

# Statistiques déductives

## I. GENERALITES SUR LES TEST D'HYPOTHESES

Dans les statistiques déductives, contrairement aux statistiques descriptives, on essaie, à partir des **observations** faites, de tirer des **conclusions**. Pour cela, les épidémiologistes utilisent des **tests d'hypothèse**.

### 1. Les tests de comparaison

Le plus souvent, les tests utilisés en statistiques déductives sont des **tests de comparaison entre 2 populations** présentant des caractères/paramètres différents. On constitue alors **2 échantillons** représentatifs et on essaie de déterminer s'il existe une différence significative entre ces 2 échantillons pour le caractère étudié. Le but étant d'extrapoler le résultat aux 2 populations primitives.

### 2. La définition des hypothèses

Avant de commencer une étude statistique on formule des hypothèses que le test permettra ensuite de **confirmer** ou **d'infirmer**. On définira au début de chaque test **2 hypothèses** jouant un rôle **symétrique** :

#### 1. **H0 = hypothèse nulle**

- Il n'y a **pas de différence** observée entre les deux groupes
- Il n'existe **pas de lien** entre les 2 caractères étudiés, les fluctuations observées sont donc dues au hasard

#### 2. **H1 = hypothèse alternative**

- Il y a une **différence significative** entre les deux groupes
- Il existe bien un **lien** entre les 2 caractères étudiés, les fluctuations observées ne sont donc pas dues au hasard

Les tests sont donc des techniques permettant de **décider si on accepte ou si on rejette H0**, en ayant fixé le **risque d'erreur  $\alpha$**  accompagnant cette décision.

### 3. Les étapes d'un test d'hypothèse

Pour mettre en œuvre un test d'hypothèse, on suivra toujours les étapes suivantes:

1. Définir **H0** et **H1**.
2. a) Déterminer le **caractère des données** à comparer (qualitative/quantitative).  
b) Choisir le **test** en fonction du type de données. On nomme **Z** le paramètre qui sera calculé.
3. Choisir **risque  $\alpha$**  à priori, généralement fixé à 5%.
4. - Recueillir les **données**  
- Calculer **Z**  
- Utiliser la **règle de rejet / décision** (définie à partir de H0 et de  $\alpha$ ) : Il s'agit de comparer Z par rapport à une valeur théorique de référence  
- Fixer le **risque d'erreur** réel à posteriori
5. **Interprétation des résultats** :  
• Au niveau de l'échantillon : Accepte-t-on H0 ?  
• Au niveau de la population : Peut-on extrapoler les résultats ?

### 4. La notion de risque

Rappel de statistique descriptive : Lors de l'estimation d'une valeur  $x$  par un IC,  **$\alpha$**  représente le **risque d'erreur** dans l'estimation de  $x$ , c'est à dire le risque pour que l'IC ne contienne pas la vraie valeur de  $x$ . Il est généralement fixé à 5%

♥ En statistique déductive on a : ♥

- **$\alpha$**  ou **risque de première espèce** : le risque de **rejeter H0 si H0 est vraie**. Ce risque d'erreur est **maîtrisé**, c'est à dire qu'il est fixé (le plus souvent à 5%) avant l'application du test statistique.
- **$\beta$**  ou **risque de seconde espèce** : le risque **d'accepter H0 si H0 est fausse**. Ce risque d'erreur est **négligé** et peut être assez **important**.

➔ Il y a une **dissymétrie** dans le traitement des deux hypothèses (parce qu'on choisit de maîtriser  $\alpha$  quitte à ignorer  $\beta$ )

- **$1 - \beta$**  ou la **puissance du test** : la probabilité de **rejeter H0 si H0 est fausse**.

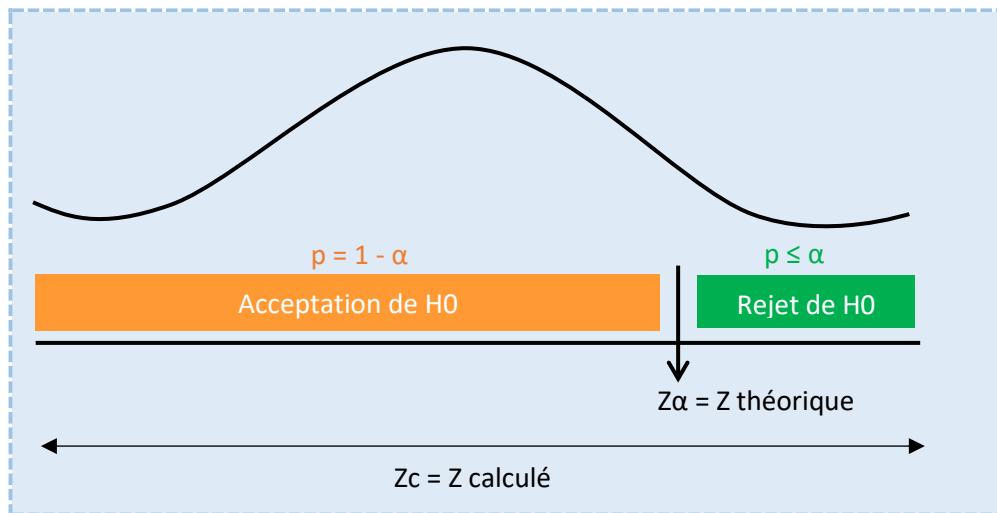
Tableau résumé ++++

Décision du statisticien

Réalité	Décision du statisticien	
	Rejet H0	Non rejet H0
H0 vraie	$\alpha$	$1 - \alpha$
H1 vraie	$1 - \beta$	$\beta$

## 5. Interprétation graphique du risque $\alpha$

Le paramètre  $Z_{\text{calculé}}$  que nous allons apprendre à calculer, suit une distribution probabiliste en forme de **courbe de Gauss**.



Pour pouvoir arriver à une conclusion après une étude statistique on :

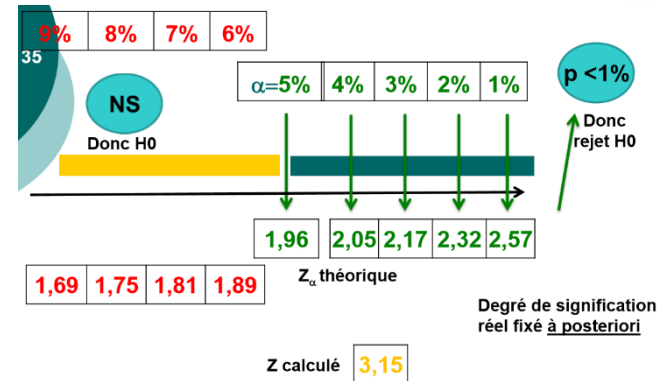
1. Fixe  $\alpha$  à priori
2. Cherche le  $Z_{\text{théorique}}$  sur la table (cf plus loin pour le trouver)
3. Calcule  $Z_c$  grâce aux formules
4. **Compare**  $Z_c$  avec  $Z_t$  et on peut arriver à deux conclusions différentes :

Acceptation de H0	Rejet de H0
$Z_c < Z_t$	$Z_c > Z_t$
$p = 1 - \alpha$	$p \leq \alpha$

5. Fixe le **degré de signification**  $p$  à postériori

**Remarque** :  $\alpha$  est fixé à priori par le statisticien (=supposition) mais le degré de signification est fixé à postériori (=réel) parce que la précision de l'étude peut s'avérer être supérieure à celle qu'on a supposé.

Exemple :



1.  $\alpha = 5\%$
2.  $Z_{\alpha} = 1,96$ .
3.  $Z_c = 3,15$ .
4.  $3,15 > 1,96$  donc on rejette H0
5. On voit sur la table (ou le schéma) que pour  $\alpha = 1\%$ ,  $Z_{\alpha} = 2,57$ . Or,  $3,15 > 2,57$  donc le degré de signification fixé à postériori est  $< 1\%$  : la précision a donc augmenté.

### a) Situation unilatérale

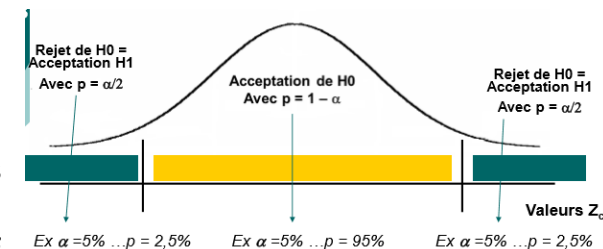
Dans une **situation unilatérale**, le rejet de H0 permet **uniquement** de dire qu'il existe une **différence significative** entre les deux situations.

**Par exemple** : si on teste l'efficacité de deux traitements A et B, le rejet de H0 en situation unilatérale permet uniquement de dire que l'efficacité de ces deux traitements est différente, mais on ne saura pas lequel est le meilleur.

### b) Situation bilatérale

**Au contraire**, dans une **situation bilatérale**, le rejet de H0 permet de dire qu'il existe bien une différence entre les deux situations, mais on peut aussi déterminer laquelle des deux est **la meilleure**.

Si l'on reprend l'exemple des **traitements**, le rejet de H0 permettra, en situation bilatérale, de déterminer lequel des deux traitements sera **le plus efficace** !



## II. L'ETUDE DE LA LIAISON ENTRE DEUX CARACTERES QUALITATIFS

A partir d'ici, les formules ne sont pas à savoir à part celles des ddl

Soient **2 groupes A et B** et une **caractéristique qualitative x** (couleur des yeux etc.).  
On se demande si le pourcentage d'individus du groupe A présentant x coïncide avec le pourcentage d'individus du groupe B présentant x.

### 1. Le test de comparaison des pourcentages (tout effectif)

Le paramètre Z est donné ici par **l'écart réduit  $\varepsilon$** . On va ainsi comparer :

-  $\varepsilon_t$  : donné par la **table de l'écart réduit** en fonction de  $\alpha$

-  $\varepsilon_c = \frac{p_A - p_B}{\sqrt{\frac{p_A q_A}{n_A} + \frac{p_B q_B}{n_B}}}$  avec  $q = 1 - p$

→ Si  $\varepsilon_{calculé} > \varepsilon_{théorique}$  Rejet de  $H_0$

#### Comment trouver $\varepsilon_t$ sur la table de l'écart réduit ?

On cherche  $\varepsilon_t$  en fonction du risque  $\alpha$ .

On regarde les **unités** et les **dizaine** d' $\alpha$  sur les **lignes** et les **centièmes** sur les **colonnes**.  $\varepsilon_t$  se trouve à l'intersection de la ligne et la colonne

Ainsi, pour  $\alpha = 5\% = 0,05$ , on est à **0,0** pour la ligne et à **0,05** pour la colonne.

		0,01	0,02	0,03	0,04	0,05	0,06	0,07	0,08	0,09
0		2,576	2,326	2,17	2,054	<b>1,96</b>	1,881	1,812	1,751	1,695
0,1	1,645	1,598	1,555	1,514	1,476	1,44	1,405	1,372	1,341	1,311
0,2	1,282	1,254	1,227	1,2	1,175	1,15	1,126	1,103	1,08	1,058
0,3	1,036	1,015	0,994	0,974	0,954	0,935	0,915	0,896	0,878	0,86
0,4	0,842	0,824	0,806	0,789	0,772	0,755	0,739	0,722	0,706	0,69
0,5	0,674	0,659	0,643	0,628	0,613	0,598	0,583	0,568	0,553	0,539
0,6	0,524	0,51	0,496	0,482	0,468	0,454	0,44	0,426	0,412	0,399
0,7	0,385	0,372	0,358	0,345	0,332	0,319	0,305	0,292	0,279	0,266
0,8	0,253	0,24	0,228	0,215	0,202	0,189	0,176	0,164	0,151	0,138
0,9	0,126	0,113	0,1	0,088	0,075	0,063	0,05	0,038	0,025	0,013

Pour l'écart réduit, si  $\alpha = 5\%$  alors  $\varepsilon_t = 1,96$  ++++

Exemple :

Soient 2 populations : la première où les enfants vont à la **crèche** et la deuxième où ils restent à la **maison**.

On cherche à savoir si le mode de garde (crèche ou domicile) modifie le **risque de rhinopharyngite** des enfants.

On fait une étude sur 2 groupes de 200 enfants :

Crèche       $n_A=200$       Nb rhino = 130

Domicile     $n_B=200$       Nb rhino = 96

Le mode de garde influe-t-il sur le risque d'avoir une rhinopharyngite ?

- $H_0$  : pas de différence entre les 2 modes de garde vis-à-vis des rhinopharyngites  
 $H_1$  : différence entre les 2 modes de garde
- Caractère 1 : garde en crèche ou à domicile = qualitatif  
Caractère 2 : avoir une rhinopharyngite ou non = qualitatif
- $\alpha = 5\%$  défini à priori donc  $\varepsilon_{théorique} = 1,96$  (on lit sur la table)
- On calcule le paramètre (donné dans l'énoncé)  $\varepsilon_{calculé} = 3,4$
- $3,4 > 1,96$  : on **rejette  $H_0$**  donc on accepte  $H_1$ 
  - Au niveau de l'échantillon, on peut en conclure que le risque de rhinopharyngites est supérieur chez les enfants gardés en crèche
  - On ne peut pas généraliser cette conclusion au niveau de tous les enfants car il n'y a pas eu TAS

## 2. Le test du $\chi^2$ (tout effectif)

Le paramètre Z est donné ici par le  $\chi^2$ . On va donc comparer :

- $\chi^2_t$  : donné par la **table du  $\chi^2$**  en fonction d'  $\alpha$  et du nombre de ddl
- $\chi^2_c = \sum \frac{(o_i - c_i)^2}{c_i}$

Cette formule permet de comparer les **chiffres calculés C** qui forment le modèle théorique, aux **chiffres observés O**

$$\text{ddl} = (\text{nb lignes} - 1) (\text{nb colonnes} - 1)$$

Le test du  $\chi^2$  permet de **prendre en compte tous les cas de figure** et pas seulement deux %.

### Comment trouver $\chi^2_t$ sur la table du $\chi^2$ ?

On cherche  $\chi^2_t$  en fonction du risque  $\alpha$  et du nb de ddl.

Sur les lignes on trouve le **nombre de ddl** et sur la colonne le **risque  $\alpha$** .  $\chi^2_t$  se trouve à l'intersection des deux.

$$\text{nb ddl} = (\text{nb lignes} - 1) (\text{nb colonnes} - 1)$$

on ne compte pas les lignes et colonnes « total »

Ici par exemple, on cherche  $\chi^2_t$  pour  $\alpha=5\%$  et  $\text{ddl}=1$

$\alpha$										
ddl	0,9	0,5	0,3	0,2	0,1	0,05	0,02	0,01	0,001	
1	0,016	0,455	1,074	1,642	2,706	<b>3,841</b>	5,412	6,635	10,827	
2	0,211	1,386	2,408	3,219	4,605	5,991	7,824	9,21	13,815	
3	0,584	2,366	3,665	4,642	6,251	7,815	9,837	11,345	16,266	
4	1,064	3,357	4,878	5,989	7,779	9,488	11,668	13,277	18,467	
5	1,61	4,351	6,064	7,289	9,236	11,07	13,388	15,086	20,515	
6	2,204	5,348	7,231	8,558	10,645	12,592	15,033	16,812	22,457	
7	2,833	6,346	8,383	9,803	12,017	14,067	16,622	18,475	24,322	
8	3,49	7,344	9,524	11,03	13,362	15,507	18,168	20,09	26,125	
9	4,168	8,343	10,656	12,242	14,684	16,919	19,679	21,666	27,877	
10	4,865	9,342	11,781	13,442	15,987	18,307	21,161	23,209	29,588	
11	5,578	10,341	12,899	14,631	17,275	19,675	22,618	24,725	31,264	
12	6,304	11,34	14,011	15,812	18,549	21,026	24,054	26,217	32,909	
13	7,042	12,34	15,119	16,985	19,812	22,362	25,472	27,688	34,528	
14	7,79	13,339	16,222	18,151	21,064	23,685	26,873	29,141	36,123	
15	8,547	14,339	17,322	19,311	22,307	24,996	28,259	30,578	37,697	
16	9,312	15,338	18,418	20,465	23,542	26,296	29,633	32	39,252	
17	10,085	16,338	19,511	21,615	24,769	27,587	30,995	33,409	40,79	

Exemple :

On cherche à savoir si l'**exposition professionnelle au benzène** peut entrainer une **leucémie**. On lance une étude dans une grande entreprise, on dénombre les salariés exposés au benzène, et ceux qui ne le sont pas. Au bout de 12 ans, on fait le bilan des leucémies apparues.

	Leucémies	Non leucémies	Total
Expo	15	485	500
Non expo	20	980	1000
Total	35	1465	1500

Existe-t-il une relation entre exposition au benzène et leucémies ?

1.  $H_0$  : pas de lien entre l'exposition au benzène et les leucémies
2. Variable 1 : leucémie ou non = qualitatif  
Variable 2 : exposition au benzène ou non = qualitatif
3.  $\alpha = 5\%$  défini à priori  
nb ddl =  $(2 - 1) (2 - 1) = 1$  }  $\chi^2_t = 3,841$
4. On calcule le paramètre (donné dans l'énoncé)  $\chi^2_c = 1,42$
5.  $1,42 < 3,84$  donc on **accepte  $H_0$**

## III. LIAISONS ENTRE CARACTERES QUALITATIFS ET QUANTITATIFS

*Problématique : En moyenne la taille des individus d'une population A coïncide-t-elle avec la taille des individus d'une population B ?*

### 1. Comparaison de moyenne : $n_1$ et $n_2 > 30$ "Grands échantillons"

On utilise : La **table de l'écart réduit** ( $\varepsilon < 1,96$  avec  $\alpha < 5\%$ ).  $\varepsilon = \frac{m_1 - m_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$

→ Si  $\varepsilon_{calculé} > \varepsilon_{théorique}$  Rejet de  $H_0$

**Exemple :** On teste un antiviral diminuant le nombre de jours de symptômes cliniques chez des patients infectés par le virus de la grippe. Soit 100 sujets non traités, atteints de grippe. Le nombre moyen de jours avec symptômes est  $m_1 = 4,74$  jours et l'écart-type :  $s_1 = 1$ . Soit 100 autres sujets traités avec l'antiviral et atteints de grippe le nombre moyen de jours avec symptômes est  $m_2 = 4,2$  jours et l'écart-type est  $s_2 = 1,7$ .

→ On fera ici un test de comparaison de moyenne qui nous permettra de répondre à la question : Peut-on accepter ou rejeter  $H_0$  ?

## 2. Série numérique : t de Student $n_1$ ou $n_2 < 30$ "Petits échantillons"

On utilise : La **table t de Student**, avec  $(n_1-1)+(n_2-1)$ ddl.  $t = \frac{m_1 - m_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$

→ Si  $t_{calculé} > t_{théorique}$  Rejet de  $H_0$

**Remarque :** Ici on calcule un écart-type « s » sur les deux échantillons au lieu d'utiliser  $s_1$  et  $s_2$  (comme pour les comparaisons de moyennes). En effet ici  $s_1$  et  $s_2$  sont moins significatifs que pour les tests de comparaison de moyenne (mais « osef » selon le prof).

$$s = \sqrt{\frac{\sum(x_i - m_1)^2 + \sum(x_j - m_2)^2}{(n_1-1) + (n_2-1)}}$$

**Précision sur les ddl :**

On prend une série de 8 valeurs donc  $n=8$  :

Toutes les valeurs :	2	3	5	12	10	4	7	8	Total=51
Avec 1 valeur manquante :	2	3	5	12	10		7	8	Total=47
Avec 2 valeurs manquantes :	2	3		12	10		7	8	Total=42

S'il nous **manque une donnée** : mais que **nous connaissons le total=51** et le total sans la valeur on peut retrouver la valeur manquante :  $51-47=4$

→ Avec  $n-1$  valeurs on peut calculer la valeur manquante à partir du total.

S'il nous manque deux valeurs et qu'on connaît seulement le total=51 et le total sans les deux valeurs on ne peut pas retrouver les deux valeurs manquantes :  $51-42=9$ . Or on ne sait pas quelles sont les deux valeurs : on connaît seulement leur somme 9.

→ Avec  $n-2$  valeur il est impossible de trouver les deux valeurs manquantes.

Ici on a donc  $ddl = n-1 = 7$

Dans le **test t de Student** on a 2 série à comparer d'où :

$$ddl = (n_1 - 1) + (n_2 - 1)$$

**Exemple :**

Soient un groupe de 15 femmes obèses, et un autre groupe de 12 femmes de poids normal. On a mesuré le taux de corticoïdes sanguins moyens à l'intérieur de ces 2 groupes. Pour le groupe 1 :  $n_1 = 15$  ;  $m_1 = 6,3$  ;  $s_1 = 1,8$  et pour le groupe 2 :  $n_2 = 12$  ;  $m_2 = 4,5$  ;  $s_2 = 1,6$ .

L'obésité a-t-elle une influence sur le taux de corticoïdes ?

**Méthode :**

1. **On pose  $H_0$**  :  $m_1$  et  $m_2$  ne sont pas différentes dans ces 2 groupes.
2. **Type de caractères étudiés** : Relation entre caractères qualitatifs (obèses et non obèses), et quantitatifs (valeurs de dosages sanguins, valeurs moyennes).
3. **Taille de l'échantillon** :  $n_1=15$  et  $n_2=12$  les deux sont  $< 30$  c'est donc un petit échantillon. **Choix du test** : test t de Student.
4. **Ecart-type** : Ici on doit calculer l'écart-type commun aux deux groupes car on est en t de Student (en pratique dans les QRU vous n'aurez pas à calculer ce « s », comme dit le Professeur Bénoliel « pas de calculatrice : pas de calcul ! »  $< 3$ ).

$$s^2 = 2,53 ; ddl = (15-1) + (12-1) = 25 ; t = 2,92$$

On cherche donc t dans la table t de Student :

Ddl/ $\alpha$	0,9	0,5	0,3	0,2	0,1	0,05	0,02	0,01	0,001
25	0,127	0,684	1,058	1,316	1,708	2,06	2,485	2,787	3,725

$$t = 2,92 \in [2,787 ; 3,725]$$

→ On a  $\alpha < 1\%$  (après lecture dans la table) donc on peut rejeter  $H_0$  et conclure à une relation entre obésité et augmentation du taux de corticoïdes au niveau de ces échantillons (on ne peut conclure que sur les échantillons et non pas sur la population générale car il n'y a pas eu de TAS).

### 3. Séries appariées ou méthode des couples

On utilise la méthode des couples lorsqu'on étudie la liaison entre deux variables **qualitatives et quantitatives** dans 2 échantillons non indépendants.

**Série indépendante :** Lorsque les deux groupes comparés sont distincts et indépendants (=sans lien).

Exemple : On tire au sort un groupe 1 puis un groupe 2, G1 consommera un placebo et G2 le nouveau médicament à tester.

**Série appariée :** Lorsque les deux groupes comparés ne sont pas distincts et indépendants (=liés).

Exemple : On compare les résultats avant traitement puis après traitement : c'est donc une observation sur le même groupe : les groupe avant traitements et après traitements sont identiques. Ils ne sont pas indépendants car ils forment un seul et même groupe.

#### Test utilisés :

Si  $n > 30$   $\varepsilon = m_d / \sqrt{\frac{s^2}{n}}$  avec un test de **comparaison de moyenne**.

Si  $n < 30$   $t = m_d / \sqrt{\frac{s^2}{n}}$  avec un **test t de Student**.

#### Exemple :

On veut comparer deux méthodes de dosage de la glycémie. On dispose de n patients, auxquels on prélève 2 tubes de sang. On dose la glycémie dans chacun de ces tubes par une méthode différente. On souhaite comparer les valeurs moyennes de ces 2 séries de n résultats. La question posée est :

Ces 2 méthodes de dosage fournissent-elles des résultats identiques ?

On calcule si  $n > 30$   $\varepsilon = m_d / \sqrt{\frac{s^2}{n}}$  si  $n < 30$   $t = m_d / \sqrt{\frac{s^2}{n}}$

Avec d = différence des résultats pour un même sujet,  $m_d$  = moyenne des d,  $n$  = nb de couples,  $s$  = variance des différences

Puis la méthodologie est identique aux tests déjà vus : on compare cette valeur calculée aux valeurs dans la table adaptée, et la conclusion se fait de la même manière en fixant un risque  $\alpha$ .

#### Autre exemple :

On souhaite évaluer l'intérêt d'une substance S capable de désintoxiquer les fumeurs. On constitue par T.A.S. 2 groupes de 40 fumeurs, un reçoit S, l'autre reçoit

un placebo P. Le traitement dure 2 mois pour les 2 groupes. La consommation de cig/jours (C) est notée avant et après traitement.

1) Quelle est la **première précaution** à prendre ? La consommation est-elle identique dans les 2 groupes ? Les 2 **groupes** doivent être **comparables** vis-à-vis des paramètres susceptibles d'influencer la réponse au traitement (âge, sexe, CSP, conso/jour etc..). Si ce n'est pas le cas, il faut en tenir compte lors des conclusions.

	S (n=40)		P (n=40)	
	$m_1$	$s_1^2$	$m_2$	$s_2^2$
C avant tt	19,5	54,2	16,5	35,6
C après tt	5,4	30,4	3,8	20,1
Variation de C	14,1	9,1	12,7	8,9

Comparaison des conso moyennes avant tt dans les 2 groupes

1.  $H_0$  = les moyennes des consommations sont équivalentes dans les 2 groupes.
2. Etude liaison entre variables qualitatives (S ou P) et quantitatives (nb cig/j) dans échantillons indépendants.
3.  $n > 30 \rightarrow$  Test de comparaison de moyennes
4.  $\varepsilon = 2,00 > 1,96$  (1,96 est le  $\varepsilon$  pour  $\alpha = 5\%$ ).

On rejette  $H_0$ , avec un risque  $\alpha = 5\%$ . Il existe donc une différence significative entre les consommations moyennes des 2 groupes : on fume plus dans le groupe S. Il faudra en tenir compte lors de l'étude de la variation de cette conso avant/après traitement.

2) Dans le groupe Placebo, la conso moyenne après traitement diffère-t-elle de la valeur avant traitement ? **Interpréter le résultat.**

1. Liaison entre variable qualitative (avant / après tt) et quantitative (nb cig/j)
2. Echantillons non indépendants (méthode des couples)
3.  $n > 30 \rightarrow$  Test de comparaison de moyennes
4.  $\varepsilon = 26,9 > 1,96$  au risque  $\alpha = 5\%$ .

On rejette  $H_0$ . Il existe une différence très significative ( $p < 0,001$ ) entre les consommations avant / après tt, dans le groupe P. Il y a sûrement eu un effet psychologique : envie de profiter de l'étude pour arrêter de fumer ?

3) Les 2 groupes diffèrent-ils pour leur consommation moyenne après traitement ?

1.  $H_0$  = les moyennes des consommations sont équivalentes dans les 2 groupes.
2. Liaison entre variables qualitatives (S ou P) et quantitatives (nb cig/j) dans échantillons indépendants.
3.  $n > 30 \rightarrow$  Test de comparaison de moyennes
4.  $\varepsilon = 1,42 < 1,96$

On accepte  $H_0$  : il n'existe pas de différence significative entre les 2 groupes pour la consommation après tt.



4) Les 2 groupes diffèrent-ils pour la variation de conso avant/après traitement ?  
Il faut comparer les variations avant / après tt dans les 2 groupes afin de prouver l'intérêt de la substance S

1.  $H_0$  : Il n'existe pas de différence entre les variations de consommation dans les 2 groupes.

2. Etude liaison entre variables qualitatives (S ou P) et quantitatives (nb cig/j) dans 2 échantillons indépendants.

3.  $n > 30 \rightarrow$  Test de comparaison de moyennes

4.  $\varepsilon = 2,09 > 1,96$  au risque de 5%

On rejette  $H_0$  : il existe une différence significative entre les variations de conso dans les 2 groupes ( $p < 5\%$ ). Conclusion : efficacité de S. Il y avait eu TAS, donc résultat généralisable.

Conclusion : Pas de différence après tt dans chaque groupe (cf Question 3). Mais le Groupe S fumait plus (cf Question 1)  $\rightarrow$  Efficacité du traitement S.

## IV. LIAISON ENTRE CARACTÈRES QUANTITATIFS

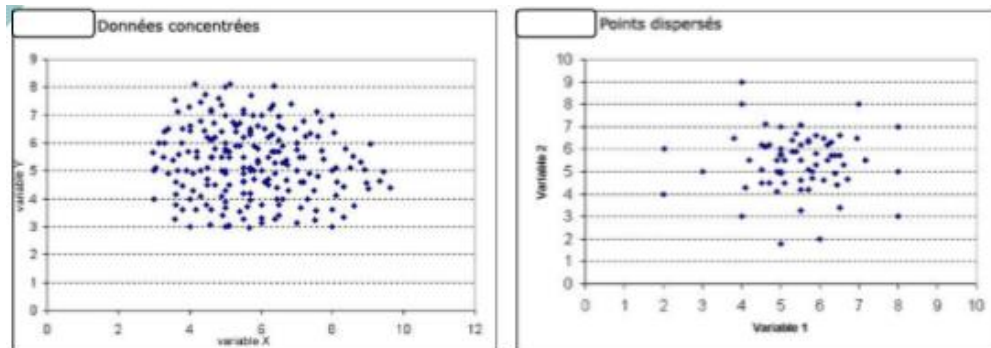
### 1. Corrélation et régression

**Corrélation** = Evaluation de la liaison entre 2 variables quantitatives

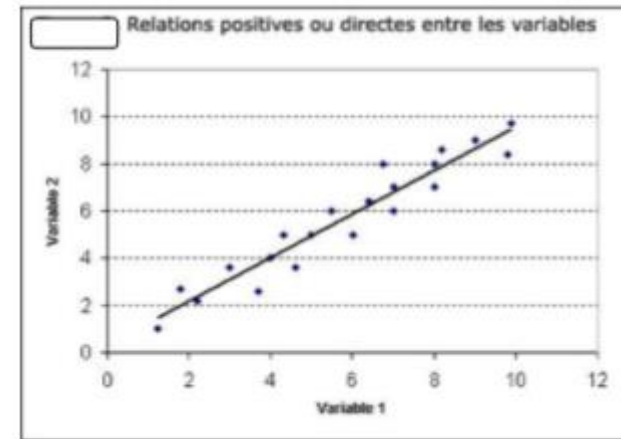
**Régression** = Méthode mathématique expliquant les relations entre variables observées.

### 2. Représentation des données

**Nuage de points :**



**Droite de Régression :** La droite de régression permet de visualiser si une des 2 variables est dépendante de l'autre :



$\rightarrow$  La droite de régression est aussi appelée **droite des moindres carrés** car elle passe au « plus près » de chaque point du graphe. (Dans ce cours, on ne parle que de régression linéaire.)

Une droite de régression peut permettre de prédire certaine valeur de y à partir d'une valeur de x.

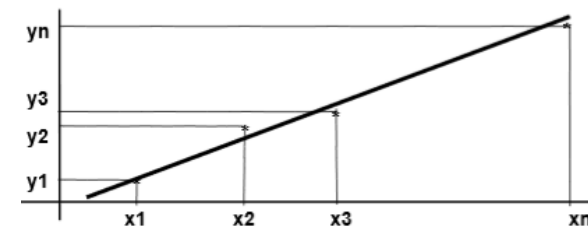
### 3. Etude de la liaison entre caractères quantitatifs

Exemple :

$\rightarrow$  La capacité respiratoire est-elle dépendante de la consommation de cigarettes ?

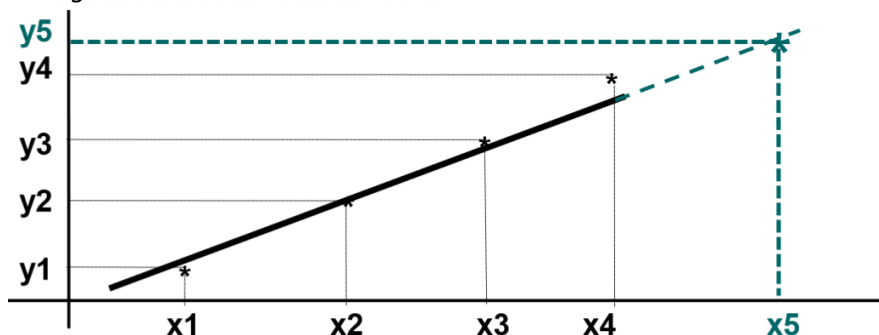
$\rightarrow$  Le poids des bébés à la naissance est-il lié à l'âge de la mère ?

Si **x et y liées** alors  $\rightarrow y=f(x)$  on a une **droite de régression de y en x** !



Remarque : On dit que y peut être peut-être « expliqué » en fonction de x.

**Prédiction :** La droite de régression permet de prédire des valeurs de y pour un certains x. Il suffit pour cela de prolonger la droite ! *ici on arrive à prédire le point en turquoise grâce à la droite.*



### Coefficient de corrélation : Pente de la droite

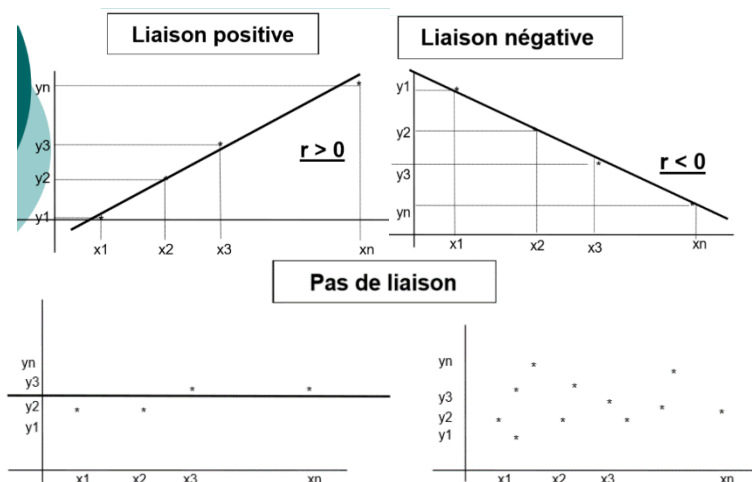
On utilise : La **Table du coefficient de corrélation**, avec **n-2 ddl**.

$$r = \frac{\sum (xi - mx)(yi - my)}{\sqrt{\sum (xi - mx)^2 \sum (yi - my)^2}} \quad r = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\sqrt{(\sum x^2 - \frac{(\sum x)^2}{n})(\sum y^2 - \frac{(\sum y)^2}{n})}} \quad (\text{osef})$$

→ Si  $r > 0$  : liaison positive donc x et y varient **dans le même sens**.

→ Si  $r < 0$  : liaison négative donc x et y varient **en sens inverse**.

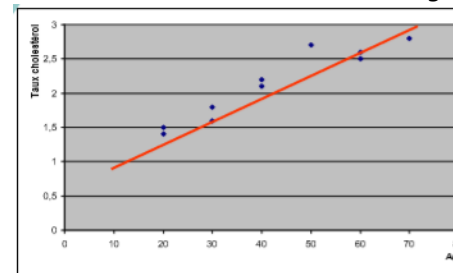
**!/ r est toujours inférieur à 1 !/**



**Exemple :** Sur un échantillon de 10 sujets d'âges différents, on recueille les données suivantes : âge (années) et concentration de cholestérol dans le sang (g/L).

X âges	30	60	40	20	50	30	40	20		70	60
Y cholestérol	1,6	2,5	2,2	1,4	2,7	1,8	2,1	1,5		2,8	2,6

Le taux de cholestérol est-il lié à l'âge ?



Existe-t-il un lien entre ces 2 séries de données ? Ou bien s'agit-il de 2 séries totalement indépendantes ?

1. H0 = Le taux de cholestérol est indépendant de l'âge

H1 = Le taux de cholestérol est lié à l'âge

2. 2 variables quantitatives > **Test du coefficient de corrélation**

3. r calculé = 0,955 > r théorique = 0,76 avec 10-2 = 8 ddl on voit dans la tables : ( $\alpha = 1\%$ )

ddl	$\alpha$			
	0,1	0,05	0,02	0,01
8	0,5494	<b>0,6319</b>	0,7155	<b>0,7646</b>

Conclusion : Rejet de H0, Il existe une relation significative ( $\alpha = 1\%$ ) entre l'âge et le taux de cholestérol : plus l'âge augmente, plus le taux de cholestérol augmente. Cependant le résultat est non généralisable (pas de TAS).

**Remarque :**

**!/ Corrélation  $\neq$  causalité !/**

Si on établit corrélation entre deux variables cela veut dire qu'il existe un lien entre les deux. *Ex : l'âge et le cholestérol sont liés (on ne dit pas que l'un cause l'autre).*

Si on établit une causalité entre deux variables cela veut dire que l'une est la cause de l'autre. *Ex : l'âge cause le cholestérol.*

**Exemple 2 :**

On teste un traitement favorisant la baisse de tension artérielle.



On cherche à savoir si l'effet de ce traitement est lié à l'âge des patients. Dans ce but, on effectue 2 prises de tension : une avant traitement et une autre 1h après le traitement. On note la différence entre les 2 valeurs.

Soit X la série des âges des patients et soit Y la série des différences de T.A avant – après traitement : la question posée peut se traduire par  $Y = f(X)$  ?

1.  $H_0$  : les 2 séries x et y sont indépendantes, et il n'existe pas de relation entre elles.

2. 2 variables quantitatives > **Test du coefficient de corrélation**

3. r calculé =  $|-0,83| > r$  théorique =  $|0,76|$  au risque de 1% :

NB : on compare les r théoriques et calculés en valeur absolue !

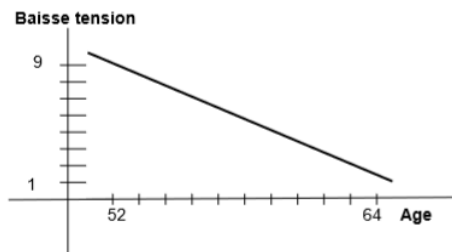
x = âge	y = diff	xy
64	1	64
52	9	468
56	9	504
60	6	360
62	2	124
58	4	232
57	3	171
61	1	61
58	6	348
55	9	495

$$\sum x = 583 \quad \sum y = 50 \quad \sum xy = 2827$$

ddl	$\alpha$			
	0,1	0,05	0,02	0,01
8	0,5494	0,6319	0,7155	0,7646

**Conclusion** : Nous pouvons rejeter  $H_0$ . Il existe une relation entre x et y (elles sont « corrélées »), avec  $p < 1\%$

r calculé -0,83 : le signe – indique que plus les valeurs d'une série augmentent, plus les valeurs de l'autre série diminuent. Elles varient en sens inverse.



La pente de la droite est négative. L'âge augmente, la baisse de tension est moins forte. Le traitement est moins efficace avec l'âge.

## V. TESTS NON PARAMETRIQUES

Ces définitions ne sont pas dans le cours c'est pour vous aider à comprendre ☺

**Test paramétrique** : test avec modèle à forte contrainte, difficile à réaliser sur de petits effectifs. (ex : t de Student ; Chi-2 ; comparaison de moyenne...)

**Test non paramétrique** : test dont le modèle ne précise pas les conditions que doivent remplir les paramètres de la population dont a été extrait l'échantillon. (ex : U Mann et Whitney ; r' de Spearman ; Wilcoxon)

⚠ On utilise obligatoirement un test non paramétrique si les effectifs sont trop faibles ( $4 < X < 12$ ) ⚠

➔ Ils sont utilisés pour des variables quantitatives lorsque les effectifs sont trop faibles (les populations ne sont pas distribuées normalement).

➔ Ces tests présentent une excellente robustesse.

### 1. Les différents tests non paramétriques

	Test	
	Paramétriques	Non paramétriques
Comparaison de 2 échantillons indépendants	T de Student Comparaison de moyenne	Mann-Whitney
Comparaison de 2 échantillons appariés	T de Student Comparaison de moyenne	Wilcoxon
Test de corrélation	Coefficient r	Coefficient r' de Spearman

**Principe** :

Pour réaliser ces tests il est nécessaire de transformer les données quantitatives en données de mesures ordinales (les rangs). On prend les données quantitatives et on les ordonne.

**Exemple** :

Si on étudie la variable quantitative « âge » on va ranger les âges des patients du plus jeune au plus âgé :

Ages	14	17	28	30
Rangs des âges	1	2	3	4

## 2. U Mann et Whitney : $4 < n < 12$

On utilise : La **table de U Mann et Whitney** ( $\alpha < 5\%$ )

→ Si  $U_{calculé} > U_{théorique}$  **Acceptation** de  $H_0$

Ces paramètres **résumant les données et traduisent leur imbrication**. Si les données sont très imbriquées il n'y a pas de différence significative entre les 2 groupes. Si les données ne sont pas ou peu imbriquées il y a une différence significative entre les 2 groupes. On compare **le plus petit des deux U** (soit  $U_{AB}$  soit  $U_{BA}$ ) à une valeur théorique lue dans la table (intersection de la ligne  $n_A - n_B$  et colonne  $n$ ). Cet U théorique est la limite max au-delà de laquelle l'imbrication est considérée comme importante. C'est-à-dire que lorsque le U calculé est plus important que le U théorique (qui est cette limite) on peut conclure qu'il y a imbrication et donc **accepter  $H_0$**  (si les données sont imbriquées c'est qu'elles sont issues d'un même ensemble).

🔍 Le test de Mann-Whitney permet de **tester si 2 groupes indépendants sont extraits d'une population unique** (2de possibilité) **ou non** (1ère possibilité).

Exemple :

1) On dispose de 2 groupes (groupe A de 6 malades, et groupe B de 5 sujets non malades). On dose une certaine hormone dans le sang de ces 11 sujets. Y a-t-il une différence significative entre ces 2 groupes du point de vue de cette hormone ?

Gr A 11 ; 21 ; 25 ; 52 ; 71 ; 79

Gr B 22 ; 43 ; 72 ; 92 ; 116

Nous constatons :

🔍 **Peu** de valeurs à comparer

🔍 Le test « Ressemble » à la comparaison de moyennes car on compare une variable quantitative à une variable qualitative.

🔍 Soit les dosages diffèrent en fonction du groupe malades/non malades, soit ils sont tous équivalents.

1.  $H_0$  : pas de différence significative entre les 2 groupes pour cette hormone.

2. Etude d'une liaison éventuelle entre données qualitatives (malades ou non), et données quantitatives (valeurs des dosages).

3. Comparaison de moyennes, effectif faible → U de Mann & Whitney

4. On range toutes les valeurs A ou B par ordre croissant :

	11	21	22	25	43	52	71	72	79	92	116
	A	A	B	A	B	A	A	B	A	B	B
$U_{AB}$	0	0	2	1	3	2	2	5	3	6	6

On va calculer le **paramètre  $U_{BA}$**  : pour chaque membre du groupe A on cumule le nombre de membres du groupe B qui sont classés avant lui. On peut calculer de même  $U_{AB}$ .

$$U_{BA} = 0 + 0 + 1 + 2 + 2 + 3 = 8$$

$$U_{AB} = 2 + 3 + 5 + 6 + 6 + 6 = 22$$

Remarque :  $U_{AB} + U_{BA} = n_A \times n_B = 6 \times 5 = 30$

On compare le plus petit des deux U : soit  $U_{BA}$  à une valeur théorique lue dans la table (intersection de la ligne  $n_A - n_B = 1$  et colonne  $n = 5$ ).

		$n_1$									
2	$n_2 - n_1$	1	2	3	4	5	6	7	8	9	10
	0	-	-	-	0	2	5	8	13	17	23
	1	-	-	-	1	3	6	10	15	20	26

→ La valeur **U théorique = 3**

$U_{BA} > 3$  donc l'imbrication des 2 groupes est considérée comme importante = les données sont issues d'un même ensemble = il n'y a pas donc de différence significative entre les 2 groupes = **on accepte  $H_0$**

Méthode :

1. **1ère possibilité** : 2 groupes indépendants sont extraits d'une population unique.
2. **2ème possibilité** : 2 groupes indépendants viennent d'une même population.
2. **Type de caractères étudiés** : Relation entre caractères qualitatifs (Groupe A ou B), et quantitatifs (dosage de l'hormone).
3. **Taille de l'échantillon** :  $n_1=6$  et  $n_2=9$  les deux sont  $< 12$  Il faut donc un test non paramétrique. **Choix du test** : test U Mann et Whitney.
4. **Calcul du paramètre et consultation de la table**

2) Analogie : lors d'une course de vélo on a :

🔍 Toute l'équipe blanche en tête de peloton : Pourquoi cette équipe est meilleure ? L'équipe est-elle dopée ? (Dans ce cas elle ne correspondrait pas à la population générale des auteurs cyclistes non dopés).

🔍 Des cyclistes de toutes les couleurs sont intercalés : Ont-ils tous le même niveau ?

### 3. $r'$ de Spearman (corrélation) :

On utilise : La **table du  $r$  de Spearman** ( $r'=0,89$  avec  $\alpha < 5\%$ ,  $r'=1$  avec  $\alpha < 1\%$ )

$$r' = 1 - \frac{6 \sum di^2}{n(n^2-1)}$$

→ Si  $r_{calculé} = r_{théorique}$  Rejet de  $H_0$

Exemple :

On a recensé pour 6 étudiants les notes obtenues au concours PACES en Biostatistique, et le classement final à ce même concours. On cherche à établir s'il existe une relation entre cette note et le classement final.

Rappel : Variable Classement = variable « pseudo quantitative ».

1.  $H_0$  : Il n'y a pas de lien entre ces 2 séries de valeurs numériques. Il s'agit de 2 séries indépendantes.

2. Etude d'une liaison éventuelle entre données quantitatives (classement et note en Biostat) : coeff de corrélation mais avec de faibles effectifs →  $r'$  Spearman

3. On calcule le  $r_4$  (vous n'aurez pas à le faire au cc) :  $r'_{calculé} = -1$

On regarde dans la table du  $r'$  de Spearman l'intersection entre l'effectif et  $\alpha$  :

$\alpha$	0.50	0.20	0.10	0.05	0.02	0.01
n						
4	0.600	1.000	1.000			
5	0.500	0.800	0.900	1.000	1.000	
6	0.371	0.657	0.829	0.886	0.943	1.000
7	0.321	0.571	0.714	0.786	0.893	0.929
8	0.310	0.524	0.643	0.738	0.833	0.881
9	0.267	0.483	0.600	0.700	0.783	0.833
10	0.248	0.455	0.564	0.648	0.745	0.794
11	0.236	0.427	0.536	0.618	0.709	0.755

$r'_{calculé} = |-1| = |r'_{théorique}|$

On repousse donc  $H_0$  ( $p < 1\%$ ) On met en évidence un lien très significatif entre ces 2 séries. Il s'agit donc de 2 séries corrélées :

→ Plus la note de Biostat est élevée, plus petit est le rang de classement (d'où le signe – pour  $r'$ ).

X Biostat	Y Classement
12,4	210
4,9	555
18,1	6
5,4	445
19,4	5
16	14

## VI. ♥♥♥+++RECAP DES TESTS+++♥♥♥

Effectif	Données Quantitatives	Données Qualitatives	Données Qualitatives - Quantitatives
$\geq 30$	Coeff de corrélation $r$ $r'$ Spearman	Comp % ou $\chi^2$	Comp moyennes $t$ Student ou U Mann & Withney
$< 30$ & $\geq 12$	Coeff de corrélation $r$ $r'$ Spearman	Comp % ou $\chi^2$	$t$ Student ou U Mann & Withney
$> 4$ & $< 12$	$r'$ de Spearman	Comp % ou $\chi^2$	U Mann & Withney

NB : la flèche qui est orientée vers le haut indique que l'on peut utiliser le test pour des effectifs supérieurs. Attention l'inverse n'est pas vrai on ne peut pas utiliser de test pour des effectifs inférieurs !

NB bis : Le test écrit en vert est celui que l'on préférera utiliser.