

Revision

1

Table 1
Baseline characteristics of patients with SARS-CoV-2 infection.

Characteristic	COVID-19 patients			P value
	Total (N = 80)	FPV (N = 35)	LPV/RTV (N = 45)	
Age, median (IQR)	47 (35.75–61)	43 (35.5–59)	49 (36–61)	0.61
Age subgroup				
15–44	36 (45.0%)	18 (51.4%)	18 (40.0%)	—
45–64	33 (41.3%)	13 (37.1%)	20 (44.4%)	—
> 65	11 (13.7%)	4 (11.4%)	7 (15.6%)	0.47
Male	35 (43.8%)	14 (40.0%)	21 (46.7%)	0.55
BMI, median (IQR)	22.9 (16.2–31.6)	22.7 (16.2–31.6)	23.1 (16.4–28.4)	0.51
Epidemiology				
History of visiting Wuhan City	46 (57.5%)	20 (57.1%)	26 (57.8%)	—
Not been to Wuhan City	34 (42.5%)	15 (20.0%)	19 (17.8%)	0.95
Onset symptoms				
Fever	59 (73.8%)	22 (62.9%)	37 (82.2%)	0.11
Cough	22 (27.5%)	12 (34.3%)	10 (22.2%)	0.23
Headache/myalgia	8 (10.0%)	3 (8.6%)	5 (11.1%)	1.00
Chill	1 (1.3%)	0 (0%)	1 (2.2%)	1.00
Diarrhea	1 (1.3%)	1 (2.9%)	0 (0%)	0.44
Stuffy nose/sore throat	8 (10.0%)	6 (17.1%)	2 (4.4%)	0.13
Laboratory test, median (IQR)				
WBC ($\times 10^9 L^{-1}$)	6.0 (3.5–5.2)	8.1 (3.8–6.6)	4.3 (3.4–4.9)	0.21
Neutrophils ($\times 10^9 L^{-1}$)	2.8 (2.1–3.4)	3.0 (2.1–3.7)	2.6 (2.1–3.1)	0.43
Lymphocyte ($\times 10^9 L^{-1}$)	1.3 (0.9–1.6)	1.5 (1.0–1.8)	1.2 (0.9–1.4)	0.06
ALT (U L ⁻¹)	22.2 (15.0–26.3)	21.6 (15.0–24.0)	22.6 (15.5–27.0)	0.54
AST (U L ⁻¹)	25.1 (18.0–28.0)	24.1 (18.0–26.0)	25.8 (19.0–31.0)	0.47
GGT (U L ⁻¹)	25.5 (14–31.1)	26.9 (14.0–33.0)	24.4 (14.4–31.1)	0.48
CRP (mg dL ⁻¹)	18.6 (5.0–20.0)	15.0 (3.0–19.2)	21.4 (5.0–23.2)	0.33
IL-6 (ng L ⁻¹)	13.4 (4.4–16.2)	14.0 (3.5–11.0)	12.9 (5.3–16.8)	0.77
T lymphocyte count	973.8 (594.3–1227.0)	1046.7 (600.8–1314.8)	925.2 (572.8–1211.5)	0.40
CD4 ⁺ T lymphocyte count	562.3 (382.5–733)	593.3 (369.0–802.75)	542.3 (388.0–689.0)	0.54
CD8 ⁺ T lymphocyte count	354.4 (206.5–496.5)	397.8 (212.3–528.5)	326.4 (207.5–423)	0.76
Ct values, median (IQR)	30.0 (26.5–33.8)	30.7 (28.0–33.3)	29 (26.0–34.0)	0.38
Chest CT score, median (IQR)	9.5 (4.0–14.0)	12 (4.0–14.0)	9 (4.5–14.0)	0.78

BMI: body mass index; WBC: white blood cell; ALT: alanine aminotransferase; AST: aspartate aminotransferase; GGT: γ -glutamyl transpeptidase; CRP: c-reactive protein; IL: interleukin; Ct: cycle threshold.

Cao Q, Yang M, Liu D, Chen J, Shu D, Xia J, Xiao X, Gu Y, Cai Q, Yang Y, Shen C, Li X, Peng L, Huang D, Zhang J, Zhang S, Wang F, Liu L, Chen L, Chen S, Wang Z, Zhang Z, Cao R, Zhong W, Liu Y, Liu L. Experimental Treatment with Favipiravir for COVID-19: An Open-Label Control Study. Engineering (Beijing). 2020 Oct;6(10):1192–1198. doi: 10.1016/j.eng.2020.03.007.

2

Zhang et al; Home Blood Pressure and the COVID-19 Outbreak

Table 1. Baseline Characteristics of Included Patients in Wuhan and Non-Wuhan Areas of China in the Present Study

Characteristics at the preepidemic phase	Included patients			P value*
	Total (n=7394)	Wuhan (n=283)	Non-Wuhan (n=7111)	
Age, y	68.6±4.8	68.3±4.8	68.6±4.8	0.35
Distribution of age, no. (%)				
60–70 y	5112 (69.1)	203 (71.7)	4909 (69.0)	0.57
≥70 y	2282 (30.9)	80 (28.3)	2202 (31.0)	
Men, no. (%)	3429 (46.4)	146 (51.6)	3283 (46.2)	0.07
Body mass index, kg/m ²	25.6±3.1	25.6±3.5	25.6±3.1	0.85
Morning SBP, mm Hg	131.2±9.8	132.1±10.0	131.1±9.8	0.10
<140 mmHg, no. (%)	6026 (81.5)	220 (77.7)	5806 (81.6)	0.24
140–149 mmHg, no. (%)	1167 (15.8)	53 (18.7)	1114 (15.7)	
≥150 mmHg, no. (%)	201 (2.7)	10 (3.1)	191 (2.7)	
Morning DBP, mmHg	79.6±7.7	80.5±7.7	79.6±7.7	0.05
Fasting glucose, mmol/L	6.2±1.7	6.0±1.5	6.2±1.7	0.02

Zhang S, Zhou X, Chen Y, Wang L, Zhu B, Xiang Y, Bu P, Liu W, Li D, Li Y, Tao Y, Ren L, Fu L, Li Y, Shen X, Liu H, Sun Q, Xu X, Bai L, Zhang W, Cai L, STEP Study Group. Changes in Home Blood Pressure Monitored Among Elderly Patients With Hypertension During the COVID-19 Outbreak: A Longitudinal Study in China Leveraging a Smartphone-Based Application. *Circ Cardiovasc Qual Outcomes*. 2021 May;14(5):e007098. doi: 10.1161/CIRCOUTCOMES.120.007098.

3

3

Table 5
Statistics of adverse reactions after medication.

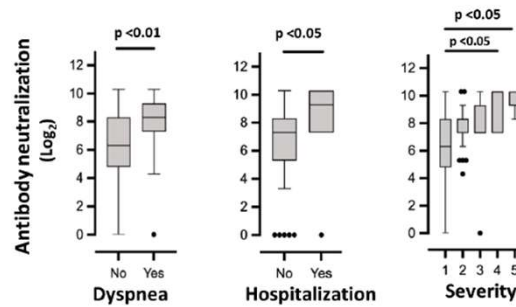
Characteristic	Treatment		
	FPV (N = 35)	LPV/RTV (N = 45)	P value
Total number of adverse reactions	4 (11.43%)	25 (55.56%)	< 0.001
Diarrhea	2 (5.71%)	5 (11.11%)	0.46
Vomiting	0 (0%)	5 (11.11%)	0.06
Nausea	0 (0%)	6 (13.33%)	0.03
Rash	0 (0%)	4 (8.89%)	0.13
Liver and kidney injury	1 (2.86%)	3 (6.67%)	0.63
Others	1 (2.86%)	2 (4.44%)	1.00

4

Cai Q, Yang M, Liu D, Chen J, Shu D, Xia J, Xiao X, Gu Y, Cai Q, Yang Y, Shen C, Li X, Peng L, Huang D, Zhang J, Zhang S, Wang F, Liu J, Chen L, Chen S, Wang Z, Zhang Z, Cao R, Zhong W, Liu Y, Liu L. Experimental Treatment with Favipiravir for COVID-19: An Open-Label Control Study. *Engineering (Beijing)*. 2020 Oct;6(10):1192–1198. doi: 10.1016/j.eng.2020.03.007.

4

Figure 1



Note: Adapted from Salazar et al. Boxplots of interquartile range, showing the log₂ transformed antibody neutralization titer by category in 68 CP donors. Black dots are outliers, whisker bars are upper and lower quartiles. Severity- highest severity score during illness. Open access journal; all content freely available.

Salazar E, Kuchipudi SV, Christensen PA, Eagan T, Yi X, Zhao P, Jin Z, Long SM, Olsen RJ, Chen J, Castillo R, Leveque C, Towers D, Lavinder J, Gollmar J, Cardona J, Ippolito G, Nisly R, Bird I, Greenwalt D, Rossi RM, Gortu A, Srinivasan S, Progeny I, Cattadori IM, Hudson RJ, Joshiy NK, Progeny L, Hsieh K, Heberts A, Bernard DM, Dye JM, Kupper V, Musser JM. Convalescent plasma with SARS-CoV-2 spike protein ectodomain and receptor binding domain IgG correlate with virus neutralization. J Clin Invest. 2020 Dec 1;130(12):6728-6738. doi: 10.1172/JCI141206.

5

Comment choisir un test statistique

Questions a répondre:

- Combien des **échantillons (groups)** on a?
- Les échantillons sont:
 - **Dépendantes/appariées?**
 - (jumeaux/ données répétées/comparaison du partie gauche et droit d'un sujet/études appariées/deux tests diagnostiques ou méthodes de mesure qui observent les mêmes sujets)
 - **Indépendantes?**
- Quel est le **type des variables?**
- Combien des sujets?
 - Pour les **données qualitatives:**
 - Tableau de contingence: % des cellules theoriques < 5?
- Quelle est la nature des données?
 - Pour les **données quantitatives:**
 - **Distribution normale?**
 - **Variances égales/ inégales?**

6

Precisions

- Ecart type = déviation standard
- Notations:
 - s – déviation standard,
 - S – déviation standard d' échantillonnage
 - $S = s * \sqrt{n/(n-1)}$
- Echantillons:
 - Indépendants
 - Dépendants (appariées)
- Distribution gaussienne = normale

7

7

Les hypothèses statistiques

- La création des hypothèses:
 - Certaines questions médicales ont deux réponses opposées
 - on force beaucoup des questions dans ce format
 - Les réponses correspondent aux deux **modèles** possibles de la réalité
 - Ces deux modèles sont nommées: **hypothèses**
 - L' hypothèse **nulle**: H_0
 - Il **n'y a pas** une **différence** statistiquement significative entre 2/>>=2 groups (ex. un **traitement** [ibuprofène vs. placebo] ou un **facteur de risque** [présent vs. absent]) en ce qui concerne la **moyenne/ médiane/ variance/ fréquence ... d'une caractéristique** (ex **résultat du traitement**: la fréquence de la guérison [oui vs. non] ou la moyenne de la température)
 - Il **n'y a pas** une **relation/lien/association/dépendance/corrélation** statistiquement significative entre 2 caractéristiques/variables: (ex. un **Facteur de risque** [présent vs. absent] – une **maladie** [oui vs. non] , ou un **traitement** [ibuprofène vs. placebo] – le **résultat du traitement** (ex: la fréquence de la guérison [oui vs. non] ou la moyenne de la **température**))
 - L' hypothèse **alternative**: H_1 (négation du H_0)
 - Il **y a** une **différence** statistiquement significative entre 2/>>=2 groups en ce qui concerne la **moyenne/ médiane/ variance/ fréquence ... d'une caractéristique**
 - Il **y a** une **relation/lien/association/dépendance/corrélation** statistiquement significative entre 2 caractéristiques/variables
- Les tests statistiques nous permettent de faire le choix entre les deux possibilités (H_1 / H_0)

8

Données	Nombre échantillons	Tests paramétriques	Compare	Tests non paramétriques	Compare
qualitatives	2 / > 2 indépendants	Chi deux <20% cellules tableau théorique/attendue <5	fréquences	exact Fisher >20% cellules tableau théorique/attendue <5	fréquences
	2 dépendants (appariés)	Mc Nemar	fréquences		
(Tests de normalité: Kolmogorov Smirnov, Shapiro Wilk, HB - il n'y a pas de différence entre la distribution et la distribution normale)					
quantitatives		Données normale distribuées Test pour normalité p>0,05		Données non normale distribuées Test pour normalité p<0,05	
	2 indépendants	Student (t) pour échantillons indépendants avec variances Egales Inégales Tests pour variances: F, Levene Bartlett HB: V1-V2	moyennes	Mann Whitney U (Wilcoxon somme des rangs)	~médianes (distributions)
	2 appariés (dépendants)	Student (t) pour échantillons appariés / dépendants	moyennes	Wilcoxon pour échantillons appariés (Wilcoxon rangs signées)	~médianes (distributions)
	> 2 indépendants	ANOVA (pour variances égales) ou ANOVA de Welch ou Brown Forsyth (pour variances inégales)	moyennes	Kruskal Wallis	~médianes (distributions)

9

Les étapes d'un test statistique

- **Étape 1. Formuler les hypothèses statistiques:**
- **Étape 2. Décider sur une statistique appropriée du test (paramètre du test)**
- **Étape 3. Sélectionner le niveau de signification - la valeur alpha. $\alpha = 0,05$**
- **Étape 4. Déterminer la valeur critique (v.c.) de la statistique du test**
 - on détermine - une région critique ou région de rejet (RR)
 - Lois t, Z: $RR = (-\infty, -v.c.] \cup [v.c., +\infty)$
 - Lois F, Khi 2: $RR = [v.c., +\infty)$
- **Étape 5. Calculer la valeur de la statistique / paramètre du test**
- **Étape 6. la décision statistique en fonction de la région critique :**
 - Si Z_0 est dans RR (région du rejet/critique) on rejette H_0 et on accepte H_1
 - Il y a une différence/ Il y a une relation –statistiquement significative
 - Si Z_0 est dans RnR (région de non rejet) on reste avec H_0 (on peut pas rejeter H_0)
 - On ne peut pas dire qu'il y a une différence/ il y a une relation statistiquement significative
- **Étape 6'. La décision avec p-value**
 - La décision avec p-value.
 - Si $p\text{-value} < \alpha (=0,05)$ on rejette $H_0 \Rightarrow$ on accepte H_1
 - Il y a une différence/ Il y a une relation –statistiquement significative
 - Si $p\text{-value} \geq \alpha (=0,05) \Rightarrow$ on reste avec H_0
 - On ne peut pas dire qu'il y a une différence/ Il y a une relation –statistiquement significative

10

10

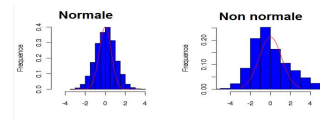
Vérification de la condition de normalité des données

Utilité:

- Importante pour appliquer des test paramétriques, avec condition de normalité:
 - Test Z pour les moyennes
 - Test t (Student)
 - Test ANOVA

Modalités de vérification (ici, conditions de normalité):

- des graphiques (les meilleures) modalités
 - Histogramme (symétrique, comme un chapeau)
 - Boite à moustaches (symétrique autour de la médiane)
 - Le graphique des quantiles (voir diapositive suivant)
- des statistiques descriptives (pas très fiables)
 - Si la moyenne est \sim médiane
 - Si le coefficient de l'aplatissement ~ 0 / appartient à $[-1, 1]$ (kurtosis)
 - Si le coefficient de symétrie ~ 0 / appartient à $[-1, 1]$ (skewness)
- des tests de normalité: (ne sont pas recommandées)
 - Test de Kolmogorov-Smirnov ($p < 0,05$ – non normale, $p > 0,05$ \sim normale)
 - Test de Shapiro-Wilk ($p < 0,05$ – non normale, $p > 0,05$ \sim normale)



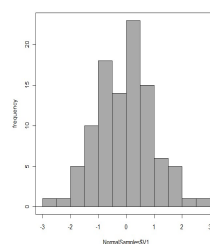
11

Comparaison des données normale/non normale distribuées

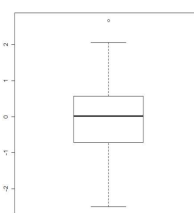
Normale

moyenne \sim médiane
 (= -0,03 = 0,015)
 c. asymétrie = 0,11
 appartient à $[-1, 1]$, ~ 0
 c. aplatissement = -0,09
 appartient à $[-1, 1]$, ~ 0
 Shapiro-Wilk test
 $p = 0,99 > 0,05$

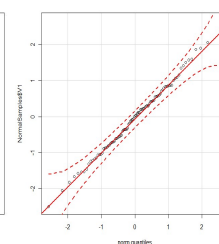
Histogramme



Boite à moustaches

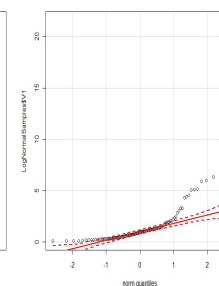
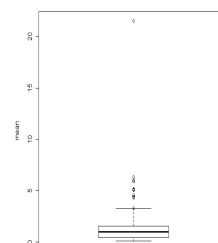
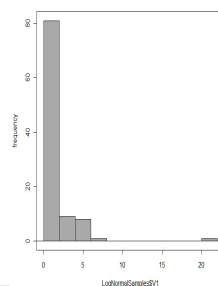


Le graphique des quantiles



Non normale

moyenne \neq médiane
 (= 1,57 = 0,98)
 c. asymétrie = 5,59
 > 1 , < 0
 c. aplatissement = 40,63
 > 1 , < 0
 Shapiro-Wilk test
 $p \sim 0 < 0,05$



12

La normalité des données

- Test de normalité des données :
 - Test de Kolmogorov-Smirnov
 - Si < 50 sujets le test Shapiro-Wilk
- H0 = aucune différence statistiquement significative entre la distribution observée et la distribution normale
- H1 = aucune différence statistiquement significative entre la distribution observée et la distribution normale
- $p < 0,05$ on rejette l'hypothèse nulle, les données ne sont pas normalement distribuées
- $p > 0,05$ on ne rejette pas l'hypothèse nulle, on n'a pas des motifs pour considérer les données anormales – on peut considérer les données normales distribuées (dans certaines conditions)

13

Récapitulatif des tests utilisés

Tests pour variables quantitatives - comparer la moyenne des deux échantillons

Type variable	Nb sujets	Nature des données	Statistique comparée	Test utilisé	Paramètre du test	Région du rejet – test bidirection
Deux échantillons indépendants Pour l'examen écrit je ne demande pas les choses en gris						
Quantitative	$n_1, n_2 \geq 30$	Normale distribuées, Variances dans la population <u>connues</u> inégales	Différence des moyennes	Test Z pour la différence entre les moyennes	$Z = \frac{m_1 - m_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$	$(-\infty, -v.c.] \cup [v.c., +\infty)$
	$n_1, n_2 \geq 30$	Normale distribuées, Variances dans la population <u>connues</u> égales	Différence des moyennes	Test Z pour la différence entre les moyennes	$Z = \frac{m_1 - m_2}{\sqrt{\frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2}}}$	$(-\infty, -v.c.] \cup [v.c., +\infty)$
	$n_1, n_2 \geq$ ou < 30	Normale distribuées, Variances dans la population <u>inconnues</u> inégales	Différence des moyennes	Test t (Student) $n_1 + n_2 - 2$ d.d.l.	$t = \frac{m_1 - m_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$	$(-\infty, -v.c.] \cup [v.c., +\infty)$
	$n_1, n_2 \geq$ ou < 30	Normale distribuées, Variances dans la population <u>inconnues</u> égales	Différence des moyennes	Test t (Student) $n_1 + n_2 - 2$ d.d.l.	$t = \frac{m_1 - m_2}{\sqrt{\frac{s^2}{n_1} + \frac{s^2}{n_2}}}$ $s^2 = \frac{(n_1 - 1) \times S_1^2 + (n_2 - 1) \times S_2^2}{n_1 + n_2 - 2}$	$(-\infty, -v.c.] \cup [v.c., +\infty)$
Deux échantillons dépendants (appariées)						
Quantitative	$n_1 = n_2 \geq$ ou < 30	Normale distribuées,	Moyenne des différences	Test t (Student) $n - 1$ d.d.l.	$t = \frac{\bar{d}}{\frac{S}{\sqrt{n}}}$	$(-\infty, -v.c.] \cup [v.c., +\infty)$
(ou n, n_1, n_2 - nombre des sujets; m, m_1, m_2 - moyennes; s, s_1, s_2 - déviations standard descriptive de l'échantillon; S, S_1, S_2 - déviation standard d'échantillonnage; $S = \sqrt{\frac{n}{n-1}} s$; $S = \sqrt{\frac{n-1}{n}} s$; $\sigma, \sigma_1, \sigma_2$ - déviation standard dans la population; pour $\alpha=0,05$, $Z_{\alpha}=1,96$; d.d.l. - degrés de liberté; v.c. - valeur critique)						

14

Récapitulatif des tests utilisés

Tests statistiques pour **comparer deux échantillons, variables quantitatives non normale distribuées**

Tests statistiques pour comparer deux échantillons indépendants					
Type variable	Nature des données	Statistique comparée	Test utilise	Paramètre du test	Région du rejet
Quantitative ou qualitative ordinale	Données non normale distribuées	~médiane (si les distributions sont similaires) - la distribution des rangs /moyenne des rangs	Test Mann Whitney U Etapas: Tableaux avec toutes les valeurs des deux échantillons Ordonnées en ascendant Numéroter du plus petit a le plus grand Donner un rang. (Valeurs identiques, le rang = moyenne de la numérotation) Calculez la somme des rangs pour group A et B	obsSR _a somme rangs group a, obsSR _b somme rangs group b, $U_a = n_a \times n_b + \frac{n_a(n_a+1)}{2}$ $U_b = n_a \times n_b + \frac{n_b(n_b+1)}{2}$ obsSR _a	Min (U _a , U _b) ≤ v.c.
Tests statistiques pour deux échantillons dépendants/ appariées					
Quantitative ou qualitative ordinale	Données non normale distribuées	~médiane (si les distributions sont similaires) - la distribution des rangs	Test Wilcoxon pour échantillons appariées Etapas: Calculer la différence des paires Ignorer les zéros Ignorer les signes Numéroter du plus petit a le plus grand Donner le rang 1, 2, 3... (Valeurs identiques, le rang = moyenne de la numérotation) Calculez la somme des rangs positifs (W+), et puis négatifs (W-).	W ₊ somme rangs négatifs W ₋ somme rangs positifs	Min (W ₊ , W ₋) ≤ v.c.

v.c. = valeur critique; n_a = nombre sujets dans group a; n_b = nombre sujets dans group b; Min = le minimum entre les deux valeurs

15

Récapitulatif des tests utilisés

Tests statistiques pour deux échantillons <u>indépendants</u> , variables qualitatives dichotomiques						
Type variable	Nb sujets	Nature des données	Statistique comparée	Test utilise	Paramètre du test	Région du rejet
Qualitative dichotomique	$n_1 p_1 > 10$, $n_2 p_2 > 10$, $n_1 (1-p_1) > 10$, $n_2 (1-p_2) > 10$		fréquences	Test Z * Equivalent au test Chi carrée ci dessous (Chi carrée est recommandé):	$Z = \frac{P_1 - P_2}{\sqrt{p(1-p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$ $p = \frac{P_1 n_1 + P_2 n_2}{n_1 + n_2}$	$(-\infty, -Z_{\frac{\alpha}{2}}] \cup [Z_{\frac{\alpha}{2}}, +\infty)$
	<20% cellules du tableau théorique sont <5		fréquences	Test Chi carrée * $v = (L-1) * (C-1) = 1$ d.d.l.	$\chi^2 = \sum_{i=1}^L \sum_{j=1}^C \frac{(f_{ij}^o - f_{ij}^t)^2}{f_{ij}^t}$	$[\chi_{\alpha, v}^2, +\infty)$
	>20% cellules du tableau théorique sont <5		fréquences	Test exact Fisher	- (test non paramétrique)	-
Tests statistiques pour deux échantillons <u>dépendants</u> , variables qualitatives dichotomiques						
Qualitative dichotomique	$b+c > 25$		fréquences	Test Chi carrée $v=1$ d.d.l.	$\chi_{1,ddl}^2 = \frac{(b-c)^2}{b+c}$	$[\chi_{\alpha, v}^2, +\infty)$

* On préfère pour ce cours (et l'examen) le test Chi carrée pour comparer les fréquences, au lieu du test Z pour comparer les fréquences

p₁, p₂ – fréquences; n₁, n₂ – nombre des sujets; L et C – nombres des lignes et des colonnes dans le tableau de contingence, f^o – fréquence observée, f^t – fréquence théorique; d.d.l. – degrés de liberté;

Pour l'examen écrit je ne demande pas les choses en gris

16

Récapitulatif des tests utilisés

Tests statistiques pour plus des deux échantillons (groups) indépendants						
Type variable	Nb sujets	Nature des données	Statistique comparée	Test utilise	Paramètre du test	Région du rejet
Qualitative	<20% cellules du tableau théorique sont <5		Fréquence	Test Chi carrée $v=(L-1)*(C-1)$ d.d.l.	$\chi^2 = \sum_{i=1}^L \frac{(f_i^o - f_i^t)^2}{f_i^t}$	$\left[\chi_{v,\alpha}^2, +\infty \right)$
	>20% cellules du tableau théorique sont <5		fréquence	Test exact Fisher	-	-
Quantitative		Normale distribuées, Variance des échantillons égaux	moyenne	Test ANOVA $v1 = \text{nb. groups} - 1$ $v2 = \text{nb. obs.} - \text{nb groups}$	$F = \frac{MCG}{MCE}$	$\left[F_{v_1, v_2, \alpha}, +\infty \right)$
		Normale distribuées, Variance des échantillons inégaux	moyenne	ANOVA de Welch ou Brown Forsyth		
		Non normale distribuées	médiane	test Kruskal Wallis	-	-

L et C – nombres des lignes et des colonnes dans le tableau de contingence, f_i^o – fréquence observée, f_i^t – fréquence théorique; d.d.l. – degrés de liberté;

MCEG = $[n1 * (m1 - mt)^2 + n2 * (m2 - mt)^2 + n3 * (m3 - mt)^2 + \dots] / (p-1)$
MCDG = $[(n1-1) * DS1^2 + (n2-1) * DS2^2 + (n3-1) * DS3^2 + \dots] / (n-p)$
n=nombre total d'observations, n1, n2, n3 – le nombre d'observations par group, p – nombre des groups
mt – la moyenne des toutes les observations, m1, m2, m3... les moyennes par groups, DS1, DS2, DS3, ... déviations standard d'échantillonnage des groups

17

Récapitulatif des tests utilisés

Tests statistiques pour comparer les variances entre deux échantillons						
Type variable	Nb sujets	Nature des données	Statistique comparée	Test utilise	Paramètre du test	Région du rejet
Quantitative	Données normale distribuées		variances	Test F, $v_1 = n_1$ d.d.l. $v_2 = n_2$ d.d.l.	$F = \frac{\begin{cases} S_2^2, pour. S_2^2 > S_1^2 \\ S_1^2, pour. S_1^2 > S_2^2 \end{cases}}{S_2^2, pour. S_2^2 > S_1^2}$	$\left[F_{v_1, v_2, \alpha}, +\infty \right)$
Tests statistiques pour comparer les variances entre > deux échantillons						
Quantitative			variances	Test Bartlet ou Test Levene		

(ou n_1, n_2 - nombre des sujets; m_1, m_2 - moyennes; s_1, s_2 - déviations standard descriptive de l'échantillon; S_1, S_2 - déviation standard d'échantillonnage;
 $S = \sqrt{\frac{n}{n-1}} s$; $s = \sqrt{\frac{n-1}{n}} S$; d.d.l. - degrés de liberté)

18

Equivalences entre tests paramétriques et non paramétriques

Données	Nombre échantillons	Tests paramétriques	Tests non paramétriques
qualitatives		Chi deux	exact Fisher
Quantitatives (ou qualitatives ordinales)	2 indépendants	Student (t) pour échantillons indépendants	Mann Whitney U (Wilcoxon somme des rangs Mann Whitney Wilcoxon)
	2 appariées (dépendants)	Student (t) pour échantillons appariées	Wilcoxon rangs signées (Wilcoxon pour échantillons appariées)
	> 2 indépendants	ANOVA (pour variances égales) ou ANOVA de Welch ou Brown Forsyth (pour variances inégales)	Kruskal Wallis
		Pour données normale distribuées	Pour données non normale distribuées

19

Intervalle de confiance

Type variable	Nombre échantillons	Estimateur ponctuel	Conditions application	Formule
qualitative	une	fréquence: f	grands échantillons: nf ≥ 10 et nq ≥ 10	$\left(f - Z_{\frac{1-\alpha}{2}} ES; f + Z_{\frac{1-\alpha}{2}} ES\right) \quad ES = \sqrt{\frac{f(1-f)}{n}}$
	une	fréquence: f	petits échantillons: nf < 10 ou nq < 10	on ne va pas calculée
	deux	différence entre les fréquences: f ₁ - f ₂	grands échantillons: f ₁ n ₁ ≥ 10, (1-f ₁)n ₁ ≥ 10, f ₂ n ₂ ≥ 10, (1-f ₂)n ₂ ≥ 10	$ES = \sqrt{\frac{f_1 \times (1-f_1)}{n_1} + \frac{f_2 \times (1-f_2)}{n_2}}$ $\left((f_1 - f_2) - Z_{\frac{1-\alpha}{2}} ES; (f_1 - f_2) + Z_{\frac{1-\alpha}{2}} ES\right)$
	deux	différence entre les fréquences:	petits échantillons: np < 10 ou nq < 10	on va pas calculée
quantitative	un	moyenne: m	grands échantillons: n ≥ 30, σ - connue	$\left[m - Z_{\frac{1-\alpha}{2}} ES; m + Z_{\frac{1-\alpha}{2}} ES\right] \quad ES = \frac{\sigma}{\sqrt{n}}$
	un	moyenne: m	grands échantillons: n ≥ 30, σ - non connue	$\left[m - t_{n-1, \frac{1-\alpha}{2}} ES; m + t_{n-1, \frac{1-\alpha}{2}} ES\right] \quad ES = \frac{S}{\sqrt{n}}$
	un	moyenne: m	petits échantillons: n < 30, σ - non connue	$\left[m - t_{n-1, \frac{1-\alpha}{2}} ES; m + t_{n-1, \frac{1-\alpha}{2}} ES\right] \quad ES = \frac{S}{\sqrt{n}}$
	deux	différence entre les moyennes: m ₁ - m ₂	variances égales	$ES = \sqrt{\frac{(n_1 - 1) \times s_1^2 + (n_2 - 1) \times s_2^2}{n_1 + n_2 - 2}}$ $\left((m_1 - m_2) - t_{n_1+n_2-2, \frac{1-\alpha}{2}} ES; (m_1 - m_2) + t_{n_1+n_2-2, \frac{1-\alpha}{2}} ES\right)$
<div>20</div> <small>(ou n, n1, n2 - nombre des sujets; f, f1, f2 - fréquence observée; q = 1 - f, m, m1, m2 - moyennes; s, s1, s2 - déviations standard descriptive de l'échantillon; S - déviation standard d'échantillonnage; $S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$; $s = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$; σ - déviation standard dans la population; ES=erreur standard; pour α=0.05, Zα=1.96)</small>				

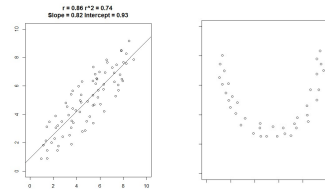
20

Evaluation graphique : diagramme de dispersion (nuage des points/scatter)

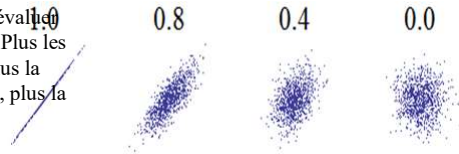
Evaluer la linearité, et l'importance de la corrélation:

Si les points semblent suggérer une droite – la relation est peut être linéaire

Si les points semblent suggérer des tendances qui ne sont pas linéaires, la relation est peut être non linéaire (exponentielle, quadratique)



Si la relation est plus probable linéaire, on peut évaluer d'une manière subjective la corrélation linéaire. Plus les points se rapprochent d'une droite de tendance, plus la corrélation est forte, plus les points sont distants, plus la corrélation est faible



21

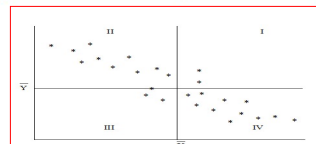
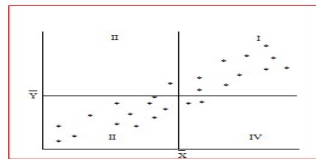
21

Evaluation graphique : diagramme de dispersion (nuage des points/scatter)

But : évaluation visuelle de la relation entre deux variables quantitatives

L'utilisation des cadrans pour identifier la **tendance/sens/direction** directe/inversement proportionnelle :

- i) presque tous les points sont dans les cadrans I et III \Rightarrow tendance croissante / pente ascendante / pente positive / lien (direct) proportionnel
- ii) presque tous les points sont dans les cadrans II et IV \Rightarrow tendance décroissante / pente descendante / pente négative / lien inversement proportionnel
- iii) les points sont distribués uniformément dans tous les cadrans \Rightarrow aucune tendance



22

22

Corrélations

Type des variables	Nature des données	Coefficient de corrélation	Formule du coefficient
quantitative	normale distribuées	Pearson (r)	$COV(X,Y) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$ $r = \frac{COV(X,Y)}{S_X \cdot S_Y}$
quantitative	non normale distribuées	Spearman (ρ - rho)	$\rho = r_s = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}, \text{ ou } d_i = R_{x_i} - R_{y_i}$
qualitative ordinales	-	Spearman (ρ - rho)	$\rho = r_s = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}, \text{ ou } d_i = R_{x_i} - R_{y_i}$

X_i , Y_i – sont les valeurs des deux séries des données. \bar{X} et \bar{Y} sont les moyennes des deux séries. R_{xi} et R_{yi} sont les rangs des valeurs X_i et Y_i après leurs rangement dans ordre croissante. n – nombre des observations. S_x et S_y sont les déviations standard d'échantillonnage. COV(X,Y) – la covariance

23

Covariance COV(X,Y) : **Corrélation linéaire Pearson - interprétations**

- > 0 tendance croissante/ pente ascendante/ lien direct proportionnel/ covariance positive
 - < 0 tendance décroissante/ pente descendante/ lien inversement proportionnel/ covariance négative
 - $\cong 0 \Rightarrow$ aucune tendance
- r** (coefficient de corrélation Pearson):
- montre la direction et l'intensité de la corrélation;
- Interprétation du direction/sens/tendance:
- > 0 tendance croissante/ pente ascendante/ lien direct proportionnel/ corrélation positive
 - < 0 tendance décroissante/ pente descendante/ lien inversement proportionnel/ corrélation négative
 - $\cong 0 \Rightarrow$ aucune tendance
- plus **r** ou COV(X,Y) est grand (en valeur absolue) plus la relation est forte
 - plus **r** ou COV(X,Y) est proche du 0, plus la relation est faible

24

24

Le coefficient Pearson - interprétation

Interprétation de l'intensité/force/degré/importance de la corrélation linéaire avec les règles empiriques de Colton [Colton T. Statistics in Medicine. Little Brown and Company, New York, NY 1974] (on préfère le mot **corrélation** ici, même si association/liens/relation peut être utilisé)

(-0.25 et 0,25)

=> une relation **négligeable** ou **aucune** corrélation linéaire entre les variables

[0.25 et 0.50) ou [-0.25 et -0.50)

=> un degré de corrélation **faible/acceptable**

[0.50 et 0.75) ou [-0.50 et -0.75)

=> un degré de corrélation **modérée à bonne**

[0.75 et 1] ou [-0.75 et -1]

=> une **très bonne à excellente** corrélation

Il y a autre divisions possibles aussi.

Ces règles doit être utilisée avec soins. Elle sont pour donner une idée, mais pour chaque problème, l'intensité de la relation est relative au domaine. Pour certain situations les valeurs en dessous de 0,8 peut être faibles.

25

25

Régression linéaire simple

- **Interprétation**

- La droite de régression $Y(X)$:

$$Y = b_0 + b_1 X$$

b_0 = est l'ordonnée à l'origine – la valeur du Y quand X est égal a 0 (d'habitude cette information n'est pas utile pour les médecins, elle présente une situation qui en réalité est impossible)

b_1 = la pente de la droite de régression.

Interprétation de b_1 - du coefficient de la variable X

chaque unité de mesure de la variable indépendante - X en plus augmente en moyenne la variable dépendante - Y avec la valeur du coefficient de la variable indépendante X - b_1

26

26

Régression linéaire multiple - Quantification de l'importance de la relation pour plusieurs variables

- L' équation du **régression linéaire multiple**:
- Variable dépendante = coefficient_1 * variable_1 + coefficient_2 * variable_2 + ... + coefficient_n * variable_n + coefficient_0
- Ex: triglycérides (mg/dL) = 23,10 * obésité (oui/non) + 1,14 * cholestérol (mg/dL)
- **L' interprétation du coefficient ajusté (adjusted – en anglais)** pour des variables **Qualitatives dichotomiques** (ex. obésité):
 - l'augmentation de la variable dépendante – les triglycérides en moyenne (ici il est de 23,1 mg/dL) pour ceux qui ont le facteur présent (être obèse - la variable indépendante) comparée à ceux qui n'ont pas le facteur (ne sont pas obèses), si on tient les autres variables constantes / si on ajuste les autres variables/ si on contrôle les autres variables) (ici – le cholestérol)
 - ceux qui ont le facteur présent (être obèse - la variable indépendante) ont la variable dépendante – les triglycérides en moyenne plus grand avec 23,1 mg/dL comparée à ceux qui n'ont pas le facteur (ne sont pas obèses), si on tient les autres variables constantes / si on ajuste les autres variables/ si on contrôle les autres variables) (ici – le cholestérol)

27

27

Variables aléatoires discrètes Esperance, variance et écart type

- **Esperance mathématique** ou la moyenne théorique d'une v.a. X

$$E(X) = \sum_{i=1}^n x_i \Pr(X = x_i) \quad \approx \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Exemple: avec le traitement

- X: 0 1 2 3 4

Pr 0,008 0,076 0,256 0,411 0,24

$$E(X) = 0 \times 0,008 + 1 \times 0,076 + 2 \times 0,256 + 3 \times 0,411 + 4 \times 0,24 = 1,31$$

- **Variance** de X:

$$V(X) = \sum_{i=1}^n [x_i - E(X)]^2 \Pr(x_i) \quad \approx \quad \sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$$

- **Ecart-type** (deviation standard) de X

28

$$\sigma(X) = \sqrt{V(X)} = V(X)^{1/2}$$

28

Tableau récapitulatif

Événements	Définitions	Notations	Calcul des probabilités
Événement contraire d'un événement A	l'événement constitué par tous les événements élémentaires qui ne sont pas dans A.	\bar{A}	Propriété : $\Pr(\bar{A}) = 1 - \Pr(A)$
Événement "A et B" ou intersection de A et B	l'événement "A et B" est constitué par tous les événements élémentaires se trouvant à la fois dans A et dans B.	$A \cap B$	
Événement "A ou B" ou réunion de A et B	l'événement "A ou B" est constitué par tous les événements élémentaires se trouvant dans l'un au moins des événements A ou B	$A \cup B$	Propriété : $\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B)$
Événements incompatibles	ils ne peuvent pas être réalisés simultanément.	$A \cap B = \emptyset$	$\Pr(A \cap B) = 0$
Événements indépendants	Deux événements sont indépendants si la survenance d'un événement n'affecte pas l'apparition d'un deuxième événement.		$\Pr(A \cap B) = \Pr(A) \times \Pr(B)$

29

Probabilité conditionnelle

- A et B évènements, $P(A) \neq 0$.

□ Notation: $\Pr(B | A)$: probabilité conditionnelle de B, sachant que l'événement A est réalisé

□ Formule de calcul: $\Pr(B | A) = \Pr(A \cap B) / \Pr(A)$

Ex.

1) la probabilité d'avoir un cancer colorectal sachant que le test Hémocult est positif, est une probabilité conditionnelle:

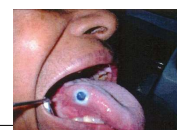
$\Pr(\text{cancer colorectal} | \text{test Hémocult Positif})$

2) la probabilité d'avoir un cancer oral sachant que le test de bleuissement avec toluidine est positif

$\Pr(\text{cancer oral} | \text{bleuissement avec toluidine Positif})$

3) la probabilité d'un homme d'avoir

$\text{PAS} > 140 \text{ mmHg} : \Pr(\text{PAS} > 140 | \text{Homme})$



Rajmohan M, Rao UK, Joshua E, Rajasekharan ST, Kannan R. Assessment of oral mucosa in normal, precancer and cancer using chemiluminescent illumination, toluidine blue supravital staining and oral exfoliative cytology. J Oral Maxillofac Pathol. 2004;56(1):42-5. <http://dx.doi.org/10.4103/0973-6229.104706>.

30

Probabilité conditionnelle et l'indépendance

□ A , B évènements **indépendantes**:

$$\Pr(B | A) = \Pr(B) = \Pr(B | \text{non } A)$$

□ A et B évènements **dépendantes**:

$$\Pr(B | A) \neq \Pr(B) \neq \Pr(B | \text{non } A)$$

$$\Pr(A \cap B) \neq \Pr(A) \times \Pr(B)$$

31

Probabilité conditionnelle - applications

Ex.: enquête des possibles facteurs de risque du cancer de poumon

Dans un échantillon de 2000 sujets, on a 1000 fumeurs parmi lesquelles 130 sujets souffrent de cancer du poumon et 1000 non fumeurs parmi lesquelles 10 sujets souffrent de cancer du poumon

Le risque relatif d'avoir le cancer de poumon = ?

Solution: on considère les événements

$A = \{\text{sujet fumeur}\}$ et $B = \{\text{sujet atteints de cancer du poumon}\}$

$$\Pr(B | A) = 130/1000 = 0,130 \qquad \Pr(B | \bar{A}) = 10/1000 = 0,010$$

$$RR = \frac{\Pr(B | A)}{\Pr(B | \bar{A})} = \frac{0,130}{0,010} = 13$$

=> un sujet qui est fumeur a un risque d'avoir le cancer du poumon de 13 fois plus grand qu'un sujet qui n'est pas fumeur

32

Probabilité conditionnelle - applications

	M+	M-	Total
T+	a	b	a+b
T-	c	d	c+d
Total	a+c	b+d	n

$$Se = \Pr(T^+ \mid M^+) = a / (a + c)$$

$$Sp = \Pr(T^- \mid M^-) = d / (b + d)$$

$$VPP = \Pr(M^+ \mid T^+) = a / (a + b)$$

$$VPN = \Pr(M^- \mid T^-) = d / (c + d)$$

33

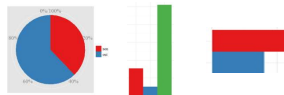
Le choix du type du graphique en fonction des types des variables et but

- Pour faire la **choix**, comptez **combien des variables** sont et quel **est le type**.

- **Description d'une seule variable**

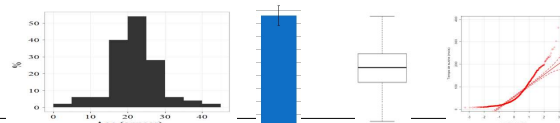
- *Qualitative*

- camembert (sectoriel – **Pie**)
- **Column** (si les noms des catégories ne sont pas très longues)
- **Bar** (si les noms des catégories sont très longues)



- *Quantitative*

- **Histogramme**, graphique des moyennes, box aux whiskers (boite à moustaches), graphique des quantiles



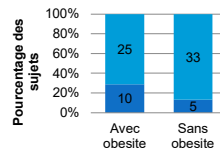
34

Le choix du type du graphique en fonction des types des variables et but

- **La relation entre deux variables**

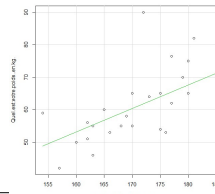
- **Qualitative**

- **Column** (Clustered Column/ Stacked Column/ 100% Stacked column), ou **Bar** (Clustered Bar / Stacked Bar / 100% Stacked Bar)



- **Quantitative**

- **Scatter** (nuage des points)



35

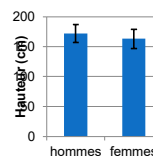
Le choix du type du graphique en fonction des types des variables et but

- **La relation entre deux variables**

- **Une variable quantitative en fonction d'une variable qualitative**

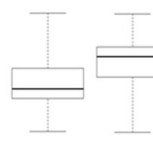
- **Si les données sont normale distribuées**

- Graphique des moyennes (avec deviation standard)
 - Graphique Colonnes avec barre d' erreur



- **Si les données sont normale distribuées**

- Graphique box plot ou whiskers/boite a moustaches (boite a moustaches)

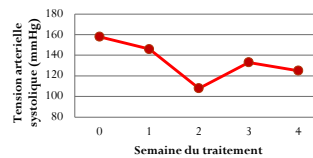


36

Le choix du type du graphique en fonction des types des variables et but

- L'évolution dans le temps d'une variable qualitative ou quantitative

- Line (Ligne)



- La relation entre trois variables quantitatives

- Bubble (nouage des sphères)
- Nouage des points tridimensionnel

- Une variable qualitative en fonction des intervalles d'une variable quantitative

- Area (Surface)

37

Mesures de symétrie: (skewness)

Coefficient d'asymétrie (α_3):

degré d'asymétrie d'une distribution

la direction de cette asymétrie (positive ou négative);

$\alpha_3 \approx 0 \Rightarrow$ une distribution symétrique.

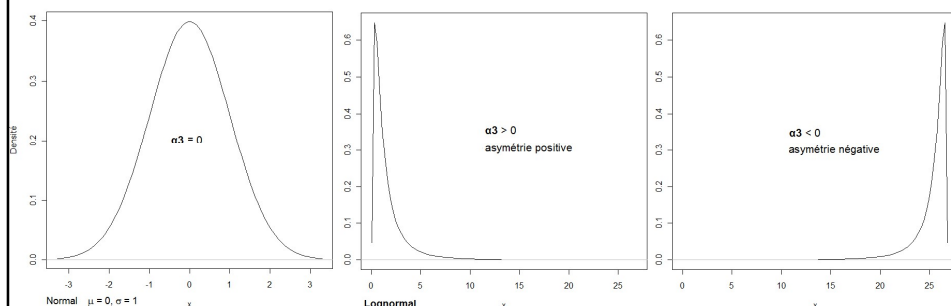
$\alpha_3 > 0 \Rightarrow$ distribution est plus allongée vers la droite – asymétrie positive

$\alpha_3 < 0 \Rightarrow$ distribution est plus allongée vers la gauche – asymétrie négative

$\alpha_3 (-0,5 - 0,5)$ approximative symétrique

$\alpha_3 (-1 - -0,5)$ ou $(0,5 - 1)$ modérée asymétrique

$\alpha_3 < -1$ ou > 1 asymétrie importante



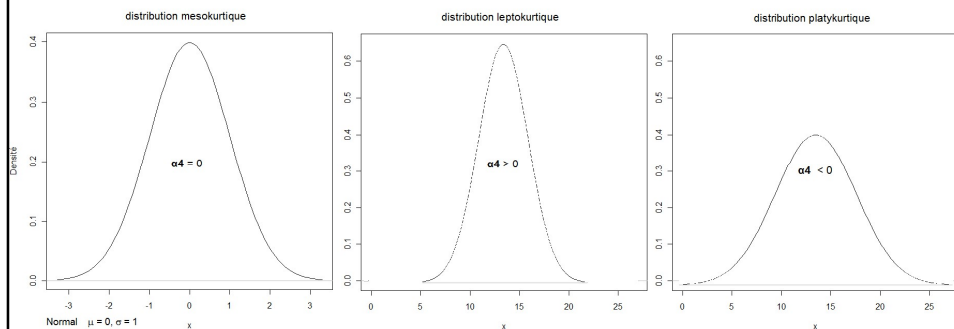
38

Le coefficient d'aplatissement (Kurtosis)

Le coefficient d'aplatissement (α_4):

l'angle de la courbe du milieu d'une distribution par rapport à une distribution normale (gaussienne)

- $\alpha_4 \approx 0 \Rightarrow$ l'angle normal \Rightarrow distribution mesokurtique
- $\alpha_4 > 0 \Rightarrow$ l'angle aigu \Rightarrow distribution leptokurtique - centre élevée
- $\alpha_4 < 0 \Rightarrow$ la pente aplati \Rightarrow distribution platykurtique - centre plus bas



39

Statistique descriptive

Mesures de tendance centrale:

- ✓ Moyenne
- ✓ Médiane
- ✓ Mode

Mesures de dispersion:

- ✓ Amplitude (entendue)
- ✓ Intervalle interquartile
- ✓ moyenne des écarts de la moyenne
- ✓ moyenne des écarts de la médiane
- ✓ Variance
- ✓ Déviation standard (écart-type)
- ✓ Coefficient de variation
- ✓ Erreur standard

Mesures de symétrie/aplatissement:

- ✓ Coefficient d'asymétrie (skewness)
- ✓ Coefficient d'aplatissement (Kurtosis)

Mesures de position:

- ✓ Quartiles
- ✓ Déciles
- ✓ Percentiles

40

Bon success!!!

41

41