

## 3ème partie: STATISTIQUE A 2 VARIABLES

Il va s'agir de la description et de la quantification de la relation entre les valeurs de 2 caractères X et Y simultanément observés sur les individus d'une population définie. Les caractères X et Y peuvent être de nature qualitative ou quantitative.

### I : Caractères qualitatifs

#### 1°) test d'indépendance

Si les caractères sont qualitatifs, les résultats sont présentes dans un tableau à double entrée noté *tableau de contingence*.

X/Y	y1	y2	yj	yk	total
x1	n <sub>11</sub>	n <sub>12</sub>	n <sub>1j</sub>	n <sub>1k</sub>	n <sub>1.</sub>
x2	n <sub>21</sub>	n <sub>22</sub>	n <sub>2j</sub>	n <sub>2k</sub>	n <sub>2.</sub>
.	.	.	.	.	.
.	.	.	.	.	.
x <sub>i</sub>	n <sub>i1</sub>	n <sub>i2</sub>	n <sub>ij</sub>	n <sub>ik</sub>	n <sub>i.</sub>
.	.	.	.	.	.
.	.	.	.	.	.
x <sub>l</sub>	n <sub>l1</sub>	n <sub>l2</sub>	n <sub>lj</sub>	n <sub>lk</sub>	n <sub>l.</sub>
total	n <sub>.1</sub>	n <sub>.2</sub>	n <sub>.j</sub>	n <sub>.k</sub>	n <sub>..</sub>

Les distributions marginales donnent pour l'ensemble des modalités d'une variable, la répartition de l'effectif total selon les modalités de l'autre variable.

Si X et Y sont **indépendants**, alors chaque n<sub>ij</sub> devrait être égal à :

$$n'_{ij} = (n_{i.} \cdot n_{.j}) / n_{..}$$

On établit le tableau des effectifs théoriques selon le modèle de l'indépendance et on effectuera le **test de l'indépendance** qui consiste à définir une règle de décision concernant la validité de l'hypothèse relative à l'indépendance des 2 caractères X et Y .

hypothèses :

Hypothèse nulle d'indépendance H0 : X et Y sont indépendants

Hypothèse alternative H1 : X et Y sont liés

critère statistique calculé:

$$X^2_{\text{calc}} = (n_{ij} - n'_{ij})^2 / n'_{ij}$$

critère statistique théorique pour un niveau de probabilité (1 - α ) choisi **à priori** : **X<sup>2</sup> théo**

interprétation du test:

**$X^2 \text{ calc} < X^2 \text{ théo}$**  on ne peut pas mettre en évidence l'existence d'un lien, on garde  $H_0$

**$X^2 \text{ calc} > X^2 \text{ théo}$**  on rejette  $H_0$  avec un risque  $\alpha$  de le faire à tort. X et Y sont liés.

Il existe une autre approche si  $\alpha$  n'est pas fixé.

**conditions d'utilisation:  $n.. \geq 20$  et chaque  $n'_{ij}$  calculé  $\geq 5$**

2°) test d'adéquation

Il s'agit d'une autre utilisation du critère  $X^2$ . On parle de test d'ajustement, il permet de comparer la forme de la distribution de 2 séries statistiques entre elles.

## II :Caractères quantitatifs

On a souvent besoin d'examiner la façon dont une variable quantitative est reliée à une autre variable quantitative, il s'agit de la **régression**.

On étudiera le cas où Y est relié linéairement à une seule variable X, il s'agit de régression simple. Y est la **variable dépendante** à expliquer ou variable de réponse et X est la **variable explicative** ou **variable indépendante** ou encore **régresseur**.

1°) la covariance

La covariance est la moyenne du produit des écarts de 2 variables statistiques à leurs moyennes respectives.

formule de définition:  $\text{cov}(x, y) = \text{SPE}(x, y) / n$

formule de calcul manuel:  $\text{Cov}(x, y) = E(XY) - E(X) \cdot E(Y)$

La covariance peut être positive, négative ou nulle.

2°) nuage de points

L'observation de l'allure du nuage de points nous donne déjà un certain nombre de renseignements.

3°) coefficients de la régression

Les coefficients de la régression sont déterminés par la méthode des moindres carrés ordinaires notée MCO.

$$a = \bar{y} - b \bar{x}$$

$$b = \text{SPE}(x, y) / \text{SCE } x \quad b = \text{cov}(x, y) / \text{var } x$$

Les coefficients a et b sont les coefficients de régression, ils s'interprètent.

La droite de régression passe nécessairement par le centre de gravité du nuage  $(\bar{y} ; \bar{x})$ .

La droite de régression de y en fonction de x (droite de y pour x fixé;  $Dy/x$ ) a pour équation:

$$\hat{y}_i = b x_i + a \quad \text{ou encore} \quad y_i = b x_i + a + e_i \quad \text{avec } \sum e_i = 0$$

## 4°) analyse de la qualité de la régression

☺ coefficient de détermination  $r^2$ 

Il mesure la part de la variance totale qui est expliquée par la régression.

Il représente le pourcentage des variations de Y expliquées par les variations de X.

variance totale des y = variance expliquée par la régression + variance résiduelle

$$s^2(y) = s^2(\hat{y}) + s^2(e)$$

$$r^2 = s^2(\hat{y}) / s^2(y)$$

**mais c'est aussi :**

$$r^2 = \text{cov}(x,y)^2 / (s^2x * s^2y)$$

☺ coefficient de corrélation  $r$ 

Le carré du coefficient de corrélation est égal au coefficient de détermination.

## ☺ analyse des résidus

*normalité des résidus*

*graphe des résidus standardisés*

*indépendance des résidus*