#### People's Democratic Republic Of Algeria Ministry Of Higher Education And Scientific Research



Centtre Univercitaire Abdelhafid Boussouf Mila



# Rappels sur les statistiques descriptive

L'objectif du cours est d'étudier sur une popolution de N individus, deux variables différentes et de chercher s'il y a un lien ou une corrélation entre ces deux caractères. Chacune des deux variables peut être, soit quantitative, soit qualitative. On étudie deux cas :

- Les deux variables sont qualitatives.
- Les deux variables sont quantitatives.

Exemples de relations possibles entre deux charactères : Taille et poids des individus, Effet et dosage, tabagisme et cancers du poumon, rendement et quantité d'engrais utilisée, ... .

• Variables qualitatives .

- Variables qualitatives .
- Données observées-Tableau de contingence .

- Variables qualitatives .
- Données observées-Tableau de contingence .
- Tableau des fréquences .

- Variables qualitatives .
- Données observées-Tableau de contingence .
- Tableau des fréquences .
- Variables quantitatives .

- Variables qualitatives .
- Données observées-Tableau de contingence .
- Tableau des fréquences .
- Variables quantitatives .
- Paramètres de position et de dispersion .

#### Variables qualitatives

Données observées-Tableau de contingence On considère une population statistique décrite selon deux caractères  $X=(x_1,x_2,...,x_k)$  et  $Y=(y_1,y_2,...,y_p)$ : Leurs valeurs distinctes sont notées respectivement :

## Données observées-Tableau de contingence

Les données observées peuvent être regroupées sous la forme d'un tableau de contingence :

	<i>y</i> <sub>1</sub>	<i>y</i> <sub>2</sub>	 Уj	 Уp	n <sub>i</sub> .
<i>x</i> <sub>1</sub>	n <sub>11</sub>	n <sub>12</sub>	 $n_{1j}$	 $n_{1p}$	<i>n</i> <sub>1.</sub>
<i>x</i> <sub>2</sub>	n <sub>21</sub>	n <sub>22</sub>	 $n_{2j}$	 n <sub>2p</sub>	n <sub>2</sub>
:					
x <sub>i</sub>	n <sub>i1</sub>	n <sub>i2</sub>	 n <sub>ij</sub>	 n <sub>ip</sub>	n <sub>i</sub> .
:					
X <sub>k</sub>	$n_{k1}$	n <sub>k2</sub>	 n <sub>kj</sub>	 n <sub>kp</sub>	n <sub>k</sub>
n.j	n. <sub>1</sub>	n.2	 n.j	 n. <sub>p</sub>	$N = n_{}$

## Données observées-Tableau de contingence

 $-n_{ij}$  représente le nombre de fois que les modalités  $x_i$  et  $y_j$  apparaissent ensemble. - N est l'effectif total de la population :

$$N = n.. = \sum_{i=1}^{k} \sum_{j=1}^{p} n_{ij}.$$

- Les  $n_i$  et  $n_{ij}$  sont appelés les effectifs marginaux.
- $n_i$ . représente le nombre de fois que la modalité  $x_i$  apparaît,

$$\sum_{j=1}^{p} n_{ij} = n_{i1} + n_{i2} + \ldots + n_{ip} = n_{i.}, \text{ pour tou } i = 1, \ldots, k$$

-  $n_{.j}$  représente le nombre de fois que la modalité  $y_j$  apparaît,

$$\sum_{i=1}^{k} n_{ij} = n_{1j} + n_{2j} + \ldots + n_{kj} = n_{.j} \text{ pour tou } j = 1, \ldots, p$$

$$\sum_{i=1}^{k} n_{i} = \sum_{i=1}^{p} n_{i} = N = n \dots$$

## Tableau des fréquences-Tableau de contingence

Fréquences partielles :

$$f_{ij}=rac{n_{ij}}{n_{..}},\quad i=1,\ldots k,\quad j=1,\ldots,p,$$

Fréquences marginales : La fréquence marginale de la variable  $\boldsymbol{X}$  est donnée par :

$$f_{i.}=\frac{n_{i.}}{n}, \quad i=1,\ldots k,$$

#### Données observées-Tableau de fréquences :

	<i>y</i> <sub>1</sub>	<i>y</i> <sub>2</sub>	 Уј	 Ур	f <sub>i</sub> .
<i>x</i> <sub>1</sub>	$f_{11}$	$f_{12}$	 $f_{1j}$	 $f_{1p}$	$f_1$ .
<i>x</i> <sub>2</sub>	$f_{21}$	$f_{22}$	 $f_{2j}$	 $f_{2p}$	<i>f</i> <sub>2.</sub>
:					
Xi	$f_{i1}$	f <sub>i2</sub>	 $f_{ij}$	 $f_{ip}$	f <sub>i</sub> .
:					
X <sub>k</sub>	$f_{k1}$	$f_{k2}$	 $f_{kj}$	 $f_{kp}$	$f_k$ .
f.;	f. <sub>1</sub>	f. <sub>2</sub>	 f.;	 f. <sub>p</sub>	1

La fréquence marginale de la variable Y est donnée par :

$$f_{\cdot j} = \frac{n_{\cdot j}}{n_{\cdot \cdot}}, \quad j = 1, \dots p,$$

$$\sum_{i=1}^{k} \sum_{i=1}^{p} f_{ij} = 1, \quad \sum_{i=1}^{k} f_{i} = 1, \quad \sum_{i=1}^{p} f_{ij} = 1.$$

## Fréquences conditionnelles :

La fréquence conditionnelle de la variable X par rapport à  $Y_j, j=1,\ldots,p$  est donnée par :

$$f_i^j = \frac{n_{ij}}{n_{i,j}}, \quad j = 1, \dots p,$$

La fréquence conditionnelle de la variable Y par rapport à  $X_i, i=1,\ldots,k$  est donnée par :

$$f_j^i = \frac{n_{ij}}{n_{i.}}, \quad i = 1, \dots k$$

$$f_{ij} = f_i \cdot \times f_j^i = f_{\cdot j} \times f_i^j$$
.

# Relations entre fréquences marginales et fréquences conditionnelles :

Exemple. On s'intéresse à une éventuelle relation entre les caractères : "être fumeur" (plus de 20 cigarettes par jour, pendant 10 ans) et "avoir un cancer de la gorge", sur une population de 1000 personnes, dont 500 sont atteintes d'un cancer de la gorge. Les résultats observés sont présentés

dans le tableau suivant :

Observé	cancer	non cancer	
fumeur	342	258	
non fumeur	158	242	

# Relations entre fréquences marginales et fréquences conditionnelles :

On calcule l'effectif total :

On calcule i effectif total :								
Observé	cancer	non cancer	Total					
fumeur	342	258	600					
non fumeur	158	242	400					
Total	500	500	200					

En utilisant les relations

précédentes, on obtient le tableau des fréquences (fréquences marginales) :

Observé	cancer	non cancer	Total
fumeur	0.342	0.258	0.6
non fumeur	0.158	0.242	0.4
Total	0.5	0.5	1

Quelques exemples :

Fréquences marginales :

$$f_{11} = \frac{n_{11}}{n} = \frac{342}{1000} = 0.342, f_{1.} = \frac{n_1}{n} = \frac{600}{1000} = 0.6$$

# Relations entre fréquences marginales et fréquences conditionnelles :

Fréquences conditionnelles de X sachant Y :

$$f_1^2 = \frac{n_{12}}{n.2} = \frac{258}{500} = 0.516, f_2^2 = \frac{n_{22}}{n.2} = \frac{242}{500} = 0.484$$

Fréquences conditionnelles de Y sachant X :

$$f_2^1 = \frac{n_{12}}{n_1} = \frac{258}{600} = 0.43, f_2^2 = \frac{n_{22}}{n_2} = \frac{242}{400} = 0.605$$

## Paramètres de position et de dispersion :

Dans le cas où X et Y sont des variables quantitatives, on peut associer à chacune des séries marginales définies par le tableau de contingence des caractéristiques de tendance centrale et de dispersion. On considère un tableau de contingence comme celui défini précédemment. Moyenne marginale de X. La moyenne marginale de X, notée  $\overline{\bar{x}}$ :

$$\overline{\overline{x}} = \frac{1}{N} \sum_{i=1}^k n_i \cdot x_i = \sum_{i=1}^k f_i x_i.$$

Moyenne marginale de Y. La moyenne marginale de Y, notée  $\overline{\bar{y}}$  :

$$\overline{\bar{y}} = \frac{1}{N} \sum_{j=1}^{P} n_{-j} y_i = \sum_{j=1}^{P} f_{\cdot j} y_j.$$

## Paramètres de position et de dispersion :

Variance marginale de X. La variance marginale de X, notée V(X) :

$$V(X) = \frac{1}{N} \sum_{i=1}^{k} n_{i} \cdot (x_{i} - \overline{\bar{x}}) = \frac{1}{N} \sum_{i=1}^{k} n_{i} \cdot x_{i}^{2} - \overline{\bar{x}}^{2}.$$

L'écart type marginale de X :  $\sigma_X = \sqrt{V(X)}$ . Varaince marginale de Y. La moyenne marginale de Y, notée V(Y) :

$$V(Y) = \frac{1}{N} \sum_{j=1}^{p} n_{.j} (y_j - \overline{\bar{y}}) = \frac{1}{N} \sum_{j=1}^{p} n_{-j} y_j^2 - \overline{\bar{y}}^2.$$

L'écart type marginale de  $Y: \sigma_Y = \sqrt{V(Y)}$  Covariance. La covriance est la mesure de la variation simultanée de deux variables X et Y.

$$Cov(X, Y) = \frac{1}{N} \sum_{i=1}^{k} \sum_{j=1}^{p} n_{ij} (x_i - \overline{\bar{x}}) (y_j - \overline{\bar{y}}) = \frac{1}{N} \sum_{j=1}^{k} \sum_{j=1}^{p} n_{ij} x_i y_j - \overline{\bar{x}} \overline{\bar{y}}$$

#### Propriétés:

- Une covariance peut être positive, négative ou nulle.
- Cov(X, Y) = Cov(Y, X)
- Cov(aX, Y) = aCov(X, Y) = acov(Y, X)
- Cov(X, X) = Var(X)
- Var(X + Y) = Var(X) + Var(Y) + 2 cov(X, Y)

#### **Exemple**

Le Tableau suivant présente la répartition d'un ensemble de consommateurs selon leurs revenus et leurs dépenses :

X/Y	25	30	35	40	Total
20	4	2	1	0	7
25	5	1	0	0	6
30	3	2	1	1	7
Total	12	5	2	1	20

 $\boldsymbol{X}$  représente les revenus, et  $\boldsymbol{Y}$ , les

dépenses de consommation. Les revenus moyens :

$$\overline{\overline{x}} = \frac{7 \times 20 + 6 \times 25 + 7 \times 30}{20} = 25.$$

La variance des revenus :

$$V(X) = \frac{1}{20} \left( 7 \times 20^2 + 6 \times 25^2 + 7 \times 30^2 \right) - 25^2 = 17.5.$$

Les dépenses de consommation qu'effectuent en moyenne chacun des consommateurs :

$$\overline{\bar{y}} = \frac{12 \times 25 + 5 \times 30 + 2 \times 35 + 1 \times 40}{20} = 28$$

#### **Exemple**

Les dépenses de consommation qu'effectuent en moyenne chacun des consommateurs :

$$\overline{\overline{y}} = \frac{12 \times 25 + 5 \times 30 + 2 \times 35 + 1 \times 40}{20} = 28.$$

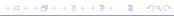
La variance des dépenses de consommation :

$$V(Y) = \frac{1}{20} \left( 12 \times 25^2 + 5 \times 30^2 + 2 \times 35^2 + 1 \times 40^2 \right) - 28^2 = 18.5.$$

La covariance de X et Y.

X/Y	25	30	35	40	Total	$\left  \sum_{i=1}^k \sum_{j=1}^p n_{ij} x_i y_j \right $
20	4	2	1	0	7	3900
25	5	1	0	0	6	3875
30	3	2	1	1	7	6300
Total	12	5	2	1	20	14075

$$Cov(X, Y) = \frac{14075}{20} - 25 \times 28 = 3.75$$



# Fin du partie 1