

Revision

Comment choisir un test statistique

Questions à 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 théoriques < 5?
- Quelle est la nature des données?
 - Pour les **données quantitatives**:
 - **Distribution normale**?
 - **Variances égales/ inégales**?

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

Les hypothèses statistiques

- La création des hypothèses:
 - Certains **questions médicales** ont deux réponses opposées
 - on force beaucoup des questions dans ce format
 - Les réponses correspondent aux deux **modèles** possibles de la réalité
 - Ce 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 du 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 du 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 permet de faire la chois entre les deux possibilités (H_1 / H_0)

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

5

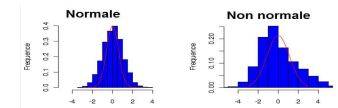
Vérification de la condition de normalité des données

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)
 - Boîte à moustaches (symétrique autour de la médiane)
 - Le graphique des quantiles (voir diapositive suivant)
- des statistiques descriptives (pas très fiables)
 - Si la moyenne est \approx 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)

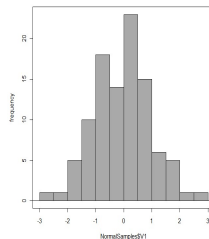


6

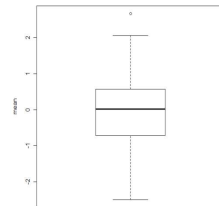
Comparaison des données normale/non normale distribuées

• **Normale**
 moyenne \approx 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
 $p = 0,99 > 0,05$

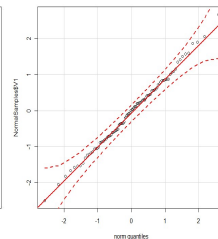
Histogramme



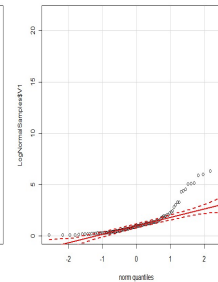
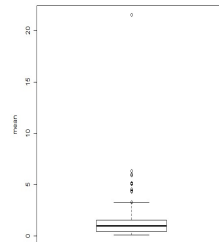
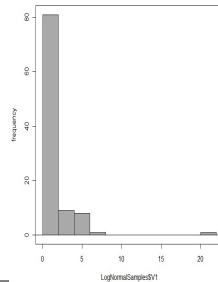
Boite à moustaches



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
 $p \approx 0 < 0,05$



La normalité des données

- Test de normalité des données :
 - Test de Kolmogorov-Smirnov
 - Si < 50 sujets le test Shapiro-Wilk
- H_0 = aucune différence statistiquement significative entre la distribution observée et la distribution normale
- H_1 = 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 normale distribuées (dans certaines conditions)

Récapitulatif des tests utilisés

Tests pour variables quantitatives - comparer la moyenne des deux échantillons

Type variable	Nb sujets	Nature des données	Statistique comparée	Test utilise	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.	$t = \frac{m_1 - m_2}{\sqrt{\frac{s^2}{n_1} + \frac{s^2}{n_2}}}$ $s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$	$(-\infty, -v.c.] \cup [v.c., +\infty)$
Deux échantillons dépendants (appariés)						
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;
 $S = \sqrt{\frac{\sum d^2}{n-1}}$; $s = \sqrt{\frac{\sum s^2}{n-1}}$; $S, \sigma, \sigma_1, \sigma_2$ - déviation standard dans la population; pour $\alpha=0,05, Z_{\alpha/2}=1,96$; d.d.l. - degrés de liberté; v.c. - valeur critique)

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 s sont similaires)	~médiane (si les distributions sont similaires) - la distribution des rangs /moyenne des rangs	Test Mann Whitney U Etapes: 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, $U_a = n_a \times n_b + n_b \times \frac{(n_b + 1)}{2} - obsSR_b$ $U_b = n_a \times n_b + n_a \times \frac{(n_a + 1)}{2} - obsSR_a$	Min (U _a , U _b) ≤ v.c.
Tests statistiques pour deux échantillons dépendants/ appariés					
Quantitative ou qualitative ordinale	Données non normale distribuées s sont similaires)	~médiane (si les distributions sont similaires) - la distribution des rangs	Test Wilcoxon pour échantillons appariés Etapes: 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.

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 *	$Z = \frac{p_1 - p_2}{\sqrt{p(1-p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$ $p = \frac{p_1 n_1 + p_2 n_2}{n_1 + n_2}$	$(-\infty, -Z_{\alpha/2}] \cup [Z_{\alpha/2}, +\infty)$
	<20% cellules du tableau théorique sont <5	fréquences		Test Chi carrée *	$\chi^2 = \sum_{i=1}^L \frac{(f_i^o - f_i^t)^2}{f_i^t}$	$[\chi_{v, \alpha}^2, +\infty)$
	>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	$\chi_{1, ad}^2 = \frac{(b-c)^2}{b+c}$	$[\chi_{v, \alpha}^2, +\infty)$

* 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

p_1, p_2 – fréquences; n_1, n_2 – 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

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	$\chi^2 = \sum_{i=1}^L \frac{(f_i^o - f_i^t)^2}{f_i^t}$	$[\chi_{v, \alpha}^2, +\infty)$
	>20% cellules du tableau théorique sont <5		fréquence	Test exact Fisher	-	-
Quantitative		Normale distribuées, Variance des échantillons égaux	moyenne	Test ANOVA	$F = \frac{MCG}{MCE}$	$[F_{v_1, v_2, \alpha}, +\infty)$
		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^o – fréquence observée, f^t – fréquence théorique; d.d.l. – degrés de liberté;

$MCEG = [n_1 * (m_1 - mt_1)^2 + n_2 * (m_2 - mt_2)^2 + n_3 * (m_3 - mt_3)^2 + ...] / (p-1)$
 $MCDG = [(n_1-1) * DS_1^2 + (n_2-1) * DS_2^2 + (n_3-1) * DS_3^2 + ...] / (n-p)$

n=nombre total d'observations, n_1, n_2, n_3 – le nombre d'observations par group, p – nombre des groups
 mt – la moyenne des toutes les observations, $m_1, m_2, m_3, ...$ les moyennes par groups, $DS_1, DS_2, DS_3, ...$ déviations standard d'échantillonnage des groups

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.	$F = \left[\begin{array}{l} \frac{S_2^2}{S_1^2},_{pour, S_2^2} S_1^2 \\ \frac{S_1^2}{S_2^2},_{pour, S_1^2} S_2^2 \end{array} \right]$	$\left[F_{v_1, v_2, \alpha}, +\infty \right)$
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)

Équivalences 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 signés (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

Intervalles de confiance				
Type variable	Nombre échantillons	Estimateur ponctuel	Conditions application	Formule
qualitative	une	fréquence: f	grands échantillons: $nf \geq 10$ et $nq \geq 10$	$\left(f - Z_{\frac{1-\alpha}{2}} ES; f + Z_{\frac{1-\alpha}{2}} ES\right) \quad ES = \sqrt{\frac{f(1-f)}{n}}$
	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: $f_1 n_1 \geq 10$, $(1-f_1)n_1 \geq 10$, $f_2 n_2 \geq 10$, $(1-f_2)n_2 \geq 10$	$ES = \sqrt{\frac{f_1 \times (1-f_1)}{n_1} + \frac{f_2 \times (1-f_2)}{n_2}}$ $\left((f_1 - f_2) - Z_{\frac{1-\alpha}{2}} ES; (f_1 - f_2) + Z_{\frac{1-\alpha}{2}} ES\right)$
quantitative	deux	différence entre les fréquences:	petits échantillons: $np < 10$ ou $nq < 10$	on va pas calculée
	un	moyenne: m	grands échantillons: $n \geq 30$, 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: $n \geq 30$, 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	$ES = \sqrt{\frac{(n_1 - 1) \times s_1^2 + (n_2 - 1) \times s_2^2}{n_1 + n_2 - 2}}$ $\left((m_1 - m_2) - t_{n_1 + n_2 - 2, \frac{1-\alpha}{2}} ES; (m_1 - m_2) + t_{n_1 + n_2 - 2, \frac{1-\alpha}{2}} ES\right)$

15

(ou n_1, n_2 - nombre des sujets; f, f_1, f_2 - fréquence observée; $q = 1 - f$; m, m_1, m_2 - moyennes; s, s_1, s_2 - déviations standard descriptive de l'échantillon; S - déviation standard d'échantillonnage; $S = \sqrt{\frac{1}{n-1} \sum (x_i - m)^2}$; $s = \sqrt{\frac{1}{n} \sum (x_i - m)^2}$; σ - déviation standard dans la population; ES=erreur standard; pour $\alpha=0.05$, $Z_{\alpha/2}=1.96$)

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

Evaluer la linearite, et l'importance de la correlation:

Si les points semble suggérer un droite – la relation est peut être linéaire

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

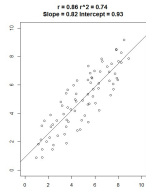
Si la relation est plus probable linéaire, on peut évaluer d'une manière subjective la corrélation linéaire. Plus les points se rapprochent à un droite de tendance, plus la corrélation est forte, plus les points sont distants, plus la corrélation est faible

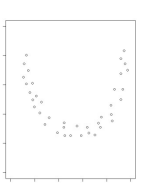
1.0

0.8

0.4

0.0





16

16

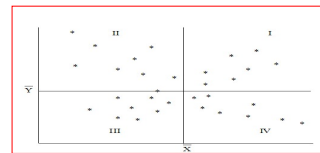
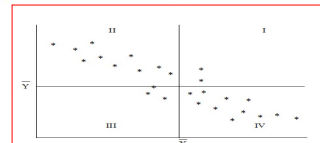
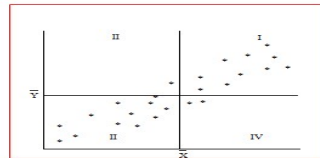
Evaluation graphique : diagramme de dispersion

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



17

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_{x_i} et R_{y_i} 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

18

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 de 0, plus la relation est faible

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/lien/relation peut être utilisé**)

(-0.25 et 0,25)

\Rightarrow une relation **négligeable** ou **aucune** corrélation linéaire entre les variables

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

\Rightarrow un degré de corrélation **faible/acceptable**

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

\Rightarrow un degré de corrélation **modérée à bonne**

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

\Rightarrow 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.

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 à 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

21

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)

22

Variables aléatoires discrètes

Esperance, variance et écart type

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

$$E(X) = \sum_{i=1}^n x_i \Pr(X = x_i)$$

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

Exemple: avec le traitement

- X: 0 1 2 3 4

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

$$E(X) = 0 \times 0,008 + 1 \times 0,076 + 2 \times 0,265 + 3 \times 0,411 + 4 \times 0,24 = 1,31$$

- **Variance** de X:

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

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

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

23

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.		$\Pr(A \cap B) = \Pr(A) \times \Pr(B)$

24

Probabilité conditionnelle

- A et B évènements, $P(A) \neq 0$.

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

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



Rajeshan M, Rao UK, Joshi E, Rajasekhar ST, Kumar 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. 2012;66(1):125-9. <http://dx.doi.org/10.4103/0972-405X.105070>.

Probabilité conditionnelle et l'indépendance

□ A, B évènements **indépendantes**:

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

□ A et B évènements **dépendantes**:

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

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

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

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

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

=>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

Probabilité conditionnelle - applications

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

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

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

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

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

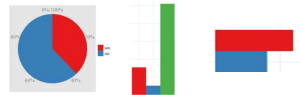
Le choix du type du graphique en fonction des types des variables et but

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

- **Description d'une seule variable**

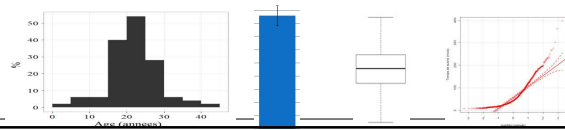
- **Qualitative**

- 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 and whiskers** (boite à moustaches), graphique des quantiles

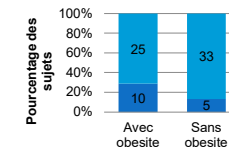


Le choix du type du graphique en fonction des types des variables et but

- **La relation entre deux variables**

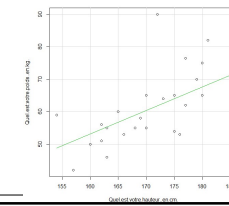
- **Qualitative**

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



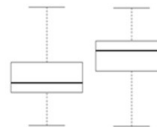
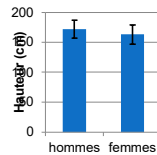
- **Quantitative**

- **Scatter** (nuage des points)



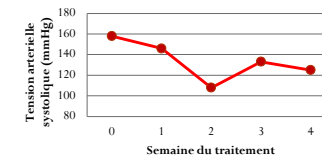
Le choix du type du graphique en fonction des types des variables et but

- **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)



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)



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.

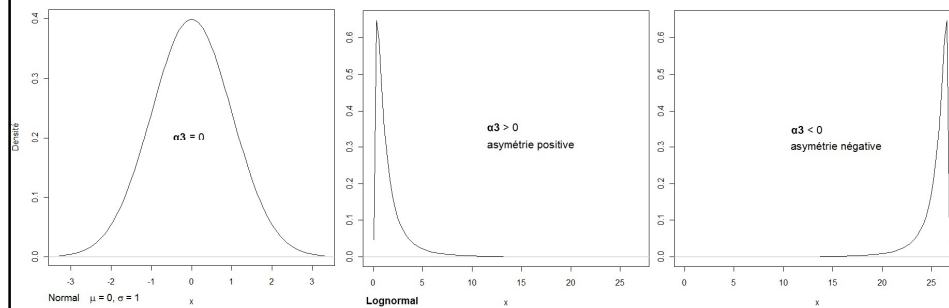
$\alpha_3 > 0 \Rightarrow$ distribution est plus allongée vers la droite – asymétrie positive

$\alpha_3 < 0 \Rightarrow$ 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



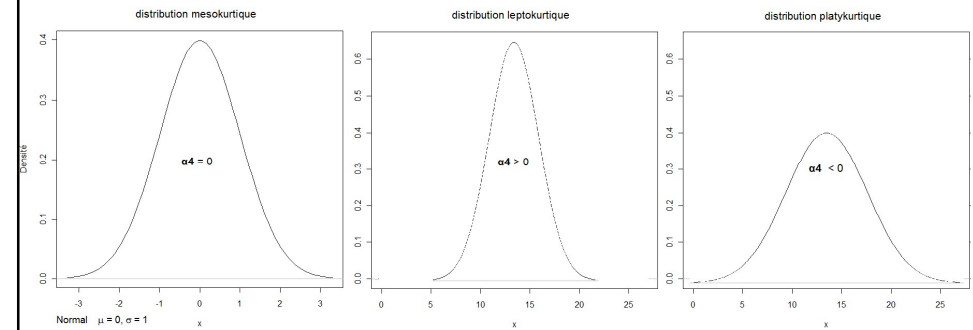
Le coefficient d'aplatissement (Kurtosis)

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 aplaté \Rightarrow distribution platykurtique – centre plus bas



Statistique descriptive

Mesures de tendance centrale:

- ✓ Moyenne
- ✓ Médiane
- ✓ Mode

Mesures de symétrie/aplatissement:

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

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 position:

- ✓ Quartiles
- ✓ Déciles
- ✓ Percentiles

Bon success!!!