

Plan du cours

1. Introduction
2. Statistique descriptive - séries univariées
3. Calcul des probabilités
4. Variables aléatoires et lois de probabilité
5. Arbres de décision
6. Statistique descriptive - séries bivariées
7. Méthodes de prévision

2017/2018

222

222

Séries bivariées

- Objectif : étudier les liens entre deux variables x et y sur base d'une série d'observations.
- Plusieurs cas possibles :
 - 2 variables quantitatives,
 - 2 variables ordinales (ou quantitatives),
 - 2 variables nominales (ou ord. ou quant.).
- Généralisation à plus de 2 variables.

2017/2018

223

223

Cas 1 : variables quantitatives

- Tableau individus \times caractères

	x	y	← 2 caractères
1	x_1	y_1	
\vdots	\vdots	\vdots	
i	x_i	y_i	
\vdots	\vdots	\vdots	
n	x_n	y_n	

n individus →

2017/2018

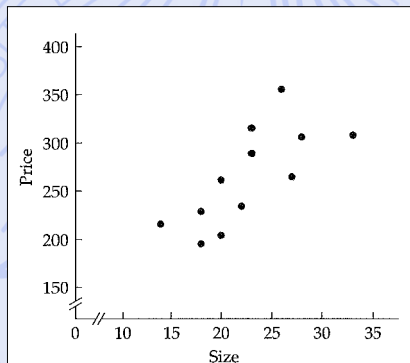
224

224

Exemple 1

- Lien entre le prix de vente d'une maison et sa surface habitable

Size	Price (\$000s)
23	315
18	229
26	355
20	261
22	234
14	216
33	308
28	306
23	289
20	204
27	265
18	195



2017/2018

225

225

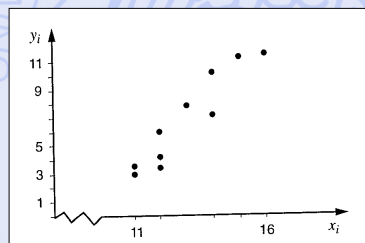
Exemple 2

- Série bivariée de 10 observations

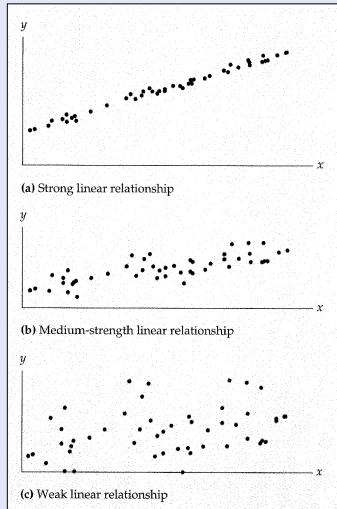
x_i	12	12	15	14	16	14	12	13	11	11
y_i	4.1	3.4	11.3	10.2	11.5	7.2	6.0	7.8	3.5	3.0

Représentation graphique

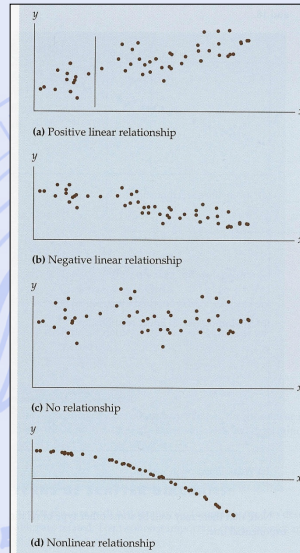
- Série bivariée :
 $\{(x_i, y_i); i = 1, 2, \dots, n\}$
- Représentation graphique - scatterplot :
 - Nuage de points dans le plan (x, y) .
 - Exemple 2 :



Représentation graphique



2017/2018



228

228

Séries marginales

- Série marginale en x :

$$\{x_i; i = 1, \dots, n\}$$

$$\Rightarrow \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad s_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \dots$$

- Série marginale en y :

$$\{y_i; i = 1, \dots, n\}$$

$$\Rightarrow \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \quad s_y^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad \dots$$

2017/2018

229

229

Distribution observée à 2 dimensions

- Lorsque les mêmes valeurs de x et de y se présentent plusieurs fois :

$$\{(x_j, y_k, n_{jk}); j = 1, \dots, J; k = 1, \dots, K\}$$

- Effectifs :

$$n_{jk} = \text{effectif associé à } (x_j, y_k)$$

- Propriété :

$$\sum_{j=1}^J \sum_{k=1}^K n_{jk} = n$$

2017/2018

230

230

Tableau de contingence

y	y_1	\dots	y_k	\dots	y_K
x					
x_1	n_{11}	\dots	n_{1k}	\dots	n_{1K}
\vdots	\vdots		\vdots		\vdots
x_j	n_{j1}	\dots	n_{jk}	\dots	n_{jK}
\vdots	\vdots		\vdots		\vdots
x_J	n_{J1}	\dots	n_{Jk}	\dots	n_{JK}

2017/2018

231

231

Exemple 3

- 80 femmes ayant eu au moins un enfant,
 - x : nombre d'enfants,
 - y : nombre de frères et sœurs.

x_j	y_k				
	0	1	2	3	4
1	4	4	2	0	0
2	9	16	4	0	0
3	4	12	9	2	0
4	1	6	1	1	2
5	0	1	0	1	1

2017/2018

232

232

Distributions marginales

- Distribution marginale en x :

$$\{(x_j, n_{.j}); j = 1, \dots, J\} \quad n_{.j} = \sum_{k=1}^K n_{jk}$$

- Distribution marginale en y :

$$\{(y_k, n_{.k}); k = 1, \dots, K\} \quad n_{.k} = \sum_{j=1}^J n_{jk}$$

- Propriété :

$$\sum_{j=1}^J n_{.j} = \sum_{k=1}^K n_{.k} = \sum_{j=1}^J \sum_{k=1}^K n_{jk} = n$$

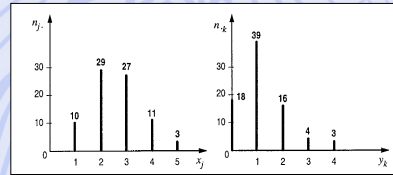
2017/2018

233

233

Exemple 3

	y_k					
x_j	0	1	2	3	4	$n_{j.}$
1	4	4	2	0	0	10
2	9	16	4	0	0	29
3	4	12	9	2	0	27
4	1	6	1	1	2	11
5	0	1	0	1	1	3
$n_{.k}$	18	39	16	4	3	



2017/2018

234

234

Distributions marginales

- Fréquences marginales :

$$f_{j.} = \frac{n_{j.}}{n} \quad f_{.k} = \frac{n_{.k}}{n}$$

- Paramètres :

$$\bar{x} = \frac{1}{n} \sum_{j=1}^J n_{j.} x_j \quad s_x^2 = \frac{1}{n} \sum_{j=1}^J n_{j.} (x_j - \bar{x})^2$$

$$\bar{y} = \frac{1}{n} \sum_{k=1}^K n_{.k} y_k \quad s_y^2 = \frac{1}{n} \sum_{k=1}^K n_{.k} (y_k - \bar{y})^2$$

2017/2018

235

235

Distributions conditionnelles

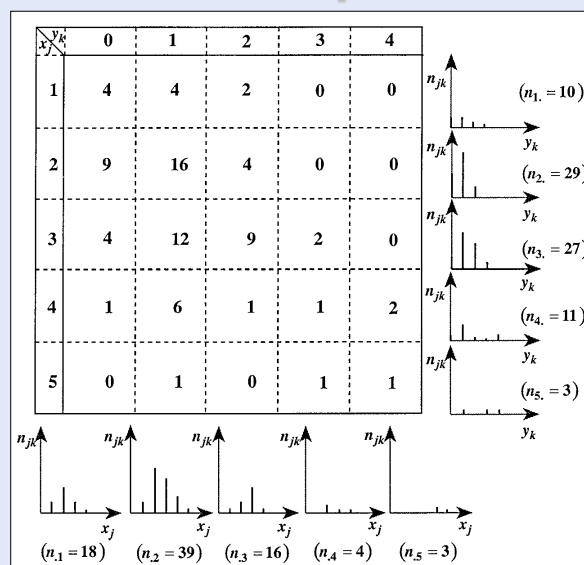
- Distribution d'une des variables lorsque la valeur de l'autre est fixée.
- Distribution conditionnelle de y en x :
 - On fixe la valeur de x : $x = x_j$
 $\{(y_k, n_{jk}); k = 1, \dots, K\}$
 - Effectif total : $n_{.j}$.
- Distribution conditionnelle de x en y :
 - On fixe la valeur de y : $y = y_k$
 $\{(x_j, n_{jk}); j = 1, \dots, J\}$
 - Effectif total : $n_{.k}$.

2017/2018

236

236

Exemple 3



2017/2018

237

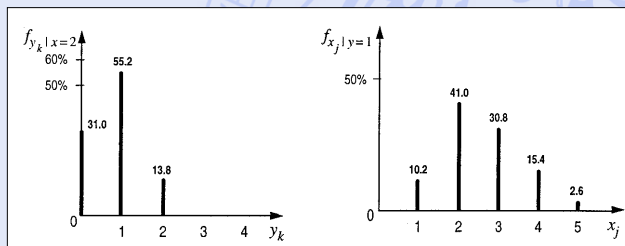
237

Distributions conditionnelles

- Fréquences conditionnelles :

$$f_{y_k|x_j} = f_{k|j} = \frac{n_{jk}}{n_{j.}} \quad f_{x_j|y_k} = f_{j|k} = \frac{n_{jk}}{n_{.k}}$$

- Exemple 3 :



2017/2018

238

238

Moyennes et variances conditionnelles

$$\bar{y}_{|x_j} = \frac{1}{n_{j.}} \sum_{k=1}^K n_{jk} y_k$$

$$s_{y|x_j}^2 = \frac{1}{n_{j.}} \sum_{k=1}^K n_{jk} (y_k - \bar{y}_{|x_j})^2$$

2017/2018

239

239

Moments

- Moments centrés :

$$m_{rs} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^r (y_i - \bar{y})^s \quad r, s \in \mathbb{N}$$

- Cas particuliers : $m_{20} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = s_x^2$
 $m_{02} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = s_y^2$

Covariance $\longrightarrow m_{11} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = s_{xy}$

Moments

- Moments par rapport à l'origine :

$$m'_{rs} = \frac{1}{n} \sum_{i=1}^n x_i^r y_i^s$$

- Cas particuliers :

$$m'_{10} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x} \quad m'_{01} = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$$

Propriétés des moments

- Changements d'origine et d'unité :

$$u_i = \frac{x_i - x_0}{d_x} \quad v_i = \frac{y_i - y_0}{d_y} \quad x_0, y_0 \in R; d_x, d_y \in R_0^+$$

$$\tilde{m}_{rs} = \frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})^r (v_i - \bar{v})^s$$

$$\Rightarrow \tilde{m}_{rs} = \frac{m_{rs}}{d_x^r d_y^s}$$

2017/2018

242

242

Propriétés des moments

- Cas particulier :

$$s_{uv} = \frac{s_{xy}}{d_x d_y} \Rightarrow s_{xy} = s_{uv} d_x d_y$$

- Calcul de la covariance (cf. variance) :

$$\begin{aligned} m_{11} = s_{xy} &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y} = m'_{11} - m'_{10} m'_{01} \end{aligned}$$

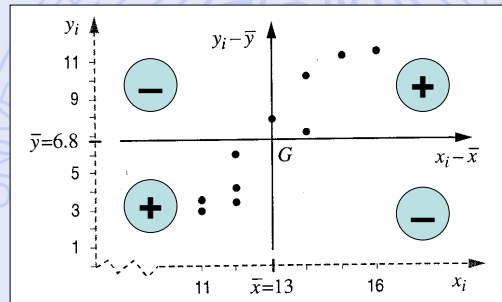
2017/2018

243

243

Covariance

- Permet de détecter une relation linéaire entre x et y :
- Positive si relation croissante,
- Négative si relation décroissante.



Covariance

- Propriétés :
 - Indépendante de changements d'origine.
 - Lien avec les variances marginales :

$$|s_{xy}| \leq s_x s_y$$

Coefficient de corrélation

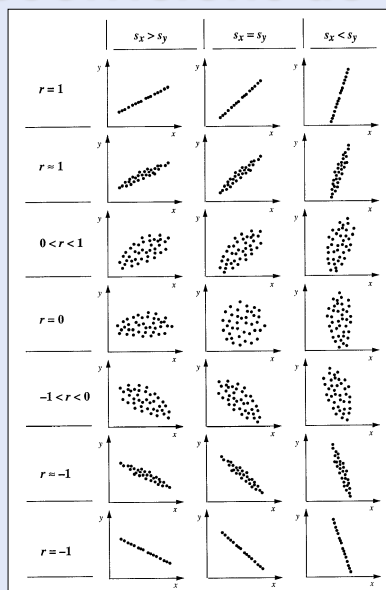
- Définition :

$$r = \frac{S_{xy}}{S_x S_y}$$

- Propriétés :

- Indépendant des changements d'origine ET d'unité.
- Même signe que la covariance.
- Compris entre -1 et +1.
- Mesure l'intensité de la relation linéaire entre x et y .

Coefficient de corrélation



Notation matricielle

- Vecteur moyenne (centre de gravité) :

$$\mathbf{g} = \begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix}$$

- Matrice de variance-covariance :

$$\mathbf{V} = \begin{pmatrix} s_x^2 & s_{xy} \\ s_{xy} & s_y^2 \end{pmatrix}$$

- Matrice symétrique

Notation matricielle

- Matrice définie par le tableau des données :

$$\mathbf{X} = \begin{pmatrix} x_1 & y_1 \\ x_2 & y_2 \\ \vdots & \vdots \\ x_n & y_n \end{pmatrix}$$

- Matrice des valeurs centrées :

$$\mathbf{X}_c = \begin{pmatrix} x_1 - \bar{x} & y_1 - \bar{y} \\ x_2 - \bar{x} & y_2 - \bar{y} \\ \vdots & \vdots \\ x_n - \bar{x} & y_n - \bar{y} \end{pmatrix}$$

Notation matricielle

- Propriété :

$$\mathbf{V} = \frac{1}{n} \mathbf{X}'_c \mathbf{X}_c$$

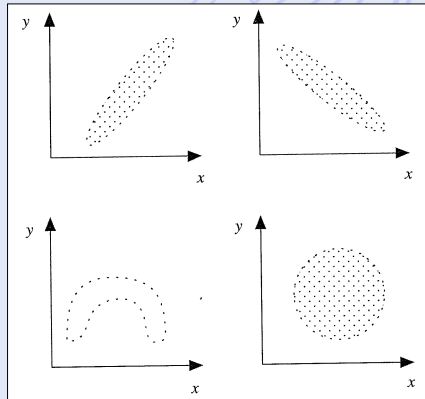
$$= \frac{1}{n} \begin{pmatrix} x_1 - \bar{x} & x_2 - \bar{x} & \cdots & x_n - \bar{x} \\ y_1 - \bar{y} & y_2 - \bar{y} & \cdots & y_n - \bar{y} \end{pmatrix} \begin{pmatrix} x_1 - \bar{x} & y_1 - \bar{y} \\ x_2 - \bar{x} & y_2 - \bar{y} \\ \vdots & \vdots \\ x_n - \bar{x} & y_n - \bar{y} \end{pmatrix}$$

Régression linéaire

- Objectif : Déterminer une relation de dépendance linéaire entre une variable y (variable dépendante) et une ou plusieurs variables explicatives.
 - Régression linéaire simple : une seule variable explicative (x),
 - Régression linéaire multiple : plusieurs variables explicatives (x_1, x_2, \dots, x_p).

Régression linéaire simple

- Attention : dépendance linéaire !
- Différents types de dépendance :



2017/2018

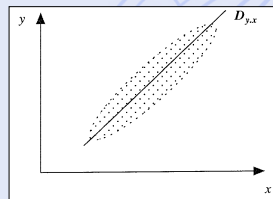
252

252

Droite de régression de y en x

- Equation :

$$D_{y.x} : y = a + bx \quad a, b \in \sim$$



- Problème : déterminer a et b à partir d'une série bivariée :

$$\{(x_i, y_i); i = 1, 2, \dots, n\}$$

2017/2018

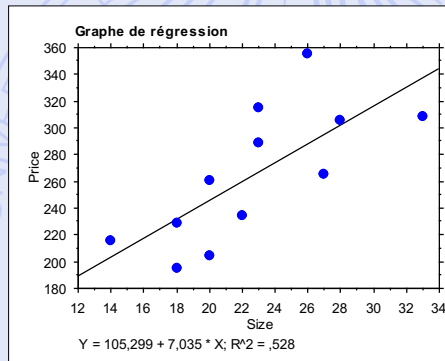
253

253

Exemple 1

- Lien entre le prix de vente d'une maison et sa surface habitable

Size	Price (\$000s)
23	315
18	229
26	355
20	261
22	234
14	216
33	308
28	306
23	289
20	204
27	265
18	195



2017/2018

254

254

Résidus

- Définition :

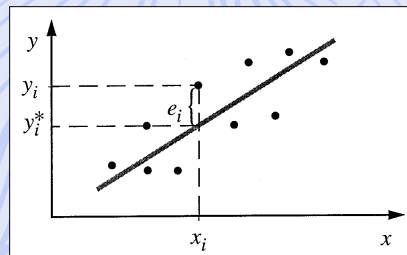
Valeurs ajustées

$$y_i^* = a + bx_i$$

⇓

$$e_i = y_i - y_i^* = y_i - a - bx_i$$

Résidus



2017/2018

255

255

Principe des moindres carrés

- Idée : choisir a et b de façon à rendre les résidus les plus petits possible.

- Minimiser $\sum_{i=1}^n |e_i|$?

Peu pratique d'un point de vue mathématique.

- Minimiser $\sum_{i=1}^n e_i^2$

→ Principe des moindres carrés !



Principe des moindres carrés

- Objectif :

$$\text{Min}_{a,b} Q(a,b) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a - bx_i)^2$$

- Solution :

$$\frac{\partial Q(a,b)}{\partial a} = 0 \quad \frac{\partial Q(a,b)}{\partial b} = 0$$

Calcul

- 1^{ère} dérivée partielle :

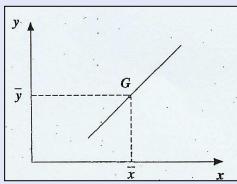
$$\frac{\partial Q(a,b)}{\partial a} = 0 \Leftrightarrow 2 \times (-1) \times \sum_{i=1}^n (y_i - a - bx_i) = 0$$

$$\Leftrightarrow \sum_{i=1}^n (y_i - a - bx_i) = 0$$

$$\Leftrightarrow \sum_{i=1}^n y_i - na - b \sum_{i=1}^n x_i = 0$$

$$\Leftrightarrow \sum_{i=1}^n y_i = na + b \sum_{i=1}^n x_i$$

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



2017/2018

258

258

Calcul

- 2^{ème} dérivée partielle :

$$\frac{\partial Q}{\partial b} = 0 \Leftrightarrow 2 \times (-1) \sum_{i=1}^n (y_i - a - bx_i) x_i = 0$$

$$\Leftrightarrow \sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i - b \sum_{i=1}^n x_i^2 = 0$$

$$\Leftrightarrow \frac{1}{n} \sum_{i=1}^n x_i y_i - \underbrace{(\bar{y} - b\bar{x})}_{a} \bar{x} - b \frac{1}{n} \sum_{i=1}^n x_i^2 = 0$$

$$\Leftrightarrow \underbrace{\left(\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y} \right)}_{s_{xy}} - b \underbrace{\left(\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \right)}_{s_x^2} = 0$$

2017/2018

259

259

Solution

$$a = \bar{y} - b\bar{x} \quad b = \frac{S_{xy}}{S_x^2}$$

- Remarques :

- Il s'agit bien d'un minimum (Cf. dérivées secondes).
- Autres expressions de la droite de régression :

$$y = \bar{y} + b(x - \bar{x})$$

$$y - \bar{y} = \frac{S_{xy}}{S_x^2}(x - \bar{x})$$

2017/2018

260

260

2 droites de régression !

- Régression de y en x :

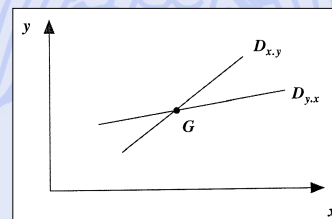
- Coefficient de régression de y en x : $b = \frac{S_{xy}}{S_x^2} = b_{y..x}$

- Régression de x en y :

- Coefficient de régression de x en y : $b_{x..y} = \frac{S_{xy}}{S_y^2}$

$$D_{x..y} : x = \bar{x} + \frac{S_{xy}}{S_y^2}(y - \bar{y})$$

$$\Leftrightarrow y = \bar{y} + \frac{S_y^2}{S_{xy}}(x - \bar{x})$$



2017/2018

261

261

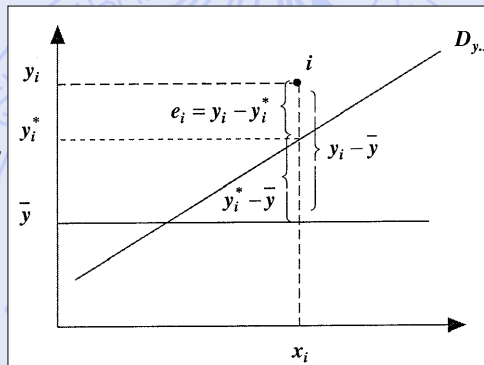
Variances résiduelle et de régression

- Décomposition de la variance de y :

$$s_y^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - y_i^*)^2}_{s_{y,x}^2} + \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i^* - \bar{y})^2}_{s_{reg}^2}$$

Variance résiduelle

Variance de régression



2017/2018

262

262

Démonstration

$$\begin{aligned} s_y^2 &= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^n \underbrace{(y_i - y_i^*)}_{\text{résiduel}} + \underbrace{(y_i^* - \bar{y})}_{\text{régression}}^2 \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - y_i^*)^2 + \frac{1}{n} \sum_{i=1}^n (y_i^* - \bar{y})^2 \\ &\quad + \underbrace{\frac{2}{n} \sum_{i=1}^n (y_i - y_i^*)(y_i^* - \bar{y})}_{=0 \text{ (à démontrer)}} \end{aligned}$$

Rappel : $y_i^* = \bar{y} + b_{y,x} (x_i - \bar{x})$

2017/2018

263

263

Corrélation et régression

- Lien entre r et le coefficient de régression :

$$b = \frac{s_{xy}}{s_x^2} = b_{y.x} \Rightarrow b_{y.x} = r \frac{s_y}{s_x}$$

$$r = \frac{s_{xy}}{s_x s_y}$$

2017/2018

264

264

Interprétation de r

$$s_{y.x}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - y_i^*)^2$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y} - b_{y.x} (x_i - \bar{x}))^2$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 + b_{y.x}^2 \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$- 2b_{y.x} \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$= s_y^2 + \frac{s_{xy}^2}{s_x^4} s_x^2 - 2 \frac{s_{xy}}{s_x^2} s_{xy} = s_y^2 - \frac{s_{xy}^2}{s_x^2 s_y^2} s_y^2$$

$$= s_y^2 \left(1 - \frac{s_{xy}^2}{s_x^2 s_y^2} \right) = s_y^2 (1 - r^2)$$

2017/2018

265

265

Interprétation de r

- Variance résiduelle :

$$s_{y.x}^2 = s_y^2 (1 - r^2)$$

- Variance de régression :

$$s_{reg}^2 = s_y^2 - s_{y.x}^2 = s_y^2 r^2$$

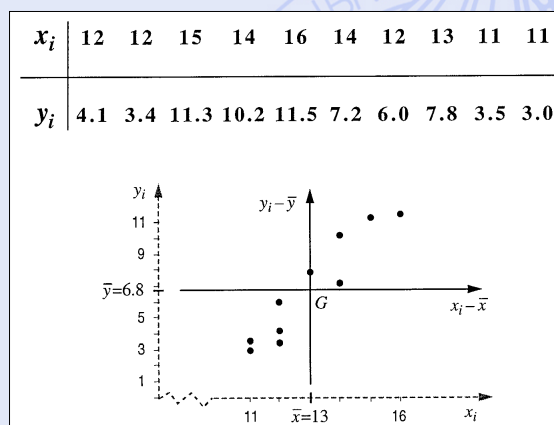
- Coefficient de détermination :

$$r^2 = \frac{s_{reg}^2}{s_y^2}$$

- Pourcentage de la variance de y expliqué par x .

Exemple 4

- Argent de poche hebdomadaire (y) en fonction de l'âge (x) - 14 observations (ados)



Exemple 4

x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
12	4.1	- 1	- 2.7	1	7.29	2.7
12	3.4	- 1	- 3.4	1	11.56	3.4
15	11.3	2	4.5	4	20.25	9.0
14	10.2	1	3.4	1	11.56	3.4
16	11.5	3	4.7	9	22.09	14.1
14	7.2	1	0.4	1	0.16	0.4
12	6.0	- 1	- 0.8	1	0.64	0.8
13	7.8	0	1.0	0	1.00	0
11	3.5	- 2	- 3.3	4	10.89	6.6
11	3.0	- 2	- 3.8	4	14.44	7.6
130	68.0	0	0	26	99.88	48.0

$$g = \begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix} = \begin{pmatrix} 13 \\ 6,8 \end{pmatrix}$$

$$V = \begin{pmatrix} s_x^2 & s_{xy} \\ s_{xy} & s_y^2 \end{pmatrix} = \begin{pmatrix} 2,6 & 4,8 \\ 4,8 & 9,99 \end{pmatrix}$$

$$r = 0,94$$

$$r^2 = 0,89$$

$$D_{y.x} : y = 6,8 + \frac{4,8}{2,6}(x - 13) = -17,2 + 1,85x$$

2017/2018

268

268

Exemples d'Ascombe

- Importance de l'examen graphique des données :
 - Forme de la relation entre x et y ,
 - Analyse des résidus (valeurs anormales).
- 4 jeux de données donnant lieu à la même droite de régression !

2017/2018

269

269

Exemples d'Ascombe

Ensemble A		Ensemble B		Ensemble C		Ensemble D	
x	y	x	y	x	y	x	y
10	8.04	10	9.14	10	7.46	8	6.58
8	6.95	8	8.14	8	6.77	8	5.76
13	7.58	13	8.74	13	12.74	8	7.71
9	8.81	9	8.77	9	7.11	8	8.84
11	8.33	11	9.26	11	7.81	8	8.47
14	9.96	14	8.10	14	8.84	8	7.14
6	7.24	6	6.13	6	6.08	8	5.25
4	4.26	4	3.10	4	5.39	19	12.50
12	10.84	12	9.13	12	8.15	8	5.56
7	4.82	7	7.26	7	6.42	8	7.91
5	5.68	5	4.74	5	5.73	8	6.89

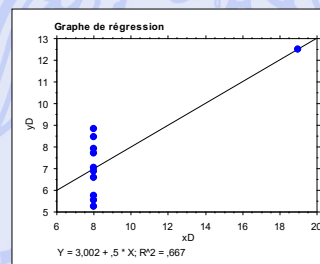
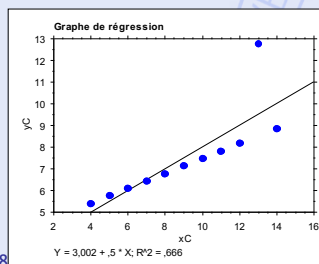
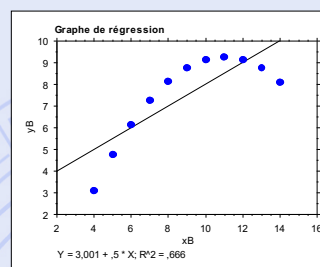
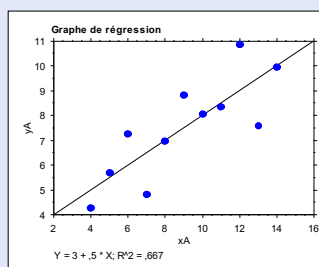
$n = 11; \bar{x} = 9; \bar{y} = 7.5; s_x^2 = 10; s_y^2 = 3.75; s_{xy} = 5$
 $r = 0.816; D_{y,x} \equiv y = 3 + 0.5x; r^2 = 0.667$

2017/2018

270

270

Exemples d'Ascombe



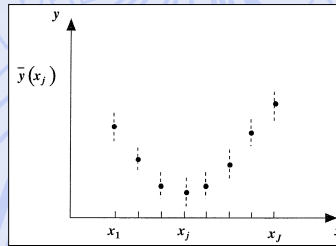
2017/2018

271

271

Régression non linéaire

- Dépendance non linéaire ? Forme ?
 - Analyse graphique :



- Ajustement d'une courbe non linéaire.
- Mesure de l'intensité de la dépendance non linéaire de y en x .

2017/2018

272

272

Ajustements non linéaires

- Ajustements linéarisables :
 - Régression exponentielle

$$y = 10^{a+bx} \Rightarrow \log_{10} y = a + bx$$

Régression linéaire de z en x , avec :

$$z = \log_{10} y \quad z_i = \log_{10} y_i$$

2017/2018

273

273

Ajustements non linéaires

- Ajustements linéarisables :
 - Régression hyperbolique

$$y = \frac{1}{a + bx} \Rightarrow \frac{1}{y} = a + bx$$

Régression linéaire de z en x , avec :

$$z = \frac{1}{y} \quad z_i = \frac{1}{y_i}$$

Ajustement polynomial

$$y = b_0 + b_1x + b_2x^2 + \dots + b_px^p$$

- Critères des moindres carrés :

$$\text{Min } Q(b_0, b_1, \dots, b_p) = \sum_{i=1}^n e_i^2$$

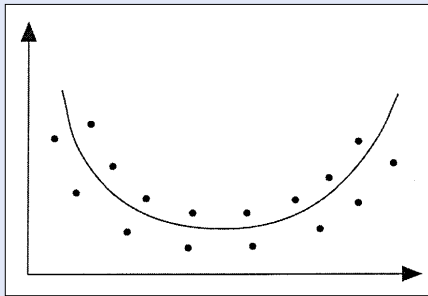
$$= \sum_{i=1}^n (y_i - b_0 - b_1x_i - b_2x_i^2 - \dots - b_px_i^p)^2$$

- Condition nécessaire pour un minimum :

$$\frac{\partial Q}{\partial b_j} = 0 \quad j = 0, 1, 2, \dots, p$$

Exemple : parabole

$$y = b_0 + b_1x + b_2x^2$$



$$\text{Min } Q(b_0, b_1, b_2)$$

$$= \sum_{i=1}^n (y_i - b_0 - b_1x_i - b_2x_i^2)^2$$

$$\rightarrow \begin{cases} \frac{\partial Q}{\partial b_0} = -2 \sum_{i=1}^n (y_i - b_0 - b_1x_i - b_2x_i^2) = 0 \\ \frac{\partial Q}{\partial b_1} = -2 \sum_{i=1}^n (y_i - b_0 - b_1x_i - b_2x_i^2)x_i = 0 \\ \frac{\partial Q}{\partial b_2} = -2 \sum_{i=1}^n (y_i - b_0 - b_1x_i - b_2x_i^2)x_i^2 = 0 \end{cases}$$

2017/2018

276

276

Rapport de corrélation

- Mesure de l'intensité de la dépendance non linéaire de y en x .
- Décomposition de la variance marginale :

$$\begin{aligned} s_y^2 &= \frac{1}{n} \sum_{k=1}^K n_k (y_k - \bar{y})^2 \\ &= \frac{1}{n} \sum_j \sum_k n_{jk} (y_k - \bar{y}(x_j) + \bar{y}(x_j) - \bar{y})^2 \\ &= \frac{1}{n} \sum_j \sum_k n_{jk} (y_k - \bar{y}(x_j))^2 + \frac{1}{n} \sum_j n_j (\bar{y}(x_j) - \bar{y})^2 \\ &= \text{moyenne des variances conditionnelles} \\ &\quad + \text{variance des moyennes conditionnelles} \end{aligned}$$

2017/2018

277

277

Rapport de corrélation

- Définition :

$$\theta_{y.x}^2 = \frac{\frac{1}{n} \sum_j n_j (\bar{y}(x_j) - \bar{y})^2}{s_y^2}$$

- A comparer avec :

$$r^2 = \frac{\text{variance de régression}}{\text{variance marginale de } y}$$

Propriétés

- Indépendant de l'origine et de l'unité.

- $0 \leq \theta_{y.x}^2 \leq 1$

$$s_y^2(x_j) = 0, j = 1, 2, \dots, J \Rightarrow \theta_{y.x}^2 = 1$$

$$\bar{y}(x_j) = \bar{y}, j = 1, 2, \dots, J \Rightarrow \theta_{y.x}^2 = 0$$

$$r^2 \leq \theta_{y.x}^2$$

- Indice de non linéarité de la régression :

$$\theta_{y.x}^2 - r^2$$

Cas 2 - Variables ordinales

- Série : 2 variables ordinales

$$\{(x_i, y_i); i = 1, 2, \dots, n\}$$

- Rangs :

$R(x_i)$ rang de x_i dans $\{x_{(1)}, x_{(2)}, \dots, x_{(n)}\}$

$R(y_i)$ rang de y_i dans $\{y_{(1)}, y_{(2)}, \dots, y_{(n)}\}$

Exemple

- Grades de 5 étudiants en 1^{ère} et 2^{ème} BA (1^{ère} session) :

1 ^{ère} BA	Sat.	GD	Aj.	LPGD	D
2 ^{ème} BA	Sat.	LPGD	D	GD	Aj.

- Rangs :

1 ^{ère} BA	2	4	1	5	3
2 ^{ème} BA	2	5	3	4	1

Coefficient de corrélation des rangs de Spearman

- Définition : r_s , coefficient de corrélation (Pearson) calculé sur les rangs.
- Exemple : $\{R(x_i), R(y_i)\} = \{(2,2), (4,5), (1,3), (5,4), (3,1)\}$
 $\Rightarrow r_s = 0,50$
- Propriété : si tous les x_i et y_i sont distincts :

$$r_s = 1 - \frac{6 \sum_i [R(x_i) - R(y_i)]^2}{n(n^2 - 1)}$$

2017/2018

282

282

Cas 3 - Variables nominales

- Tableau de contingence
 $\{(x_j, y_k, n_{jk}); j = 1, \dots, J; k = 1, \dots, K\}$
- Exemple : enquête sur 100 personnes
 – x = sexe, y = attitude vis-à-vis du sport

x_j	y_k			$n_{j.}$
	P	D	I	
F	21	15	9	45
H	39	13	3	55
$n_{.k}$	60	28	12	100

2017/2018

283

283

Profils lignes et profils marginaux

$$f_{k|j} = \frac{n_{jk}}{n_{j.}} \quad f_{.k} = \frac{n_{.k}}{n} \quad (j \text{ fixé}; k = 1, \dots, K)$$

$f_{k j}$	P	D	I	
F	0.47	0.33	0.20	1
H	0.71	0.24	0.05	1
$f_{.k}$	0.60	0.28	0.12	1

Profils colonnes et profils marginaux

$$f_{j|k} = \frac{n_{jk}}{n_{.k}} \quad f_{j.} = \frac{n_{j.}}{n} \quad (k \text{ fixé}; j = 1, \dots, J)$$

$f_{j k}$	P	D	I	$f_{j.}$
F	0.35	0.54	0.75	0.45
H	0.65	0.46	0.25	0.55
	1	1	1	1

Indépendance

- Si les profils lignes et les profils colonnes sont égaux aux profils marginaux.

$$f_{kj} = f_{.k} \Leftrightarrow \frac{n_{jk}}{n_{j.}} = \frac{n_{.k}}{n}$$

$$f_{jk} = f_{j.} \Leftrightarrow \frac{n_{jk}}{n_{.k}} = \frac{n_{j.}}{n}$$

- Condition d'indépendance :

$$n_{jk} = \frac{n_{j.} n_{.k}}{n}$$

2017/2018

286

286

Effectifs théoriques

- Définition :

$$n_{jk}^* = \frac{n_{j.} n_{.k}}{n}$$

- Exemple :

n_{jk}^*	P	D	I	
F	27	12.6	5.4	45
H	33	15.4	6.6	55
	60	28	12	100

2017/2018

287

287

Ecart à l'indépendance

- Différence entre effectifs observés et effectifs théoriques :

$$e_{jk} = n_{jk} - n_{jk}^*$$

- Exemple :

e_{jk}	P	D	I
F	- 6	2.4	3.6
H	6	- 2.4	- 3.6

2017/2018

288

288

Mesure du chi-carré (χ^2)

- Mesure de « distance » entre effectifs observés et effectifs théoriques :

$$D^2 = \sum_j \sum_k \frac{(n_{jk} - n_{jk}^*)^2}{n_{jk}^*} = \sum_j \sum_k \frac{e_{jk}^2}{n_{jk}^*}$$

- Exemple :

e_{jk}^2/n_{jk}^*	P	D	I	
F	1.33	0.46	2.40	4.19
H	1.09	0.37	1.96	3.42
	2.42	0.83	4.36	7.61

2017/2018

289

289