

DULCLAUDIAX

STATISTIQUES DEDUCTIVES





Sommaire

I. Test d'hypothèse

II. Lien entre 2 variables qualitatives

III. Lien entre variable qualitative et quantitative

IV. Lien entre 2 variables quantitatives

V. Tests non paramétriques

TEST D'HYPOTHÈSE



Statistiques déductives → tirer des conclusions à partir des observations

Hypothèse nulle (H_0)

Pas de différence entre les 2 groupes

Les fluctuations observées sont dues au hasard

Hypothèse alternative (H_1)

Différence significative entre les 2 groupes

Les fluctuations observées ne sont pas dues au hasard



ETAPES D'UN TEST

- Définir H_0 et H_1
- Choisir le test (en fonction des données)
- Fixer le risque α (souvent 5%)

Théorie

- Recueillir les données
- Calculer le paramètre Z
- Comparer Z_c à Z_t

Pratique

- Accepter ou rejeter H_0
- Fixer le risque d'erreur réel β
- Interpréter les résultats : statistiquement et cliniquement

Conclusion





NOTION DE RISQUE

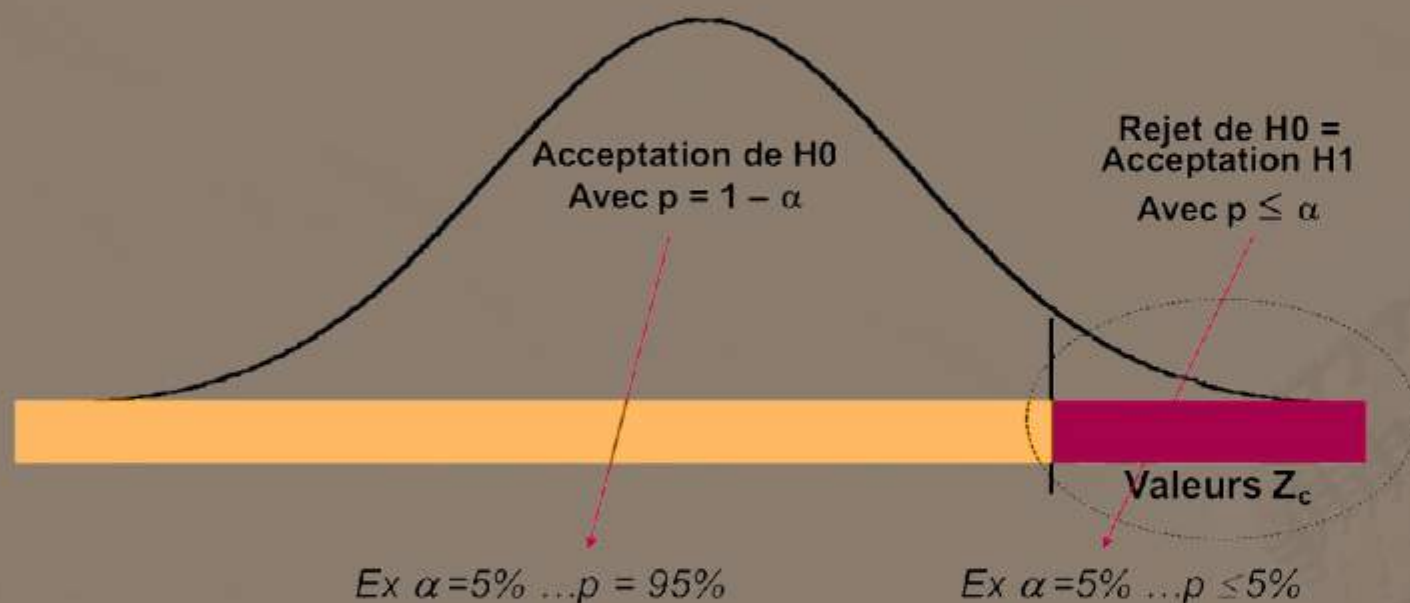
	Rejet H0	Non rejet H0
H0 vraie	α	$1-\alpha$
H1 vraie	$1-\beta$	β

Risque de première espèce Risque α	Risque de seconde espèce Risque β
Probabilité de rejeter H0 si H0 est vraie	Probabilité d'accepter H0 si H0 est fausse
Ce risque est maîtrisé	Ce risque est négligé Il peut être très élevé (en général $\beta=20\%$)
Fixé à l'avance	Fixé à postériori
La puissance du test vaut $1-\beta$: probabilité de rejeter H0 avec H1 vraie	

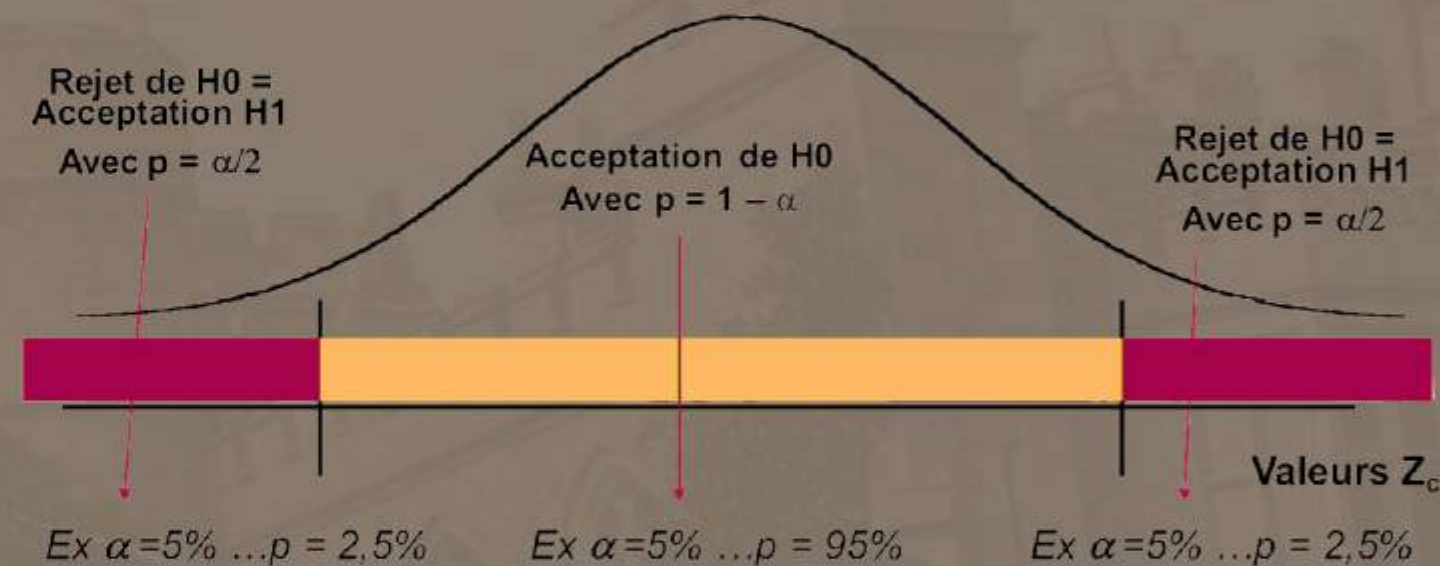




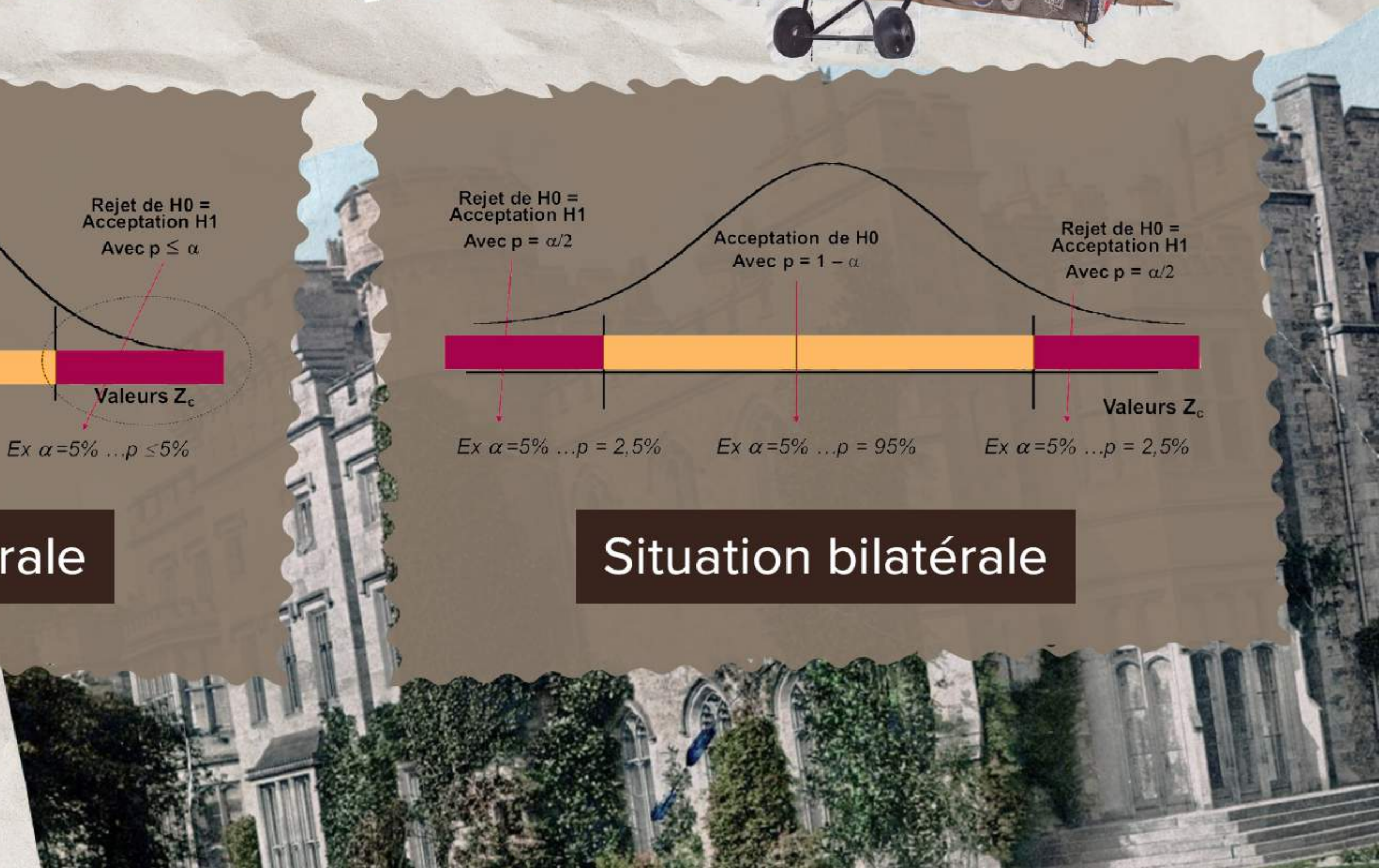
INTERPRÉTATION GRAPHIQUE



Situation unilatérale



Situation bilatérale



Lien entre 2 variables qualitatives

TEST DE COMPARAISON DE POURCENTAGES

Tout effectif

Paramètre Z → écart-réduit ε

ε_t vient de la table de l'écart-réduit

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

Si $\varepsilon_c > \varepsilon_t \rightarrow$ **rejet de H0**

Méthodologie

Comment trouver le paramètre théorique ?

Table de l'écart réduit

		α								
		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	1,96	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

Table pour les petites valeurs de la probabilité

0,001	0,000 1	0,000 01	0,000 001	0,000 000 1	0,000 000 01	0,000 000 001
3,2905	3,8905	4,41717	4,89164	5,32672	5,73073	6,10941

Exemple

Soient 2 groupes de 200 enfants :

→ Crèche : 200 enfants, 130 atteints de rhinopharyngite

→ Maison : 200 enfants, 96 atteints de rhinopharyngite

Le mode de garde influe-t-il sur le risque de rhinopharyngite ?



Exemple

	Crèche	Domicile
Sain	70	104
Malade	130	96

→ H0 : pas de différence entre les 2 modes de garde vis-à-vis du développement de rhinos

→ H1 : il y a une différence entre les 2 modes de garde

Variables qualitatives → **test de comparaison de pourcentages**

- $Z_c = 3,4$
- $Z_t = 1,96$

$Z_c > Z_t$ donc on REJETTE H0 au seuil 5%

Ccl : le risque de rhinopharyngite est supérieur chez les enfants gardés en crèche que chez les enfants gardés à domicile.

Lien entre 2 variables qualitatives

Tout effectif

TEST DU KHI²

Test préféré si tableau de données a + de 2 lignes (ou colonnes)

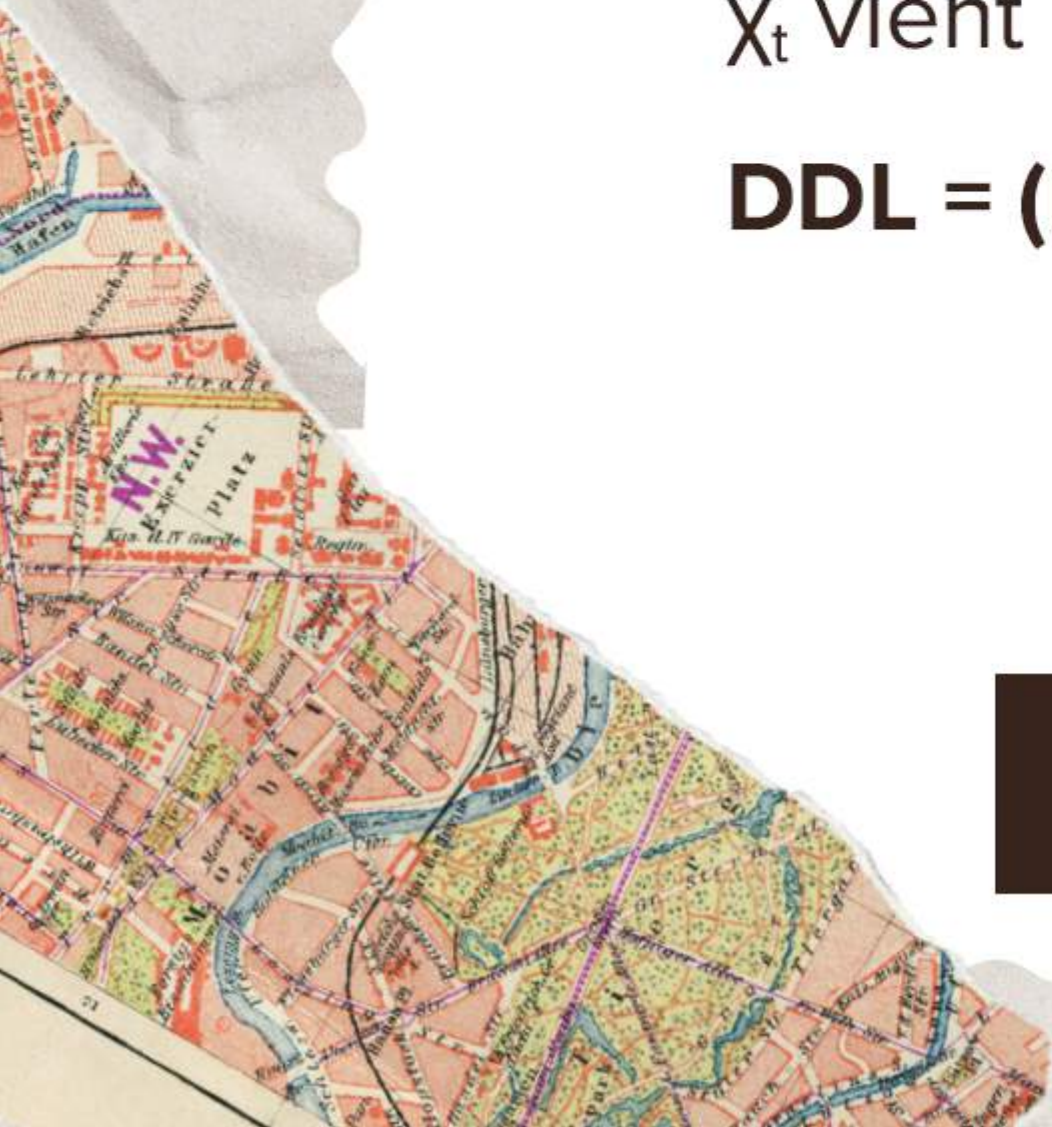
Paramètre Z → χ^2

χ_t vient de la table du χ^2

DDL = (nb de lignes - 1) * (nb de colonnes - 1)

$$\chi^2 = \sum \frac{(o_i - c_i)^2}{c_i}$$

Si $\chi^2_c > \chi^2_t \rightarrow$ rejet de H0



Exemple

Existe-t-il un lien entre l'exposition au benzène et la survenue d'une leucémie ?

	Leucémie	Non leucémie	Total
Expo	15	485	500
Non expo	20	980	1000
Total	35	1465	1500

Exemple

Hypothèses :

H0 : il n'existe pas de lien entre l'exposition au benzène et les leucémies

Choix du test :

Variable 1 : leucémie ou non → qualitatif

Variable 2 : exposé au benzène ou non → qualitatif

→ **Test du χ_2**



Calcul du χ^2

Valeurs observées :

	Leucémie	Non leucémie	Total
Expo	15	485	500
Non expo	20	980	1000
Total	35	1465	1500

Valeurs calculées:

	Leucémie	Non leucémie	Total (environ)
Expo	11,65	488,3	500
Non expo	23,35	976,7	1000
Total	35	1465	1500

$$\chi^2 = \frac{(15 - 11,65)^2}{11,65} + \frac{(20 - 23,35)^2}{23,35} + \frac{(485 - 488,3)^2}{488,3} + \frac{(980 - 976,7)^2}{976,7} = 1,42$$

$$DDL = (\text{nombre de lignes} - 1) * (\text{nombre de colonnes} - 1) = (2-1) * (2-1) = 1$$

Conclusion

$\chi^2_c < \chi^2_t$ donc on accepte H_0 au seuil 0.05

Ccl : Il n'existe pas de relation entre l'exposition au benzène et les leucémies



Lien entre variables qualitatives et quantitatives

TEST DE COMPARAISON DE MOYENNES

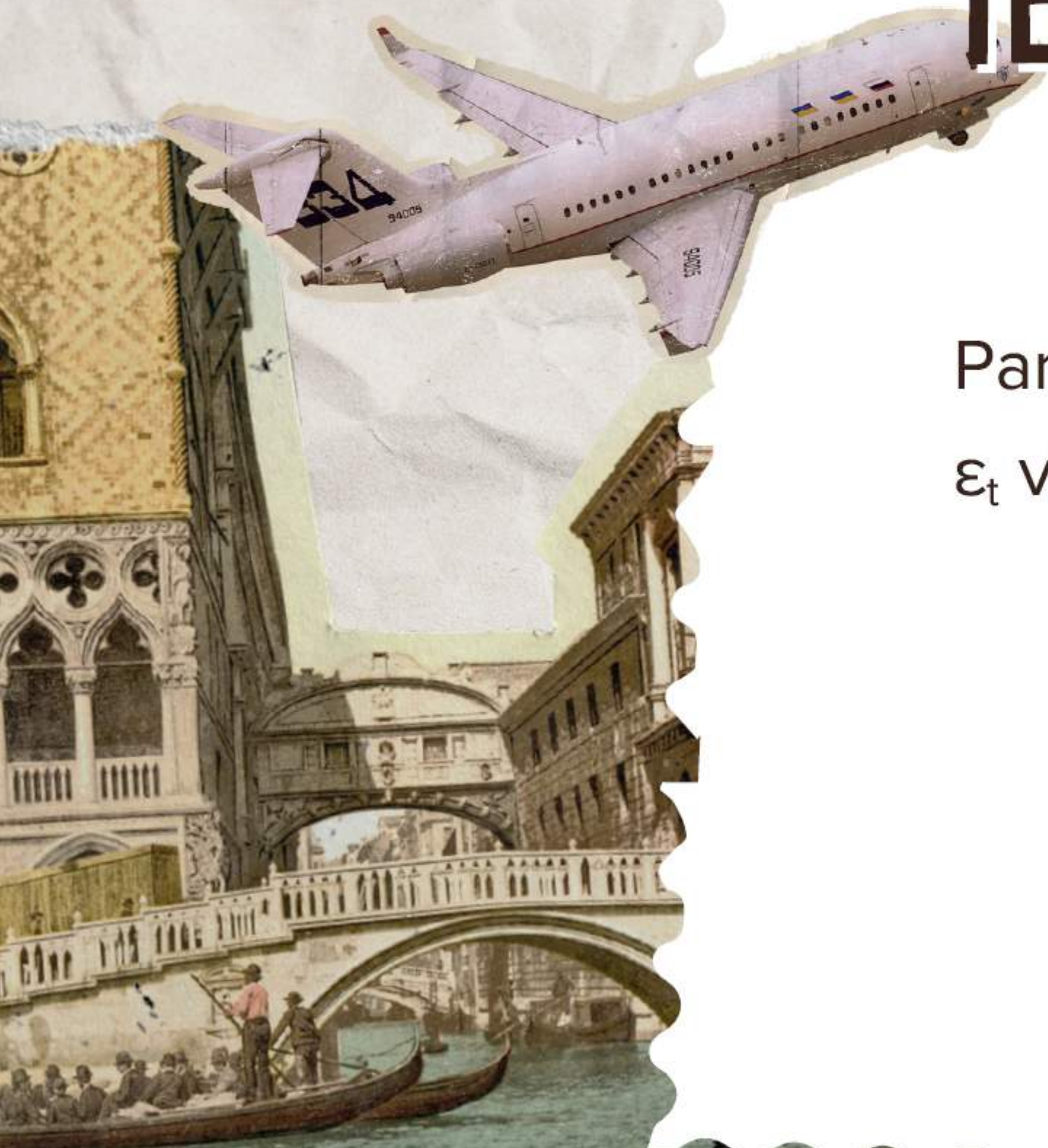
Grands échantillons
(n_1 et $n_2 > 30$)

Paramètre Z \rightarrow écart-réduit ϵ

ϵ_t vient de la table de l'écart-réduit

$$\epsilon_c = \frac{m_1 - m_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Si $\epsilon_c > \epsilon_t \rightarrow$ **rejet de H0**



Exemple

On cherche à **comparer les taux de T3 libre chez les femmes prenant un contraceptif oral (c.o.) et chez celles qui n'en prennent pas.** Après un tirage au sort, on obtient :

Femmes sans c.o. : $n_1 = 50$; $m_1 = 2$ nmol ; $s_1 = 0,35$ nmol

Femmes avec c.o. : $n_2 = 33$; $m_2 = 2,5$ nmol ; $s_2 = 0,3$ nmol

Exemple

Hypothèses :

H0 : les moyennes ne sont pas différentes, ce sont 2 estimateurs du taux de T3 libre chez la femme en général

Choix du test :

Variable 1 : prise ou non de la pilule → qualitatif

Variable 2 : dosage de T3 → quantitatif

→ **Test de comparaison de moyennes**

$$\underline{\varepsilon_t} = 1,96$$

$$\underline{\varepsilon_c} = 6.94$$

$\varepsilon_c > \varepsilon_t$ donc on rejette H0 avec $p < 0.0001$

Lien entre variables qualitatives et quantitatives

Petits échantillons
(n_1 et $n_2 < 30$)

TEST T DE STUDENT

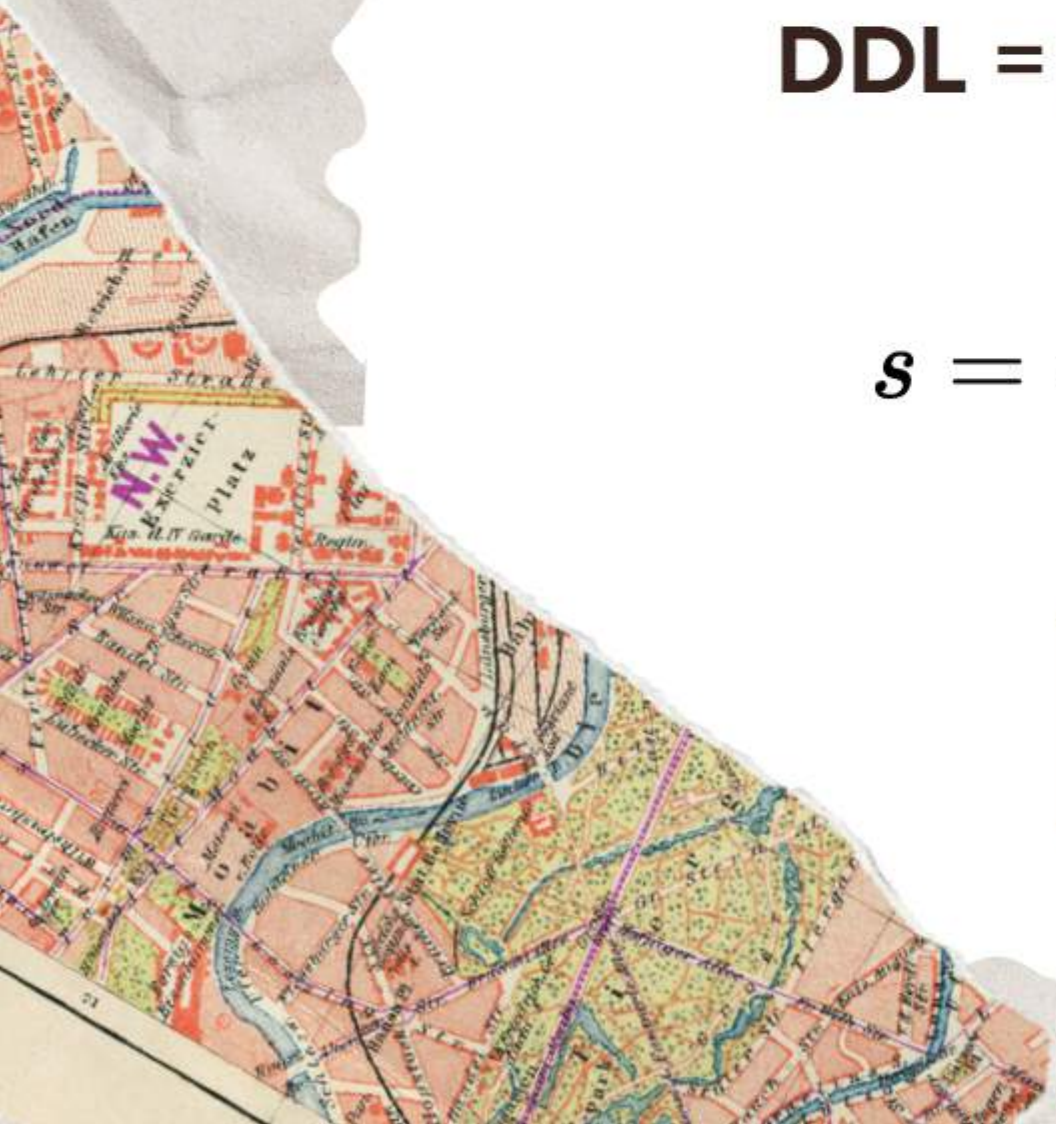
Paramètre Z \rightarrow t

t_t vient de la table du t de Student

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

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

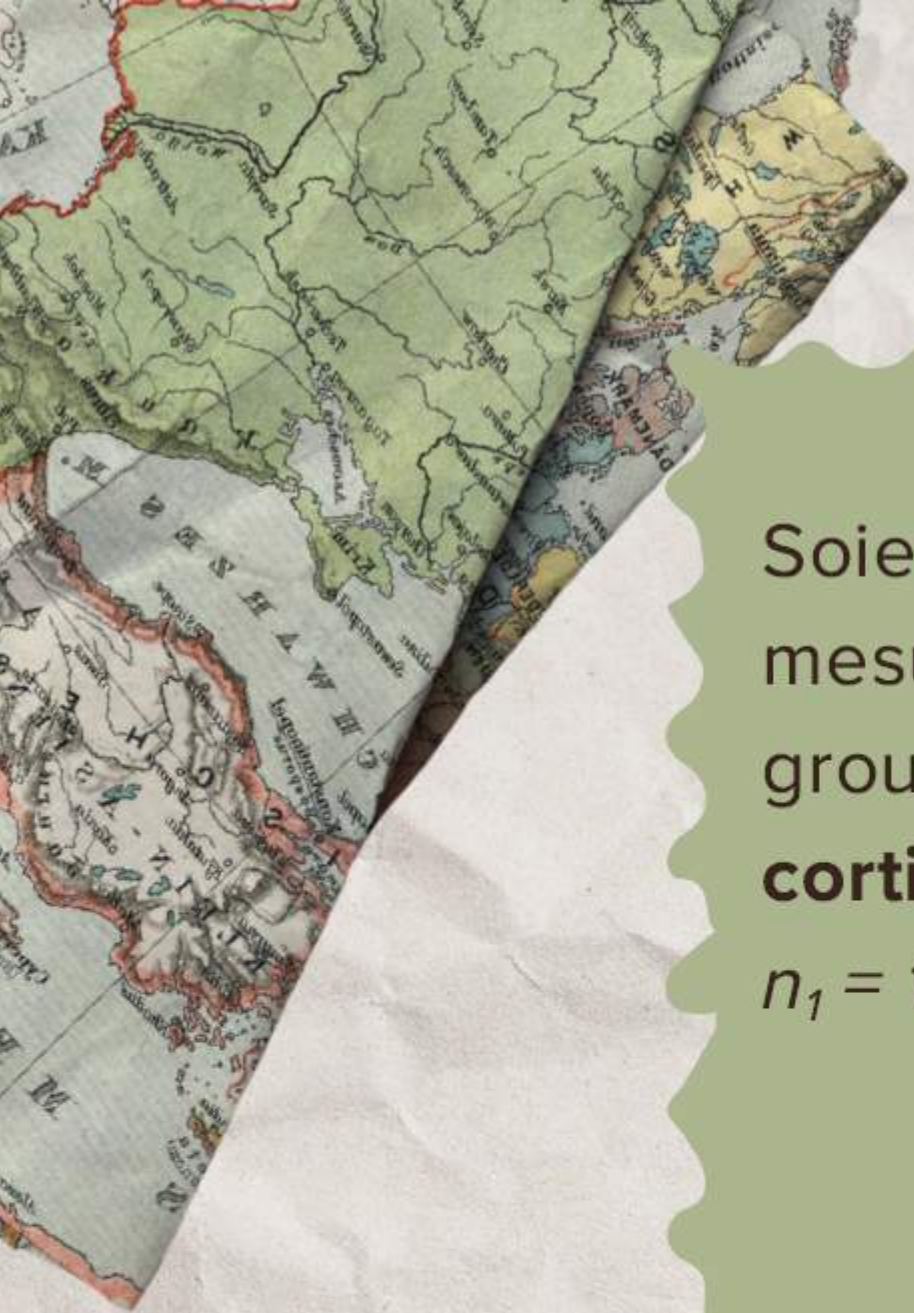
Si $t_c > t_t \rightarrow$ **rejet de H0**



Exemple

Soient 15 femmes obèses et 12 femmes de poids normal. On mesure le taux de corticoïdes sanguins moyens dans chaque groupe. **L'obésité a-t-elle une influence sur le taux de corticoïdes ?**

$$n_1 = 15 ; m_1 = 6.3 ; s_1 = 1.8 / n_2 = 12 ; m_2 = 4.5 ; s_2 = 1.6$$



Exemple

Hypothèses :

H_0 : m_1 et m_2 ne sont pas différents dans les 2 groupes

Choix du test :

Variable 1 : obèse ou non → qualitatif

Variable 2 : taux de corticoïdes → quantitatif

n_1 et $n_2 < 30$ → **test t de Student**

$t_c = 2,92$

DDL = $15 + 12 - 2 = 25$ donc d'après la table **$t_t = 2,06$**

$t_c > t_t$ donc on rejette H_0 au seuil 5%

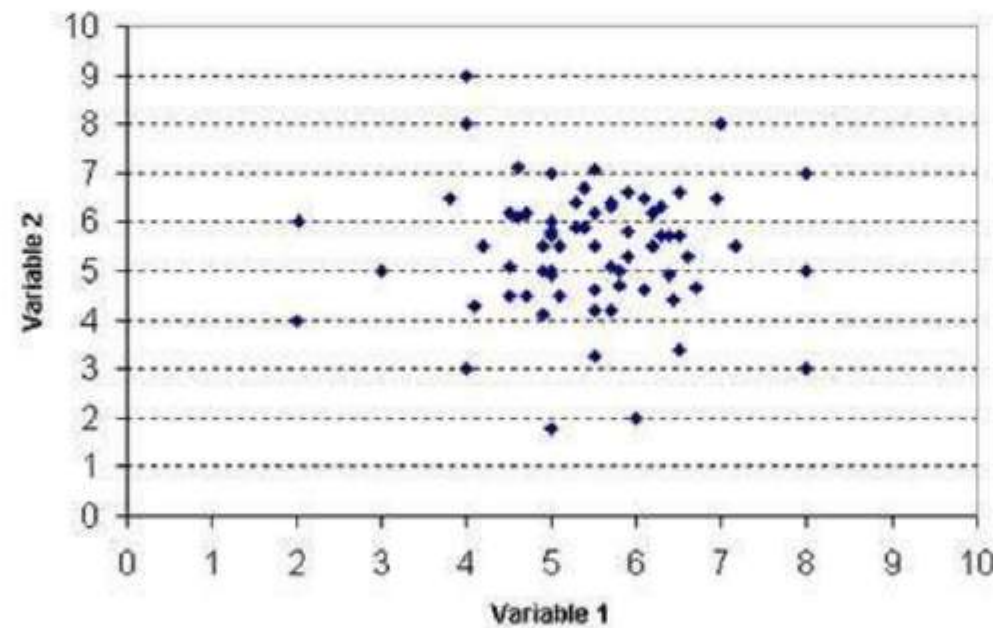
$p < 1\%$ après lecture dans la table. On rejette H_0 à 1% à postériori

Lien entre 2 variables quantitatives

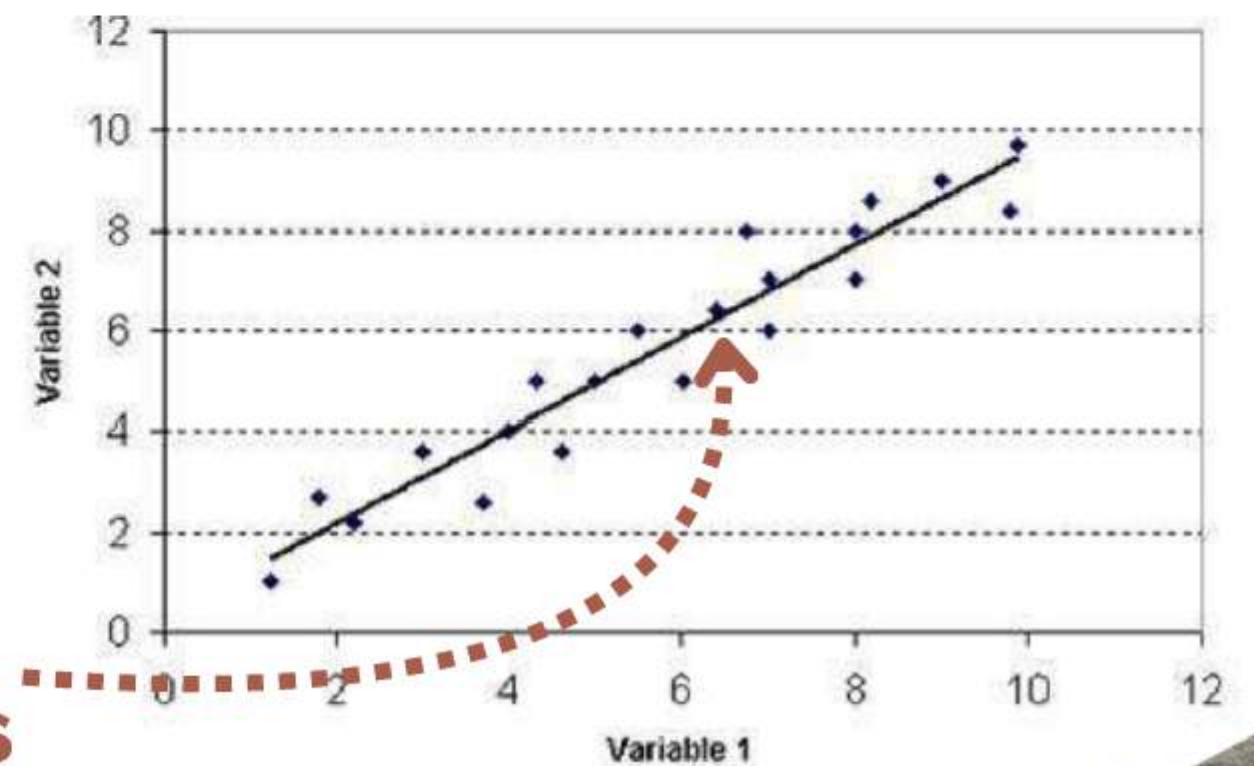
CORRÉLATION ET RÉGRESSION

Corrélation = évaluation de la liaison entre 2 variables quantitatives

Régression = méthode mathématique expliquant les relations entre variables observées



Corrélation \neq causalité !



**Droite de régression =
droite des moindres carrés**

Lien entre 2 variables quantitatives

TEST DE CORRÉLATION DE PEARSON

Paramètre $Z \rightarrow r$

r_t vient de la table du coefficient de corrélation

DDL = $n - 2$

Si $r_c > r_t \rightarrow$ rejet de H_0

$4 < n < 12$

TEST U DE MANN ET WHITNEY

Hypothèse testée : **les moyennes de deux groupes de données sont proches**

Lien entre variables quantitatives et qualitatives

Paramètre $Z \rightarrow u$

u_t vient des tables du test de Mann et Whitney.

Si $u_c > u_t \rightarrow$ acceptation de H_0



Méthodologie

Comment trouver le paramètre théorique pour un test non paramétrique ?

n_1 est le plus petit des 2 effectifs, U le plus petit des 2 U calculés

$n_2 - n_1$	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
2		-	0	2	5	8	12	17	23	29
3		-	0	3	6	10	14	19	26	33
4		-	1	4	7	11	16	22	28	36
5		-	2	4	8	13	18	24	31	39
6		0	2	5	9	14	20	26	34	42
7		0	3	6	11	16	22	29	37	45
8		0	3	7	12	17	24	31	39	48
9		0	4	8	13	19	26	34	42	52
10		1	4	9	14	21	28	36	45	55
11		1	5	10	15	22	30	38	48	
12		1	5	11	17	24	32	41	50	
13		1	6	11	18	25	34	43		
14		1	6	12	19	27	36	45		
...										
18		2	8	16	24	33				
19		3	9	17	25					
20		3	9	17	27					

Exemple

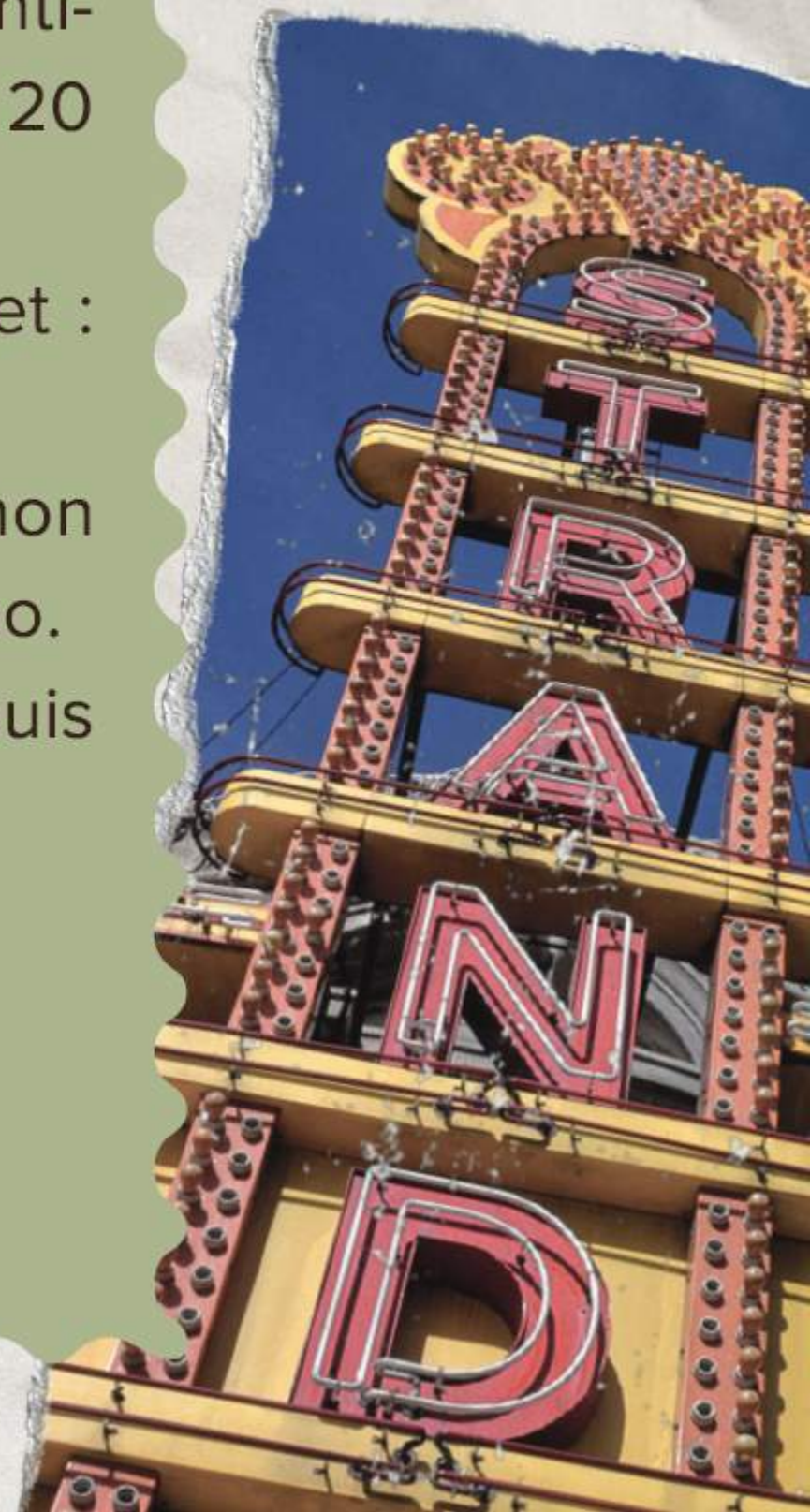
On veut savoir si une nouvelle molécule présente un effet anti-dépresseur. Pour cela, on organise un essai portant sur 20 malades dépressifs, répartis en 2 groupes.

Les 20 malades sont répartis par TAS en 2 groupes de 10 sujet : l'un recevant la **nouvelle molécule**, l'autre recevant le **placébo**.

On évalue les patients à l'aide d'une échelle numérique de **0** (non déprimé) à **50** (très déprimé). Le groupe témoin reçoit le placébo.

Les patients des 2 groupes sont évalués **avant le traitement** puis **après le traitement** au bout de 28 jours.

Le traitement est-il efficace ?



Exemple

Hypothèses :

H0 : il n'y a pas de différence de niveau de déprime entre les participants prenant le placebo et ceux prenant la molécule

Choix du test :

Variable 1 : placebo ou molécule → qualitatif

Variable 2 : score de déprime → pseudo-quantitatif

→ **test U de Mann et Whitney**

$$\underline{u_c} = 9$$

$$\underline{u_t} = 23$$

$u_c < u_t$ donc on rejette H0 au risque 5%

CCI : le traitement est efficace contre la dépression.

$4 < n < 12$

TEST R' DE SPEARMAN

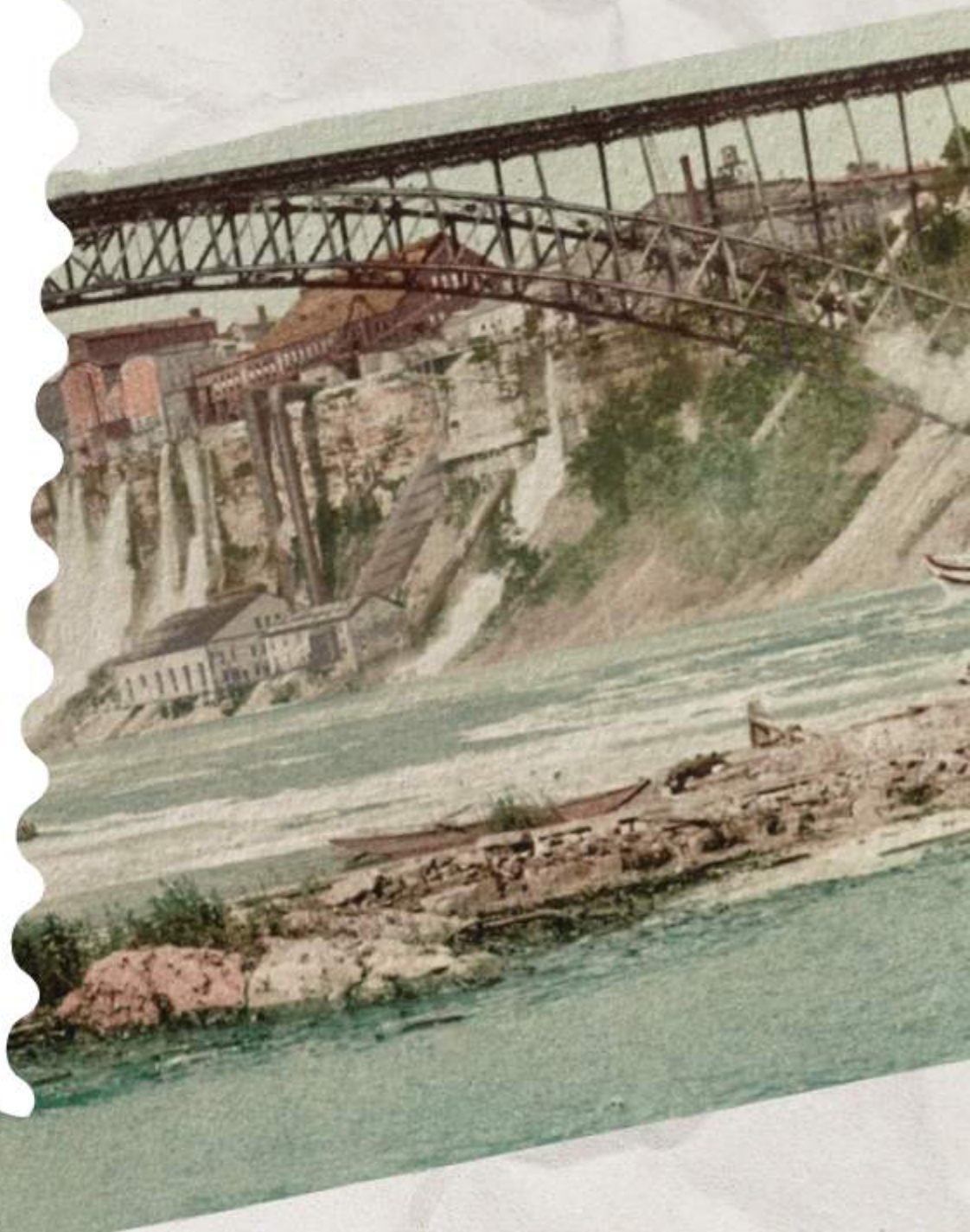
Paramètre $Z \rightarrow r'$

Lien entre variables quantitatives

r'_t vient de la table du r' de Spearman

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

Si $r'_c > r'_t \rightarrow$ **acceptation de H_0**



Exemple

On a recensé pour 6 étudiants les notes obtenues au concours de PACES en biostatistiques, et le classement final à ce même examen.

On cherche à établir s'il existe une relation entre cette note et le classement final.

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

Exemple

Hypothèses :

H0 : il n'y a pas de lien entre ces 2 séries de valeurs numériques, il s'agit de 2 séries indépendantes

Choix du test :

Variable 1 : note → quantitatif

Variable 2 : classement → pseudo-quantitatif
→ **test r' de Spearman**

$$\underline{r'_c = -1}$$

$$\underline{r'_t = 0,89} \text{ (}\alpha = 5\%) \text{ et } \underline{r'_t = 1} \text{ (}\alpha = 1\%)$$

$r'_c < r'_t$ donc on rejette H0 au risque 1%

Il s'agit de 2 séries corrélées. Plus la note de biostat est élevée, plus petit est le rang de classement (d'où le signe - pour r').

Tableau récap

Effectif	Données quantitatives	Données qualitatives	Données quantitatives et qualitatives
$4 < n < 12$	r' de Spearman	Comp % ou χ^2	U de Mann & Whitney
$12 \leq n < 30$	Coeff de corrélation r	Comp % ou χ^2	t de Student
$n \geq 30$	Coeff de corrélation r	Comp % ou χ^2	Comp moyennes