i - m$  » les écarts.

- Il y a n écarts
- Il y a (n - 1) écarts indépendants à la moyenne, ou degrés de liberté

(Ca c'est du cours pur, si vous comprenez l'exemple c'est carrée)

Ex : Un élève a eu 4 notes : 12, 15, 16 et une copie perdue (grr) dont il veut connaître la note.


Il connaît sa moyenne de 15.

Donc on fait une petite équation, et hop !

$$> (12 + 15 + 16 + ?)/4 = 15$$

$$> 43 + ? = 60$$

Sa dernière note est donc 17

 <b>INTERVALLE DE CONFIANCE</b>	C'est l'estimation de la <b>moyenne vraie</b> $\mu$ à partir de la moyenne m calculée sur l'échantillon.
--	--

On donne un intervalle auquel  $\mu$  appartient :



$$\mu \in \left[ m \pm \frac{\varepsilon s}{\sqrt{n}} \right]$$

IC

L'IC est aussi appelé **intervalle au risque  $\alpha$** .

<b>RISQUE <math>\alpha</math></b>	C'est le <b>risque d'erreur</b> dans l'estimation de $\mu$ (le risque que notre IC ne contienne pas $\mu$ )
-----------------------------------	---

> On prend en général  **$\alpha = 5\%$**  (on a **95%** de chance que la moyenne vraie soit dans notre IC)

<b>L'ÉCART-RÉDUIT <math>\varepsilon</math></b>	C'est une valeur qui dépend du risque $\alpha$ : ils varient en <b>sens inverse</b> , si $\alpha$ augmente, $\varepsilon$ diminue
--	---

> Un écart-réduit mesure de combien d'écarts-types une observation particulière est éloignée de la population.



**CA C'EST PAR  $<3$**

J'insiste ça va vous servir pour les autres cours !!

Pour  $\alpha = 5\%$  ;  $\varepsilon = 1,96$   
 Pour  $\alpha = 1\%$  ;  $\varepsilon = 2,60$

## 4) PRECISION DE L'ESTIMATION

IC Large	IC Resserré
Si $\alpha \searrow$ alors $\epsilon \nearrow$ donc l'IC $\nearrow$	Si $\alpha \nearrow$ alors $\epsilon \searrow$ donc l'IC $\searrow$
<ul style="list-style-type: none"> <li>→ On a plus de chances que <math>\mu</math> soit comprise dans l'IC</li> <li>→ Par contre on perd en précision</li> </ul>	<ul style="list-style-type: none"> <li>→ On a moins de chance que <math>\mu</math> soit dans l'IC</li> <li>→ Mais on diminue l'IC, on gagne en précision</li> </ul>

> Les variations du risque  $\alpha$  vont conditionner la précision de l'estimation et la largeur de l'intervalle de confiance.

> Si on prend moins de risque, on a un intervalle de confiance plus grand, on a plus de chances que la moyenne soit dedans, (et inversement).



## L'INDICE DE PRÉCISION I

Il permet de calculer la **précision de l'estimation de  $\mu$** . Cette valeur représente la **largeur de l'IC**.




$$i = \frac{\epsilon S}{\sqrt{n}}$$

- > D'après la formule de l'IC vu avant l'IC est donc compris entre  **$[m + i]$  et  $[m - i]$**
- > Plus la taille de **l'échantillon augmente**, plus la **précision augmente**
- > Quand **l'indice de précision diminue** la **précision augmente**.

D'après la formule de **l'indice de précision** :

$$n \nearrow, i \searrow \text{ donc l'IC } \searrow \text{ donc la précision } \nearrow$$

Le **nombre de sujets** nécessaires «**n**», pour une précision donnée :



$$n = \frac{\epsilon^2 s^2}{i^2}$$

### RECAP DU TURFU :

- ★ **L'IC** c'est l'estimation de la **moyenne vraie**  $\mu$  à partir de la **moyenne m** calculée sur l'échantillon. Il est aussi appelé "**intervalle au risque  $\alpha$** ".
- ★ Le **risque  $\alpha$**  c'est le risque d'erreur dans l'estimation de  $\mu$ .
- ★  $\epsilon$  représente **l'écart-réduit**.
- ★ Les variations du **risque  $\alpha$**  déterminent la **précision de l'estimation**
- ★ **i** représente la **largeur de l'IC**
- ★ **IC**= **[m±i]**

### DONC :

(encrez moi ça dans vos petites têtes)

- ★ Si **n**  $\nearrow$ , **i**  $\searrow$  donc **l'IC**  $\searrow$  donc la **précision**  $\nearrow$
- ★ Si  **$\alpha$**   $\nearrow$  alors  **$\epsilon$**   $\searrow$  donc **i**  $\searrow$  donc **l'IC se resserre** donc la **précision**  $\nearrow$

## 5) LOI DE GAUSS OU LOI NORMALE

En sciences humaines, on observe souvent des distributions des variables assez symétriques autour de la moyenne : c'est **la courbe de Gauss**

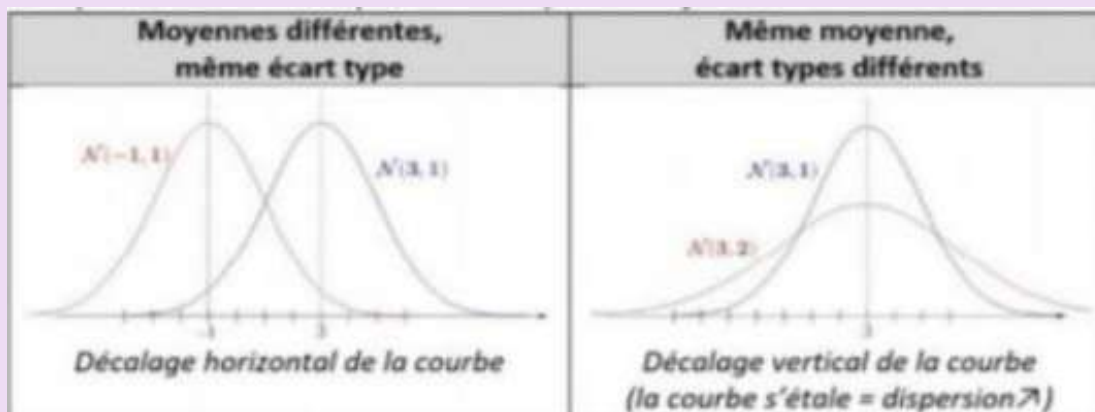
La représentation graphique de données suivant la courbe de Gauss est une courbe en cloche avec :

- En abscisse  $[m \pm \epsilon s]$  donc l'IC
- En ordonnée ni : l'effectif pour chaque valeur
- L'aire sous la courbe, le % de la population concerné

La courbe de Gauss permet de **visualiser l'IC** autour de la moyenne, **l'écart-type**, la dispersion autour de cette valeur moyenne et **la moyenne**.

Pour pouvoir faire des calculs on suppose que notre variable X (quantitative continue) suit une distribution modèle : **la loi Normale**.

Ainsi, pour chaque couple  $(\mu, s)$ , il existe une loi normale de moyenne  $\mu$  et d'écart-type  $s$  notée  **$N(\mu, s)$**

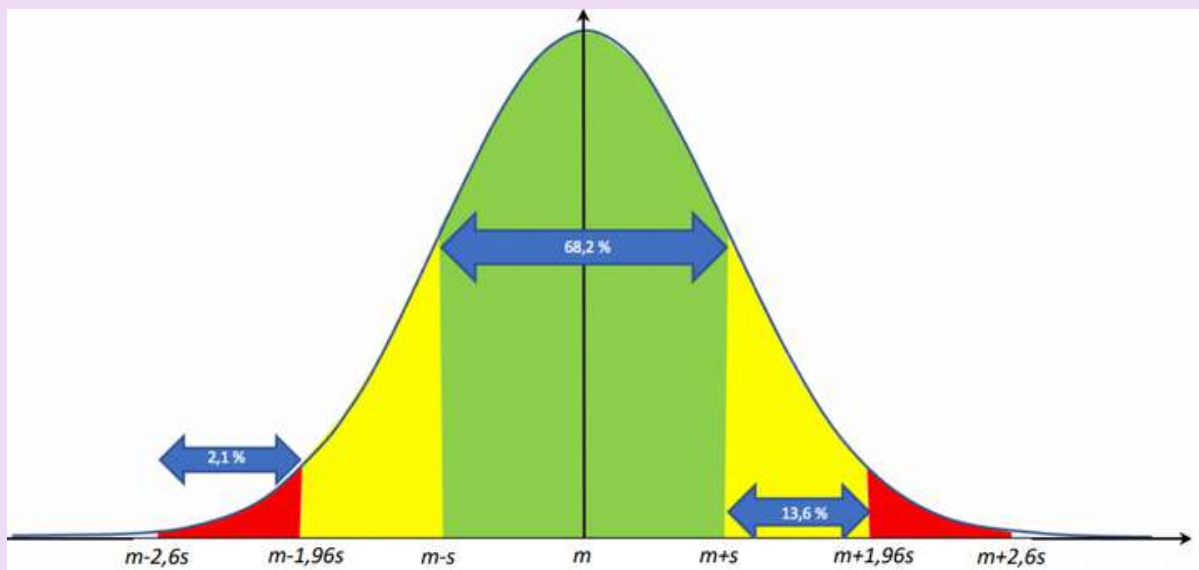


A partir de la Loi Normale ou de GAUSS, on précise les intervalles de confiance

$[m - 1s ; m + 1s]$  contient 68,2% de la population

$[m - 1,96s ; m + 1,96s]$  contient 95,4% de la population

$[m - 2,6s ; m + 2,6s]$  contient 99,6% de la population



Et voilaaaaaa pour ce cours de ttr, il manque juste une petite partie sinon ça faisait trop, mais tkt c'est que deux petites pages qui sortiront après la ttr du coup ;)  
Je vous mets les dédis à part pour pas que vous ayez à tout imprimer !

Dédi à Ines qui va vivre à mis temps chez oim, tqt la brosse à dent t'attends,  
Dédi à mon chat, mon petit raton d'amour <3  
Dédi à moi pour avoir réussi ma P1 non mais oh quand même un peu d'amour propre,  
Dédi au super CT de SE, ALLEZ LE VOIR OH !!!  
Dédi à la BIOSTARRRR cette matière phare qui met tout le monde d'accord hehe...  
Dédi à mon amoureux qui a gratté cette dédi comme pas possible, il aime trop la lumière,  
Dédi enfin à ma maman, la meilleure et que j'aime de tout mon coeur <3 (lachez une larme please)

Petit passage d'Echalote, pour vous soutenir ;)



BISOUS MES PETITS CHOUX !