

The background is a vibrant, cartoonish scene from a Mario game. It features a blue sky with white clouds, a green ground with a brown brick wall, and various characters and items. A gold coin is in the top left. A red and white mushroom is on the brick wall. A yellow question mark block is on the brick wall. A yellow Koopa is flying. A Piranha Plant is on a green pipe. Mario is riding Yoshi. Luigi is jumping. Toad is on the left. A pink Piranha Plant is on the ground. A gold coin is on the ground. A brown Goomba is on a green pipe. A green pipe is on the right. The scene is framed by a dashed black border.

# STATISTIQUES DESCRIPTIVES

BIENVENUE DANS LE NIVEAU 1 DU MONDE DES STATS ! VOTRE OBJECTIF : ARRIVER À LA FIN DE CE COURS. VOUS CROISEREZ TOUT PLEIN DE PERSONNAGES SUR VOTRE TRAJET, BON COURAGE !

START

Et coucou les pouss1 ! On se retrouve aujourd'hui pour un petit cours avec des notions que vous connaissez déjà alors pas d'impasse ! Comme vous avez pu le remarquer, vous serez plongé dans le thème des jeux vidéos. J'espère que la mise en page vous plaira ! Sur ce, la partie peut commencer ! (mes remarques seront de cette couleur).

## Introduction



La **biostatistique** nous permet d'appliquer des **théories statistiques** au domaine du vivant : le domaine de la **santé publique** permet de décrire l'état de santé de la population. Pour cela, il faut évaluer les traitements, les techniques, les coûts et mettre en place des **observations épidémiologiques** pour en tirer des conclusions. On procédera d'abord à une analyse **descriptive** puis à une analyse **déductive** (vous verrez ça avec ma cotut incr).

### QUELQUES DEFINITIONS :

° **Statistique** : art de collecter et d'interpréter des 'données'

- statistique **descriptive** : description d'une situation à l'aide de **paramètres** (ex: on peut calculer la **moyenne** des étudiants à l'épreuve de biostatistiques).
- statistique **déductive** : conclusions à partir d'**observations** et de **mesures** → ce que l'on observe est-il dû au **hasard ou pas** ? (ex : 2 traitements anti cancéreux donnent à 5 ans une survie de 42% pour l'un et de 48 % pour l'autre. Hasard ou efficacité plus grande pour l'un des deux ?)

° **Données** : résultat de l'observation d'un individu, par l'utilisation d'un instrument de mesure, ou par les **sens de l'observateur** (signes cliniques, biologiques,...). Cette donnée n'est intéressante que si on peut l'observer/la **comparer** sur plusieurs individus. Elle ne sera **pas strictement équivalente** d'un individu à l'autre. On parle donc de **variable** (ex : la taille, le poids, le groupe sanguin...). Cette variabilité peut être due au **hasard** ou être **physiologique** (intra ou inter sujets).

° **Population** : série **exhaustive** de tous les individus étudiés, sur lesquels on veut inférer des décisions (ex: population de France, les étudiants en médecine de Nice...)

° **Echantillon** : ensemble **fini** et d'effectif **limité**, extrait de la population

## Pourquoi échantillonner ?

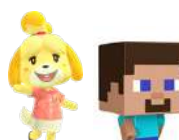
→ car la population est **inaccessible** dans son entièreté pour des raisons d'organisation et de **moyens limités**.

→ on procède à une étude sur l'**échantillon** et on regarde si on peut appliquer ces résultats à la **population**.

→ l'échantillon doit être **représentatif** de la population : on procède alors à un **tirage au sort** (randomisation)

→ on se retrouve alors avec un **échantillon connu** et une **population inconnue** ++

## Quelques variables



Qualitative binaire	sexe...
Qualitative nominale	couleur des yeux...
Qualitative ordinale	consommation de tabac (faible 0-9, moyenne 10-19, forte > 19), douleur articulaire (absente, modérée, intense), rang dans un classement...
Quantitative discrète (nombre entier)	âge, nombre d'otites dans la dernière année...
Quantitative continue (nombre à virgule, plus précis)	déficit auditif moyen (ex :11,5 dB), note de biostatistiques arrondi au dixième...



Une variable **qualitative ordinale** peut être considérée comme variable 'pseudo quantitative' pour certains tests ++ (ex : variation de la douleur sur une échelle de 1 à 3, taux de satisfaction de 1 à 5... ) ⚠ Les **nombre**s affectés aux modalités qualitatives ne peuvent **pas** faire l'objet d'**opérations arithmétiques** (logique si sur une échelle de 1 à 10, ma douleur se place à 5, ce chiffre donne un ordre d'idées de l'intensité de la douleur, mais ne veut rien dire en tant que tel).



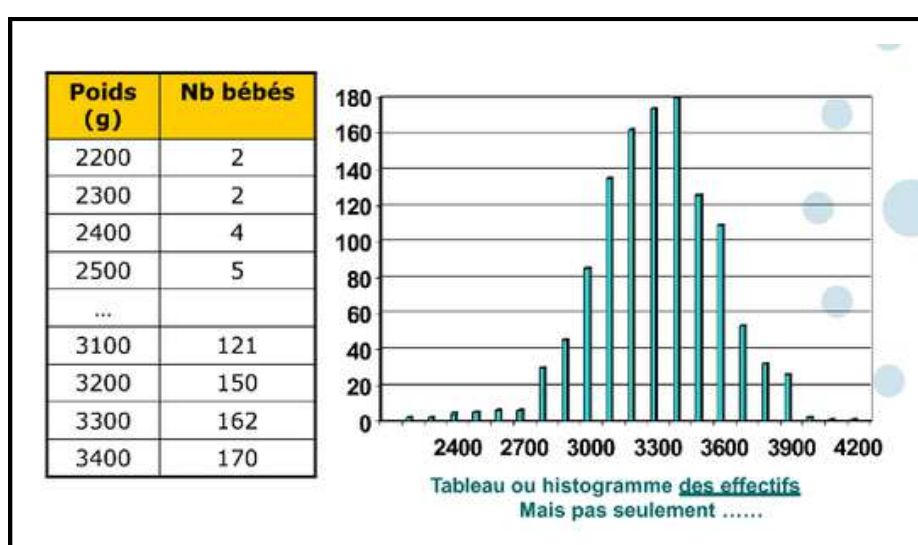
## Représentation de variables quantitatives



Exemple : on s'intéresse aux poids des nouveaux nés dans une maternité :

a) Echantillon : **TAS** des mères ayant accouché dans cette maternité pendant une période donnée (effectif  $n=1165$ )

b) Variable étudiée : **poids** du nouveau né --> **variables quantitatives**



On peut résumer en quelques **paramètres** les caractéristiques de la série de **données quantitatives** :

## ° Moyenne :

- dans le cas d'une variable quantitative **discrète** (valeur entière)
- dans le cas d'une variable quantitative **continue** (valeur à virgule)

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

$$m = \frac{\sum_{i=1}^n n_i x_i}{n}$$

° **Variance** : paramètre indiquant la **dispersion** des données autour de la moyenne.

° **Médiane** : valeur **centrale** si les données sont rangées en ordre **croissant** (attention aux pièges les gars) → 50% des valeurs < médiane ainsi que 50% des valeurs > médiane (ex : {3,5,12,18,21}).

° **Quartiles** : ils partagent la série en 4 groupes de **4 groupes de même effectif** (série en ordre croissant) ex : le quartile = 25% de la série est < à cette valeur.

Je vous mets ici la diapo du prof, c'est très bien expliqué et vous avez pas mal de paramètres qui sont calculés : si vous avez des questions, se sera sur le forum :)

Exemple : Soit une série de poids de bébés (n = 15 valeurs)

1	3400	1	1890
2	2570	2	2140
3	3210	3	2350
4	4070	4	2470
5	3840	5	2570
6	4180	6	2640
7	3480	7	3000
8	3990	8	<b>3210</b>
9	2640	9	3400
10	3000	10	3480
11	3830	11	3830
12	1890	12	3840
13	2350	13	3990
14	2140	14	4070
15	2470	15	4180

Valeurs rangées par ordre croissant

Médiane : n=15 donc  $(n+1)/2 = 8$   
**8<sup>ème</sup> valeur = 3210**  
 Si n pair, médiane située entre les n° n/2 et n/2+1. **Moyenne des 2 valeurs correspondantes**

**3<sup>ème</sup> quartile  $Q_{75} = 0,75 \times 15 = 11,25$**   
 Le 3<sup>ème</sup> quartile est situé entre le n°11 et le n°12 soit  $(3830+3840)/2 = 3835$

**Moyenne de la série = 3137,3**

# Estimation statistique



Le **problème** est le suivant : déterminer un **paramètre** au niveau d'une **population** à partir d'observations réalisées sur un **échantillon** de cette population. Après l'étude, on réfléchit à la **légitimité** des résultats et à leur **extrapolation** à la population. Exemple, comment connaître la durée de séjour moyenne des patients hospitalisés en France, pour une pathologie donnée ?

**Echantillon**  
effectif =  $n$   
moyenne =  $m$   
écart type =  $s$



**ESTIMATION**



**Population cible**  
effectif =  $N$   
moyenne =  $\mu$   
écart type =  $\sigma$

Il y a 2 types d'estimation :

## ESTIMATION PONCTUELLE



## ESTIMATION PAR INTERVALLE



--> L'ESTIMATION PAR INTERVALLE EST MOINS PRÉCISE MAIS PLUS JUSTE ++

° **Intervalle de confiance (IC)** : C'est l'estimation de la **moyenne vraie  $\mu$**  à partir de la **moyenne  $m$**  calculée sur l'échantillon. L'IC est aussi appelé intervalle au **risque  $\alpha$** .

On donne un intervalle auquel  $\mu$  appartient :



$$\mu \in \left[ m \pm \frac{\varepsilon s}{\sqrt{n}} \right] \rightarrow \text{INTERVALLE AU RISQUE } \alpha$$

° **Risque  $\alpha$**  : C'est le **risque d'erreur** dans l'estimation de  $\mu$  (le risque que notre IC ne contienne pas  $\mu$ ). On prend en général un risque  **$\alpha = 5\%$**  (on a 95% de chance que la moyenne vraie soit dans notre IC). Plus  **$\alpha$**  est **petit**, plus l'intervalle de **confiance** est **grand** : on réussit plus souvent. On s'expose aussi au risque de **rater** la "bonne" estimation.

° **Ecart réduit  $\varepsilon$**  : C'est une valeur qui dépend du **risque  $\alpha$**  : ils varient en **sens inverse**, si  $\alpha$  augmente,  $\varepsilon$  diminue. Un écart-réduit mesure de combien d'écart-types une observation particulière est **éloignée** de la population.

A CONNAÎTRE PAR COEUR ++

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

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



## Précision de l'estimation

L'intervalle de confiance peut être vu comme une **cible** (ici les centres des cibles seront les têtes de Reyna et de Killjoy pour les connaisseurs ;).

IC large	IC resserré
<p>→ plus de chances d'atteindre la cible, mais on a une <b>mauvaise précision</b> de l'estimation (au lieu de tirer dans la tête, on tire sur son corps donc moins précis)</p>	<p>→ on a un <b>risque de rater</b>, certaines balles seront à l'extérieur, mais on a une <b>meilleure précision</b> de l'estimation (donc plus de balles toucheront la tête)</p>
	

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 ++

$$i = \frac{\varepsilon s}{\sqrt{n}}$$

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 :

Quand  $n \nearrow$ ,  $i \searrow$  donc l'IC  $\searrow$  donc la précision  $\nearrow$  +++

## Loi de Gauss ou loi normale



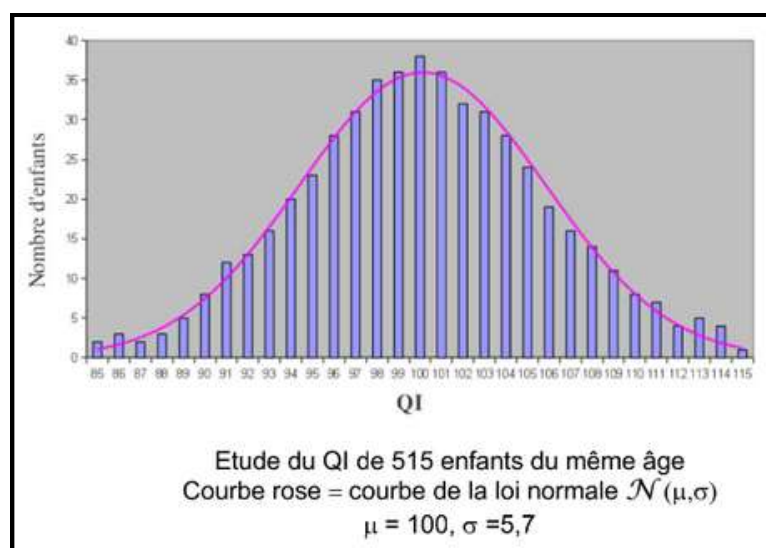
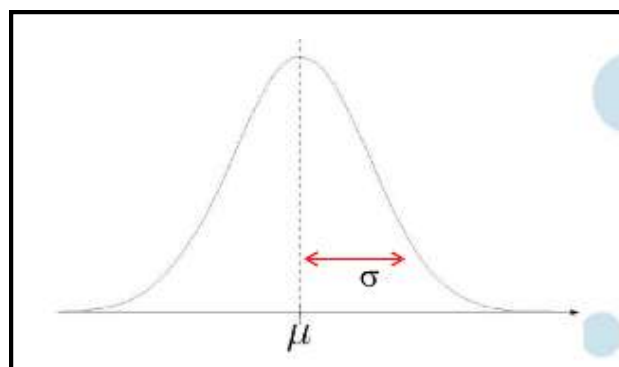
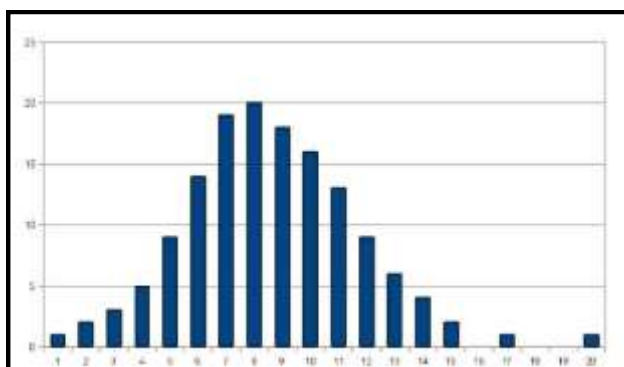
En **sciences humaines** on observe souvent des **distributions** (X) plutôt symétriques autour de la moyenne avec une forme de **cloche** pour pouvoir faire des calculs, on va supposer que X suit une distribution « modèle », pour des variables quantitatives continues : **la Loi Normale**.

- En **abscisse**  $[m \pm \varepsilon s]$  donc l'IC
- En **ordonnée** : 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**.

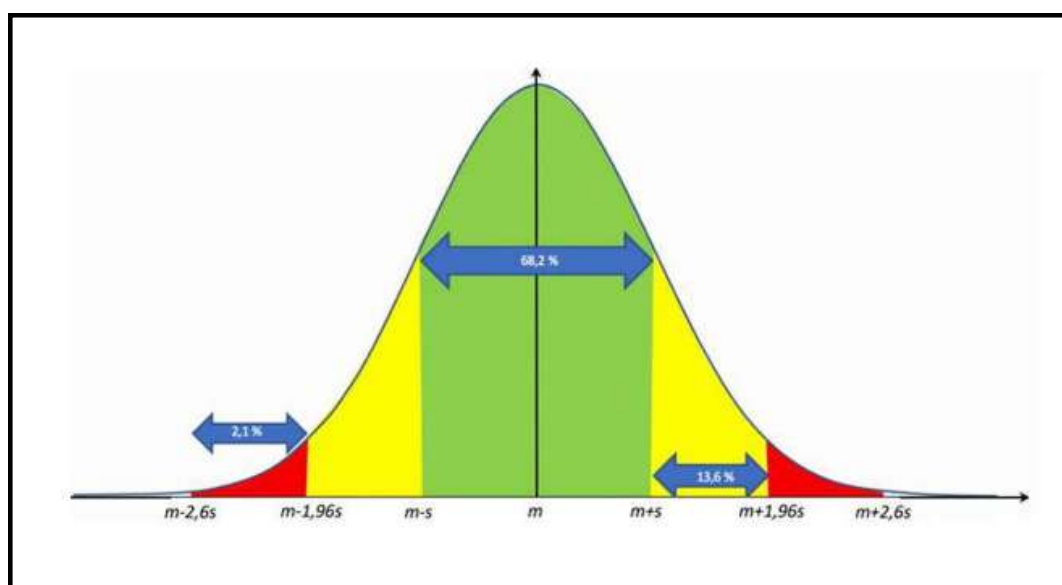
Pour pouvoir faire des calculs on suppose que notre **variable X** (quantitative continue) suit une distribution modèle : **la loi Normale**.

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



A partir de la loi normale (= loi 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 ++



## TUT RECAP :

effectif n augmente → IC se resserre → précision augmente

Contre exemple :

Un chirurgien écrit à 1000 de ses patients afin de connaître leurs suites chirurgicales et sur 100 réponses : 75 vont très bien, 25 ont des séquelles handicapantes. Or, il y a eu 900 non-réponses. On ne peut pas préjuger de l'état de ces 900 patients. Ils sont peut être décédés des suites opératoires, ou bien très mécontents du chirurgien, ou tout au contraire sont très satisfaits et ne jugent pas utile de répondre. Cet échantillon est **BIAISÉ**.

RECAP :

- ★ 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  $\mu$ .
- ★  $\varepsilon$  représente l'écart-réduit.
- ★ Les variations du risque  $\alpha$  déterminent la précision de l'estimation
- ★  $i$  représente la largeur de l'IC
- ★ IC =  $[m \pm i]$
- ★ Si  $n \nearrow$ ,  $i \searrow$  donc l'IC  $\searrow$  donc la précision  $\nearrow$  ++
- ★ Si  $\alpha \nearrow$  alors  $\varepsilon \searrow$  donc  $i \searrow$  donc l'IC se resserre donc la précision  $\nearrow$  ++

Et voilà les loulous le cours version TTR est terminé ! Je sortirai très rapidement la fiche complète car il manque des infos mais vous avez vu le plus gros ! Prenez le temps de bien comprendre, entraînez-vous bien, posez des questions fin des trucs qu'on vous rabâche depuis le début finalement 😊 Bonne continuation et ne lâchez jamais rien les bg !!