

Chapitre 3. Les distributions à deux variables

Jean-François Coeurjolly

<http://www-ljk.imag.fr/membres/Jean-Francois.Coeurjolly/>

Laboratoire Jean Kuntzmann (LJK), Grenoble University



Notes

- 1 Autour des tableaux de contingence
 - Définition
 - Distributions conditionnelles
 - Relations entre les différentes fréquences
 - Moyennes et Variances conditionnelles
- 2 Etude de la liaison entre deux variables
 - Mesure de la dépendance entre deux variables
 - Définition de l'indépendance totale
 - Définition de la dépendance totale
 - χ^2 et coefficient de Cramer
 - Mesure de la liaison fonctionnelle
 - Courbes de régression
 - Rapport de corrélation
 - Régression linéaire

Notes

Autour des tableaux de contingence

Etude de la liaison entre deux variables

Définition

Tableau de contingence

- = tableau statistique permettant de présenter séries statistiques et de
- exemple : dans une entreprise de 200 salariés, on étudie les variables $X=\text{âge}$ et $Y=\text{salaires}$.

$X=\text{Age} \setminus Y=\text{Salaire}$	$[800, 1000[$ $(j = 1)$	$[1000, 1200[$ $(j = 2)$	Total
$[20, 22[$ ($i = 1$)	14	6	20
$[22, 24[$ ($i = 2$)	28	46	74
$[24, 26[$ ($i = 3$)	20	86	106
Total	62	138	200

- X et Y sont des variables continues (regroupées en classes)
- On note I le nombre de modalités de X (ici) et J le nombre de modalités de Y (ici).

Notes

Autour des tableaux de contingence

Etude de la liaison entre deux variables

Définition

Tableau de contingence (2)

$X=\text{Age} \setminus Y=\text{Salaire}$	$[800, 1000[$ $(j = 1)$	$[1000, 1200[$ $(j = 2)$	Total
$[20, 22[$ ($i = 1$)	14	6	20
$[22, 24[$ ($i = 2$)	28	46	74
$[24, 26[$ ($i = 3$)	20	86	106
Total	62	138	200

- i désigne l'indice d'une et j désigne l'indice d'une .
- désigne l' .
Exemple : $n_{12} = 6$ salariés sont âgés entre 20 et 22 ans **et** ont un salaire compris entre 1000 et 1200 €.
- on note l' de X (eff. total en lignes) et l' de Y (effectif total en colonnes).
Exemple : $n_{2\bullet} = 74$ salariés sont âgés entre 22 et 24 ans ;
 $n_{\bullet 1} = 62$ salariés ont un salaire ente 800 et 1000€.
- correspond à l'effectif total.

Notes

Tableau de contingence (3)

X=Age \ Y=Salaire	[800, 1000[(j = 1)	[1000, 1200[(j = 2)	Total
[20, 22[(i = 1)	14	6	20
[22, 24[(i = 2)	28	46	74
[24, 26[(i = 3)	20	86	106
Total	62	138	200

Formules : Pour $i = 1, \dots, I$ et pour $j = 1, \dots, J$

$$n_{i\bullet} =$$

$$n_{i\bullet} =$$

$$n = n_{..} = \sum_{i=1}^I = \sum_{j=1}^J = \sum_{i=1}^I \sum_{j=1}^J .$$

Notes

[illegible]

Fréquences partielles et marginales

⚠ Les fréquences sont notées entre parenthèses.

X=Age \ Y=Salaire	[800, 1000[(j = 1)	[1000, 1200[(j = 2)	Total
[20, 22[(i = 1)	14 (10 %)	6 (5 %)	20 (15 %)
[22, 24[(i = 2)	28 (20 %)	46 (35 %)	74 (55 %)
[24, 26[(i = 3)	20 (15 %)	86 (65 %)	106 (80 %)
Total	62 (45 %)	138 (100 %)	200 (100%)

- n_i désigne la **fréquence**

Exemple : $f_{12} = 3\%$ des salariés sont âgés entre 20 et 22 ans **et** ont un salaire compris entre 1000 et 1200 €.

- on note la **fréquence** de X (fréq. totale en lignes) et la **fréquence** de Y (fréq. totale en colonnes).

Exemple : $f_{2\bullet} = 37\%$ des salariés sont âgés entre 22 et 24 ans ;
 $f_{\bullet 1} = 31\%$ des individus ont un salaire entre 800 et 1000€.

Notes

[illegible]

- Une distribution conditionnelle est une distribution statistique obtenue en la population à un (une classe par exemple).
- $J = 2 \Rightarrow$ il y a conditionnelles de X par rapport à Y .
 - 1 la distribution de X sachant $Y \in [800, 1000[$.
 - 2 la distribution de X sachant $Y \in [1000, 1200[$.
- $I = 3 \Rightarrow$ il y a distributions conditionnelles de Y par rapport à X
 - 1 la distribution de Y sachant $X \in [20, 22[$.
 - 2 la distribution de Y sachant $X \in [22, 24[$.
 - 3 la distribution de Y sachant $X \in [24, 26[$.

Notes

X=Age \ Y=Salaire	[800, 1000[(j = 1)	[1000, 1200[(j = 2)	Total
[20, 22[(i = 1)	14 (%)	6 (%)	20
[22, 24[(i = 2)	28 (%)	46 (%)	74
[24, 26[(i = 3)	20 (%)	86 (%)	106
Total	62 (100%)	138 (100%)	200

- On calcule les fréquences des âges en se restreignant à la sous-population des individus ayant un salaire entre **800 et 1000 €**, puis à la sous-population des individus ayant un salaire entre **1000 et 1200 €**.
- Les fréquences conditionnelles sont en général notées
- Interprétation :
 - **22.6%** des employés ayant un salaire entre 800 et 1000 €sont âgés entre 20 et 22 ans.
 - Parmi les employés ayant un salaire entre 1000 et 1200 €, **62.4%** d'entre eux sont âgés entre 24 et 26 ans.

Notes

Formules : Pour $i = 1, \dots, I$ et pour $j = 1, \dots, J$

$$f_{ij} = \frac{n_{ij}}{n_{\bullet j}}$$

(ex : 22.6% = $\frac{14}{62}$)

Notes

X=Age \ Y=Salaire	[800, 1000[(j = 1)	[1000, 1200[(j = 2)	Total
[20, 22[(i = 1)	14 (70%)	6 (30%)	20 100%
[22, 24[(i = 2)	28 (45.2%)	46 (74.2%)	74 100%
[24, 26[(i = 3)	20 (32.3%)	86 (70.7%)	106 100%
Total	62	138	200

- Ces fréquences conditionnelles sont en général notées
- Interprétation :
 - 70% des employés âgés entre 20 et 22 ans ont un salaire compris entre 800 et 1000 €.
 - Parmi les employés âgés entre 22 et 24 ans, 62.2% d'entre eux ont un salaire compris entre 1000 et 1200 €.

Notes

Formules : Pour $i = 1, \dots, I$ et pour $j = 1, \dots, J$

$$f_{ji} = \left(\text{ex : } \mathbf{30\%} = \frac{\mathbf{6}}{\mathbf{20}} \right)$$

En utilisant les précédentes définitions des fréquences conditionnelles, on peut obtenir

$$f_{ij} = f_{ij} \times f_{\bullet j}$$

De la même façon on peut obtenir

$$f_{ij} = f_{ji} \times f_{i\bullet}$$

Notes

X=Age \ Y=Salaire	[800, 1000[(j = 1)	[1000, 1200[(j = 2)	Total
[20, 22[(i = 1)	14	6	20
[22, 24[(i = 2)	28	46	74
[24, 26[(i = 3)	20	86	106
Total	62	138	200

Concentrons-nous sur la variable X : on notera \bar{x}_1 (ou $\bar{x}_{|Y \in [800, 1000[}$) et \bar{x}_2 (ou $\bar{x}_{|Y \in [1000, 1200[}$) les deux moy. cond. de X sachant Y :

- La moyenne de $X=$ la moyenne des moyennes conditionnelles

$$\bar{x} = \frac{1}{n} \sum_{j=1}^J n_{\bullet j} \bar{x}_j.$$

- Vérification :

- En utilisant la distribution marginale : $\bar{x} \simeq$ ans .
- En utilisant les fréq. conditionnelles, $\bar{x}_1 \simeq$ ans et $\bar{x}_2 \simeq$ ans .
- En combinant ans.

Notes

Décomposition de la variance

- Notons $\text{Var}_j(X)$ les variances conditionnelles de X sachant Y . Rappelons la formule de décomposition de la variance (qui peut s'exprimer en fonction des variances conditionnelles) :

$$Var(X) = \underbrace{\frac{1}{n} \sum_{j=1}^J n_{\bullet j} Var_j(X)}_{\text{within group}} + \underbrace{\frac{1}{n} \sum_{j=1}^J n_{\bullet j} (\bar{x}_j - \bar{x})^2}_{\text{between group}}$$

- La vérification sur l'exemple considéré est laissée en exercice.
- Des résultats tout à fait similaires sont bien évidemment valables pour la variable Y (notez que ceci est possible car Y est quantitative).

Notes

[illegible]

Généralités

Il y a deux extrêmes du niveau de liaison entre deux variables (quelles que soient la ou les natures des variables) :

- l' (ou liaison nulle).
- la (ou liaison fonctionnelle).

Le but de cette section est de mesurer la dépendance, et de quantifier en particulier le niveau de proximité par rapport aux deux cas précédents.

Notes

[illegible]

Définition

- 1 La variable Y est totalement indépendante de la variable X si les variations de X n'entraînent pas de variations de Y .
- 2 La variable X est totalement indépendante de la variable Y si les variations de Y n'entraînent pas de variations de X .

Théorème

- 1 Y est totalement indépendante de X si et seulement si

(c-a-d les fréquences conditionnelles ne dépendent pas des lignes du tableau de contingence et sont égales aux fréquences marginales).
- 2 X est totalement indépendante de Y si et seulement si
- 3 L'indépendance est .

Notes

[illegible]

Indépendance et tableau de contingence

Théorème

Les variables X et Y sont indépendantes si et seulement si

--

Corollaire

Un tableau de contingence est associé à deux variables X et Y indépendantes si et seulement si les \dots sont \dots entre elles.

Exemple : tableau associé à deux var. indépendantes

$X \mid Y$	y_1	y_2	y_3	Total
x_1	2	4	12	18
x_2	4	8	24	36
Total	6	12	36	54

On peut par exemple vérifier que

$$\frac{n_{2\bullet} \times n_{\bullet 3}}{n} = \frac{36 \times 36}{54} = 24 = n_{23}.$$

Notes

[illegible]

Dépendance totale

Définition

- 1 Y est fonctionnelle de X (ou fonctionnellement liée à X) si à chaque valeur x_i de X correspond une unique valeur y_j de Y , autrement dit si chaque ligne du tableau de contingence ne contient qu'un seul effectif n_{ij} non nul.
- 2 X est fonctionnelle de Y (ou fonctionnellement liée à Y) si à chaque valeur y_j de Y correspond une unique valeur x_i de X , autrement dit si chaque colonne du tableau de contingence ne contient qu'un seul effectif n_{ij} non nul.
- 3 ⚠ La dépendance totale n'est pas une relation .

Notes

Application à la notion de dépendance

Exemple 1 :

$X \backslash Y$	y_1	y_2
x_1	2	0
x_2	1	0
x_3	0	1

\Rightarrow est fonctionnelle de fonctionnelle
et la réciproque est fonctionnelle .

Exemple 2 :

$X \backslash Y$	y_1	y_2	y_3
x_1	2	0	0
x_2	0	1	4

\Rightarrow est fonctionnelle de fonctionnelle
et la réciproque est fonctionnelle .

Exemple 3 :

$X \backslash Y$	y_1	y_2
x_1	2	0
x_2	0	1

\Rightarrow est fonctionnelle de fonctionnelle
et la réciproque est fonctionnelle .

Notes

Autour des tableaux de contingence
oooooooooooooooo

Etude de la liaison entre deux variables
oooooooo●oooooooooooooooo

Mesure de la dépendance entre deux variables

χ^2 et Coefficient de Cramer

Définition

Le χ^2 est un nombre mesurant l'écart entre la situation observée et la situation si les variables avaient été théoriquement .

Méthodologie :

- construction du tableau de contingence sous hypothèse d'indépendance, c-a-d calcul des
- on calcule ensuite

$\chi^2 =$

Notes

Autour des tableaux de contingence
oooooooooooooooo

Etude de la liaison entre deux variables
oooooooo●oooooooooooooooo

Mesure de la dépendance entre deux variables

χ^2 et Coefficient de Cramer (2)

Théorème

La quantité χ^2_{\max} est la valeur du χ^2 si la dépendance entre X et Y était totale et réciproque.

Définition

Le coefficient de Cramer $C \in [0, 1]$ est défini par

- Si C est proche de alors les variables X et Y sont presque .
- Si C est proche de , alors les variables X et Y sont fortement (pas nécessairement liées fonctionnellement)
- Le C de Cramer peut être calculé pour n'importe quel type de variables X et Y .

Notes

Autour des tableaux de contingence
Mesure de la dépendance entre deux variables

Etude de la liaison entre deux variables

χ^2 et Coefficient de Cramer (3)

X=Age \ Y=Salaire	[800, 1000[(j = 1)	[1000, 1200[(j = 2)	Total
[20, 22[(i = 1)	14 ()	6 ()	20
[22, 24[(i = 2)	28 ()	46 ()	74
[24, 26[(i = 3)	20 ()	86 ()	106
Total	62	138	200

1 calcul des effectifs théoriques n'_{ij} .

Exemple : $n'_{32} = \frac{n_{3\bullet} \times n_{\bullet 2}}{n} = \frac{138 \times 106}{200} \approx 73.14$.

2 Calcul du χ^2

$$\chi^2 = \frac{(14 - \quad)^2}{\quad} + \frac{(6 - \quad)^2}{\quad} + \dots + \frac{(86 - \quad)^2}{\quad} \approx \quad.$$

3 $\chi^2_{\max} = 200 \times$

4 $C = \sqrt{\quad} \approx \quad \%$ (dépendance modérée).

Notes

Autour des tableaux de contingence
Mesure de la dépendance entre deux variables

Etude de la liaison entre deux variables

χ^2 et Coefficient de Cramer (4)

Question

Quels sont les couples (x_i, y_j) qui contribuent le plus au χ^2 ?

Réponse : il suffit de calculer pour chaque case le rapport

X=Age \ Y=Salaire	[800, 1000[(j = 1)	[1000, 1200[(j = 2)	Total
[20, 22[(i = 1)	14 (42.4%)	6 (19.1%)	20
[22, 24[(i = 2)	28 (4.8%)	46 (2.2%)	74
[24, 26[(i = 3)	20 (21.8%)	86 (9.8%)	106
Total	62	138	200

• Exemple 1ère case : $((6.2 - 14)^2/6.2)/23.13 \approx 42.4\%$.

• La case des individus

s'écarte le plus de l'hypothèse d'indépendance.

Notes

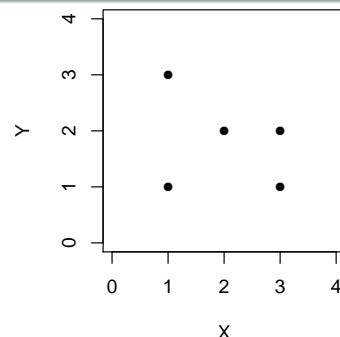
Généralités

- pour savoir si X et Y sont liées fonctionnellement, on trace le nuage de points (x_i, y_i) .
- \Rightarrow section valable uniquement pour X et Y
- \Rightarrow il faut disposer des données brutes, autrement dit chaque couple (x_i, y_i) est observée une et une seule fois. Autrement dit, la table de contingence correspondante ne contient que des
- On trace alors le nuage de points (x_i, y_j) et on essaie d'estimer la fonction de lien éventuelle.

Notes

Exemple et définition

$X \backslash Y$	1	2	3	Total
1	1	0	1	2
2	0	1	0	1
3	1	1	0	2
Total	2	2	1	5



Définition

- 1 est obtenue en faisant correspondre à chaque valeur de x_i de X la moy. conditionnelle de Y sachant $X = x_i$. Cette courbe est notée
- 2 est obtenue en faisant correspondre à chaque valeur de y_j de Y la moy. conditionnelle de X sachant $Y = y_j$. Cette courbe est notée

Notes

Propriétés

Théorème

- Si X et Y sont deux variables indépendantes alors $C_{Y/X}$ est parallèle à l'axe des abscisses et la courbe $C_{X/Y}$ est parallèle à l'axe des ordonnées (Δ réciproque fausse).
- Si aucun point ne s'écarte de , Y totalement dépendante de X ().
- Si aucun point ne s'écarte de , X totalement dépendante de Y ().

Notes

Concept basé sur la formule de décomposition de la variance

Définition

- 1 Le rapport de corrélation de Y en X est défini par

$$\eta_{Y/X}^2 = \frac{\frac{1}{n} \sum_i n_{i\bullet} (\bar{Y}_i - \bar{Y})^2}{\text{Var}(Y)}$$

- 2 Le rapport de corrélation de X en Y est défini par

$$\eta_{X/Y}^2 = \frac{\frac{1}{n} \sum_j n_{\bullet j} (\bar{X}_j - \bar{X})^2}{\text{Var}(X)}$$

- et
- Plus η^2 est (resp.) et plus la liaison fonctionnelle est (resp.)

Notes

X=Age \ Y=Salaire	[800, 1000[(j = 1)	[1000, 1200[(j = 2)	Total
[20, 22[(i = 1)	14	6	20
[22, 24[(i = 2)	28	46	74
[24, 26[(i = 3)	20	86	106
Total	62	138	200

Démarche pour calculer le rapport de corrélation de X en Y :

- calcul des moyenne et variance marginale de X : $\bar{x} \approx$ (ans) et $Var(X) \approx$ (ans²).
- calcul des moyennes conditionnelles de X sachant $Y \in [800, 1000[$ et de X sachant $Y \in [1000, 1200[$: $\bar{x}_1 \approx$ (ans) et $\bar{x}_2 \approx$ (ans).
- calcul de la variance interpopulation (var. moy. cond.)

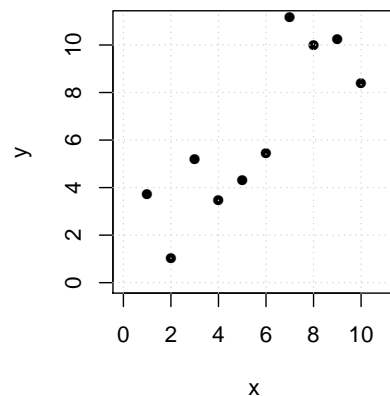
$$Var.Inter = \frac{62 \times (\quad - \quad)^2 + 138 \times (\quad - \quad)^2}{200} \approx \quad (ans^2).$$

- $\eta^2_{X/Y} \approx \frac{\quad}{\quad} \approx \quad \%$
($\quad \%$ de la variance de X est expliquée par la variable Y).

Notes

Régression linéaire

- Si le nuage de points observé est "presque" linéaire, il y a de fortes chances que la liaison entre X et Y soit linéaire (et que celle de Y à X soit linéaire).
- Exemple : imaginons observer le nuage suivant :



⇒ On peut suspecter une
Pour mesure ceci on
utilise le coefficient de

Notes

Soit (x_i, y_i) pour $i = 1, \dots, n$ un nuage de points. Ce coefficient est défini par

où

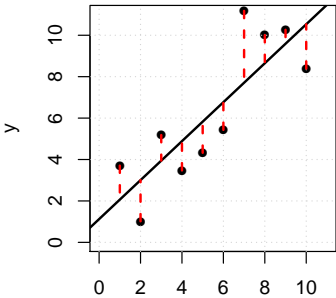
$$Cov(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \overline{xy} - \bar{x} \times \bar{y}.$$

- Si r est proche de 1 , X et Y sont (certainement) $Y = f(X)$
- Si r est proche de -1 , la pente de la droite est -1 . Si r est proche de 0 , la pente de la droite est 0 .
- Si r est proche de 0 , l'ajustement linéaire n'est pas adéquat (ce qui ne signifie pas que X et Y ne puissent pas être liées par une fonction).
- r^2 est appelé coefficient de détermination $(0 \leq r^2 \leq 1)$.

Notes

Si le coefficient r est jugé acceptable, on peut tenter d'estimer la droite de régression (de Y en X) en utilisant la

on se donne une droite d'équation $y = ax + b$, la MMC consiste à minimiser la somme des écarts rouges au carré.



Autrement dit, on va chercher le minimum en a et b de la fonction

Notes

Solutions au problème

La droite de régression ...

- ... de Y en X a pour équation $y = \widehat{a}x + \widehat{b}$ avec

$$\widehat{a} = \quad \text{et } \widehat{b} =$$

- ... de X en Y a pour équation $x = \widehat{a'}y + \widehat{b'}$ avec

$$\widehat{a'} = \quad \text{et } \widehat{b'} =$$

- les deux droites de régression passent par le point
- On peut remarquer que

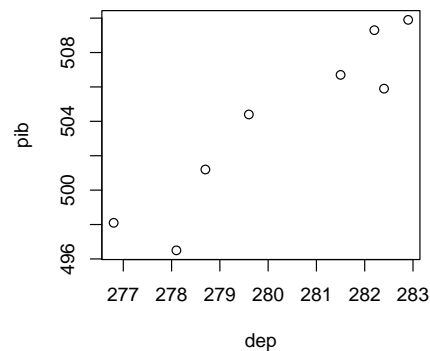
$$\widehat{a} \times \widehat{a'} =$$

Notes

Exemple d'application

Le tableau suivant présente les dépenses (dep) des ménages et PIB (pib) en milliards d'euros pour les 4 trimestres de 2011 et 2012. Peut-on expliquer l'évolution du PIB en fonction des dépenses ?

dep	278.1	276.8	278.7	279.6	282.4	281.5	282.2	282.9
pib	496.5	498.1	501.2	504.4	505.9	506.7	509.3	509.9



⇒ L'ajustement linéaire semble

Notes

Autour des tableaux de contingence
Mesure de la liaison fonctionnelle

Etude de la liaison entre deux variables

Exemple d'application (2)

dep	278.1	276.8	278.7	279.6	282.4	281.5	282.2	282.9
pib	496.5	498.1	501.2	504.4	505.9	506.7	509.3	509.9

Démarche

- 1 Calculez \overline{dep} , \overline{pib} , $Var(dep)$ et $Var(pib)$

$$\overline{dep} \approx (Me), \overline{pib} \approx (Me), Var(dep) \approx (Me)^2, Var(pib) \approx (Me)^2$$

- 2 Calcul intermédiaire

$$\overline{dep \times pib} = \frac{1}{8}(278 \times 496 + \dots + 283 \times 510) = (Me)^2.$$

- 3 Calcul de la covariance

$$Cov(dep, pib) = \overline{dep \times pib} - \overline{dep} \times \overline{pib} \approx (Me)^2.$$

- 4 Calcul du coefficient de corrélation linéaire

$$R = \sqrt{\quad} \approx$$

- 5 Puisque l'ajustement linéaire est très bon, calculons la droite de régression

$$\widehat{a} = \quad \text{et} \quad \widehat{b} = (Me).$$

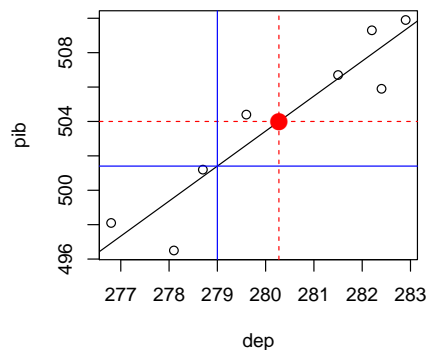
Notes

Autour des tableaux de contingence
Mesure de la liaison fonctionnelle

Etude de la liaison entre deux variables

Exemple d'application (3)

dep	278.1	276.8	278.7	279.6	282.4	281.5	282.2	282.9
pib	496.5	498.1	501.2	504.4	505.9	506.7	509.3	509.9



- La droite de régression

$$pib = 2.04 \times dep - 67.77.$$

 passe par le point $(\overline{dep}, \overline{pib})$.

- Quelle estimation du PIB proposer pour une $dep = 279$ (M€) ? \Rightarrow

$$\widehat{pib} = \quad = (Me).$$

Notes
