

Tableau de contingence

y_i
 x_i n_{ij} (effectif partiel)

Distribution marginales :

$$n_{i.} = \sum_{j=1}^q n_{ij}$$

effectifs marginaux

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

Exemple : On considère une distribution avec 10 000 jeunes salariés en fonction de l'âge, variable x , et selon le salaire mensuel noté y et en milliers d'euros.

On a le tableau de contingence suivant:

x	y	[0,8-1,2[[1,2-1,6[[1,6-2,0[Σ (distributions marginales de x)
19-21		700	300	100	1100 (n1.)
21-23		1400	2000	300	3700
23-25		1000	2900	1300	5200
Σ (distributions marginales de y)		3100	5200	1700	n=10 000

On a : $n_{12} = 300 \rightarrow$ jeunes qui ont entre 19 et 21 ans et dont le salaire est compris [1,2-1,6[

...

$$n_{21} = 1400$$

$$n_{32} = 2900$$

(effectifs partiels)

On a aussi les fréquences partielles :

$$f_{32} = \frac{n_{32}}{n} = \frac{2900}{10\,000} = 29\%$$

Distribution marginale de x :

$n_{1.} = 1100$ nombre de jeunes qui ont entre 19 et 21 ans

$n_{2.} = 3700$

$n_{3.} = 5200$

Fréquence marginale :

$$f_{1.} = \frac{n_{1.}}{n} = \frac{1100}{10\,000}$$

Distribution marginale de y :

$n_{.1} = 3100$ nombre de jeunes dont le salaire est compris entre 0,8 et 1,2 milliers d'euros

$n_{.2} = 5200$

$n_{.3} = 1700$

$$f_{.1} = \frac{3100}{10\,000} = 31\%$$

Remarques :

$$\sum_{i=1}^p \sum_{j=1}^p f_{ij} = 1$$

$$\text{car } f_{ij} = \frac{n_{ij}}{n}$$

$$\sum_{i=1}^p f_{i.} = 1$$

$$\text{car } \sum_{i=1}^p n_{i.} = n$$

$$\sum_{j=1}^p f.j = 1$$

$$\text{car } \sum_{j=1}^q n.j = n$$

c. Distribution conditionnelles :

Ce sont les distributions d'un caractère pour une modalité fixée de l'autre.

Par exemple :

- distribution conditionnelle de x pour y=y_j, ce sont (x_i, effectif n_{ij}) qui correspondent à la "jème" colonne du tableau.
- distribution conditionnelle de y pour x=x_j, ce sont (y_i, effectif n_{ij}) qui correspondent à la "ième" ligne du tableau.
- Fréquence conditionnelle :

$$f_{i/j} = \frac{n_{ij}}{n.j}$$

n_{ij} : nombre d'unités ayant le caractère (x_i, y_j)

n.j : nombre d'unités ayant le caractère y_j

et aussi :

$$f_{j/i} = \frac{n_{ij}}{n.i}$$

On a f_{i/j} = proportion d'unités ayant x_i parmi ceux ayant y_j
= fréquence conditionnelle de x selon y

Exemple des salariés :

$$f_{1/2} = \frac{n_{12}}{n.2} = \frac{300}{5200} = 5,8\%$$

58% de ceux qui gagnent entre 1,2 et 1,6 milliers d'euros ont entre 19 et 21 ans

Propriété :

$$f_{i/j} = \frac{f_{ij}}{f.j} \text{ ou } f_{j/i} = \frac{f_{ij}}{f.i}$$

Remarque:

$$\sum_{i=1}^p f_{i/j} = 1 \text{ et } \sum_{j=1}^q f_{j/i} = 1$$

2. Caractéristiques des séries à 2 caractères :

a. Caractéristiques des séries marginales : (x_i, n_{i.}) et (y_j, n._j)

Moyenne de x :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i n_i. = \frac{1}{n} \sum_{j=1}^q x_i n_{ij}$$

Moyenne de y :

$$\bar{y} = \frac{1}{n} \sum_{j=1}^q y_j n.j = \frac{1}{n} \sum_{j=1}^q \sum_{i=1}^p y_j n_{ij}$$

Exemple : tableau des 10 000 jeunes salariés

$$\bar{x} = \frac{1}{10\,000} (20 * 1100 + 22 * 3700 + 24 * 5200) = 22,82$$

20, 22 et 24 : centres des classes

$$\bar{y} = \frac{1}{10\,000} (1,0 * 3100 + 1,4 * 5200 + 1,8 * 1700) = 1,344$$

Variance de x :

$$\sigma x^2 = \frac{1}{n} \sum_{i=1}^p n_i. (x_i - \bar{x})^2 = \bar{x}^2 - (\bar{x})^2$$

$$\sigma y^2 = \frac{1}{n} \sum_{j=1}^q n.j (y_j - \bar{y})^2 = \overline{y^2} - (\bar{y})^2$$

Exemple : tableau des 10 000 jeunes salariés

$$\overline{x^2} = \frac{1}{10\,000} (20^2 * 1100 + 22^2 * 3700 + 24^2 * 5200) = 522,6$$

d'où

$$\sigma x^2 = 522,6 - (22,82)^2 = 1,85$$

$$\sqrt{\sigma x^2} = \sqrt{1,85} = 1,36$$

$$\overline{y^2} = \frac{1}{10\,000} (1,0^2 * 3100 + 1,4^2 * 5200 + 1,8^2 * 1700) = 1,88$$

d'où

$$\sigma y = 0,27$$

b. Caractéristiques des distributions conditionnelles de x si y_j et de y si x_i :

Moyennes :

de x si y_j, $\bar{x}_j = \frac{1}{n.j} \sum_{i=1}^p x_i n_{ij}$ d'après "jème" colonne

de y si x_i, $\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^q y_j n_{ij}$ d'après "ième" colonne

Variance de x si y_j :

$$\sigma_j^2(x) = \frac{1}{n.j} \sum_{i=1}^p n_{ij} (x_i - \bar{x}_j)^2 = \overline{x_j^2} - (\bar{x}_j)^2$$

Variance de y si x_i :

$$\sigma_i^2(y) = \frac{1}{n_i} \sum_{j=1}^q n_{ij} (y_j - \bar{y}_i)^2 = \overline{y_i^2} - (\bar{y}_i)^2$$

⇒ Démonstration en TD

Exemples : tableau des salariés

$$\overline{x_1} = \frac{1}{3100} (20 * 700 + 22 * 1400 + 24 * 1000) = 22,194$$

Moyenne d'âge de ceux qui gagnent entre 0,8 et 1,2 milliers d'euros

$$\sigma_1^2(x) = \frac{1}{3100} (20^2 * 700 + 22^2 * 1400 + 24^2 * 1000) - (22,194)^2 = 2,15$$

$$\sigma_1(x) = 1,47$$

σ₁(x) représente la dispersion d'âge de ceux qui gagnent entre 0,8 et 1,2 milliers d'euros (j=1)

$$\overline{y_2} = \frac{1}{3700} (1,0 * 1400 + 1,4 * 200 + 1,8 * 300) = 1,281$$

i=200, x=x₂

Il s'agit du salaire moyen de ceux qui ont entre 21 et 23 ans.

$$\sigma_2^2(y) = \frac{1}{3700} (1,0^2 * 1400 + 1,4^2 * 200 + 1,8^2 * 300) - (1,281)^2 = 0,0596$$

$$\sigma_2(y) = 0,24$$

σ₂(y) représente la dispersion du salaire de ceux qui ont entre 21 et 23 ans

c. Relations entre les caractéristiques marginales et conditionnelles

Sur les moyennes :

On a calculé les \bar{x}_j pour y=y_j fixé et d'effectif n.j alors :

$$\bar{x} = \frac{1}{n} \sum_{j=1}^q \bar{x}_j n.j$$

Idem pour y

d. Covariance de x et y

Définition :

$$cov(x, y) = \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^q nij (xi - \bar{x})(yj - \bar{y})$$

Cette notion de covariance sert à savoir si les caractères x et y sont dépendants ou indépendants.

Propriété :

$$cov(x, y) = \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^q nij xiyj - \bar{x}\bar{y} = \bar{xy} - \bar{x}\bar{y}$$

donc :

$$cov(x, y) = \bar{xy} - \bar{x}\bar{y}$$

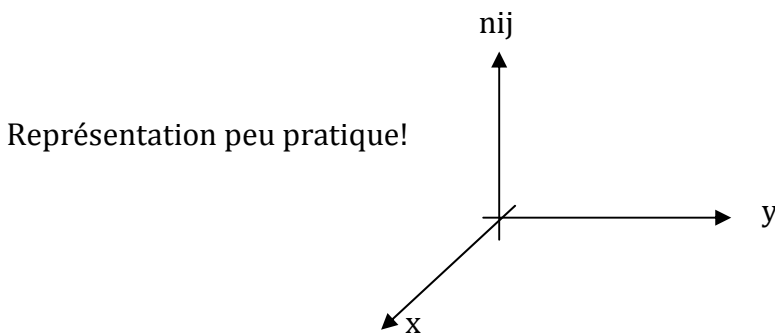
Remarque :

$$cov(x, x) = \overline{x^2} - (\bar{x})^2 = \sigma x^2$$

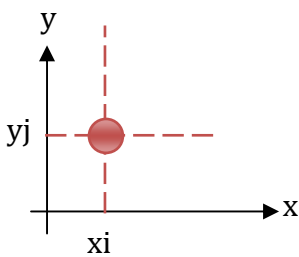
3. Représentation graphique des séries à 2 caractères :

On considère une série (xi, yj) d'effectifs nij.

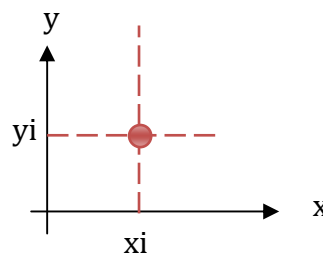
- Graphe à 3 dimensions:



- Graphe à 2 dimensions:



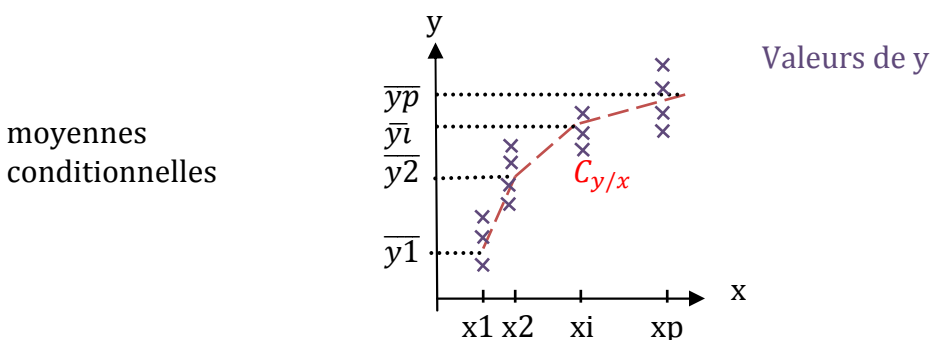
point dont la grosseur est proportionnelle à l'effectif nij



avec effectif plus faible

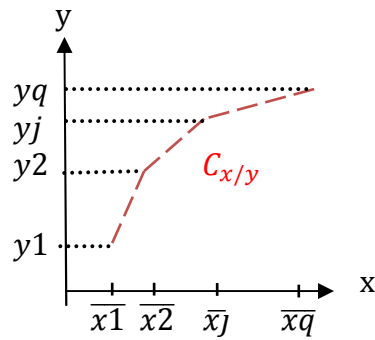
→ courbe de régression

ou encore lignes polygonales



Ici, \bar{y}_i synthétise l'information des (yj, nij)

ET



Ici, \bar{x}_j synthétise l'information des (x_i, n_{ij})

Que représente ces courbes de régression?

On peut dire que la Σ du carré des distances des points à $C_{y/x}$ ou $C_{x/y}$ est minimale!

Si y est en fonction de x, \bar{x} , on obtient alors un nuage de points.

On va chercher la meilleure courbe passant au milieu de ce nuage de points.

Cette courbe sera notée $y'_i = f(x_i)$

Donc, on doit avoir la somme S des distances au carré aux points y_j minimale (distance y'_i, y_j minimale).

$$S = \sum_i \sum_j n_{ij} (y_j - y'_i)^2$$

On veut que S soit le plus petit possible!

D'où : ... développement (la flemme) ...

S est minimale si :

$$y'_i = \bar{y}_i$$

⇒ courbe de régression.

4. Série à caractères x et y indépendants :

Définition : Les caractères sont indépendants si une variation de l'un n'entraîne pas une variation de l'autre.

Conséquences : Les fréquences conditionnelles $f_{i/j}$ ne dépendent pas j ($f_{i/j} = \frac{n_{ij}}{n.j}$)

Les fréquences conditionnelles $f_{j/i}$ ne dépendent pas i ($f_{j/i} = \frac{n_{ij}}{n.i}$)

$$f_{i/j} = f_{i.} \quad \text{et} \quad f_{j/i} = f_{.j}$$

→ fréquence marginale

Démonstration :

$$f_{i/j} = \frac{n_{ij}}{n.j} = \frac{n_{i1}}{n.1} = \frac{n_{i2}}{n.2} = \frac{n_{i3}}{n.3} = \dots = \frac{n_{iq}}{n.q} = \frac{n_{i.}}{n} = f_{i.}$$

On a vu que

$$f_{i/j} = \frac{f_{ij}}{f_{.j}}$$

donc

$$\boxed{f_{ij} = f_{i.} * f_{.j}}$$

Si indépendant.

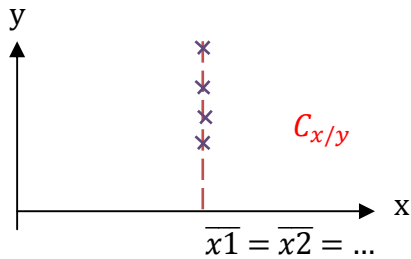
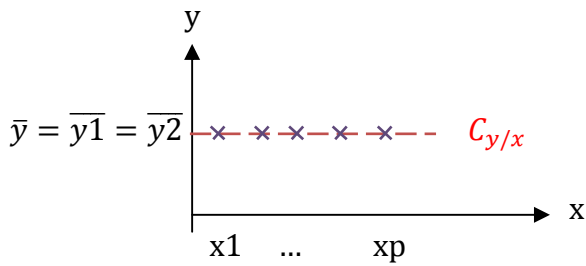
Connaissance de $f_{ij} \Rightarrow$ connaissance de $f_{i.}$ et $f_{.j}$

⇐ réciproque juste si x et y indépendants

On a :

$$\bar{x} = \bar{x}_j \quad \text{et} \quad \bar{y} = \bar{y}_i$$

Conséquence pour les courbes de régression :



Les courbes de régression sont des droites //

On a $\overline{xy} = \bar{x} \bar{y}$

d'où $cov(x,y) = 0$

Remarque : Si $cov(x,y) = 0$, cela n'entraîne pas forcément l'indépendance de x et de y

Exemple : Soit la relation $y=x^2$ (pas dépendante)

où $x = -1$, $x = 1$, $y = +1$

x/y	+1	Σ
-1	1	1
+1	1	1
Σ	+2	2

$$\begin{aligned} \bar{x} &= -1 * \frac{1}{2} + 1 * \frac{1}{2} = 0 \\ \overline{xy} &= \frac{1}{n} \sum_i \sum_j n_{ij} x_i y_i \\ &= \frac{1}{2} (-1 * (+1) + 1 * 1) = 0 \end{aligned}$$

donc $\overline{xy} - \bar{x} \bar{y} = 0$

$\Rightarrow cov(x,y) = 0$

5. Ajustement, corrélation et régression

Objectifs :

- évaluer quantitativement le degré d'indépendance entre x et y
- savoir si il existe une dépendance linéaire entre x et y ($y = ax + b$)

a. Ajustement linéaire – Droites de régression

Problème :

• si visuellement le "nuage" de points fait apparaître une droite, les courbes de régression ("meilleures" courbes représentant le nuage) sont les lignes brisées autour de droites que l'on va déterminer avec la méthode des moindres carrés (MMC).

On va parler de droites de régression

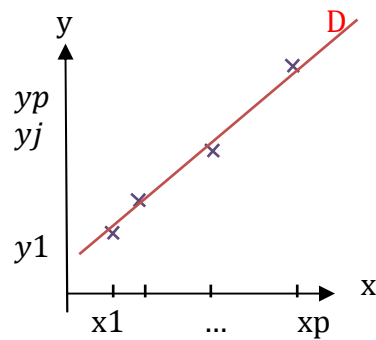
• On pourra alors déterminer l'angle entre ces deux droites de régression qui va quantifier le degré de dépendance entre x et y :

Si angle de 90° : x et y indépendantes

Si angle de 0° : on verra que y est fonction de x

d'où les graphes dans les deux cas:

1^{er} cas :

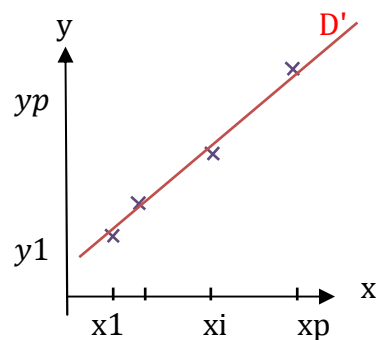


On va chercher $y = ax + b$

Cela va remplacer la courbe de régression $C_{x/y}$

D : droite de régression de y en x

2^{ème} cas :



On va chercher $x = a'y + b'$

D' : droite de régression de x en y

- MMC dans le cas 1 :

On veut minimiser la grandeur :

$$S = \sum_i \sum_j nij (yj - axi - b)^2$$

$$\frac{dS}{db} = 0 \Rightarrow$$

$$\sum_i \sum_j nij (yj - axi - b) = 0$$

$$\sum_i \sum_j nij yj = a \sum_i \sum_j nij xi + b \sum_i \sum_j nij$$

$$n\bar{y} = a n \bar{x} + nb$$

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

On sait que la droite de régression passe par (\bar{x}, \bar{y})

$$\frac{dS}{da} = 0 \Rightarrow$$

$$\sum_i \sum_j nij (yj - axi - b) xi = 0$$

$$n \bar{x}\bar{y} = a n \bar{x}^2 + nb\bar{x}$$

$$\bar{x}\bar{y} = a \bar{x}^2 + b\bar{x}$$

d'où

$$a = \frac{\bar{x}\bar{y}}{\bar{x}^2 - (\bar{x})^2} = \frac{cov(x, y)}{\sigma_x^2}$$

et

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

- MMC dans le cas 2 :

x en fonction de y

On obtient pour D' :

$$x = a'y + b'$$

On a alors

$$a' = \frac{\text{cov}(x, y)}{\sigma_y^2}$$

et

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