

Série statistique à deux variables

Niveau: 2^{ème} année

Pr : KEHAILI ABDELKADER

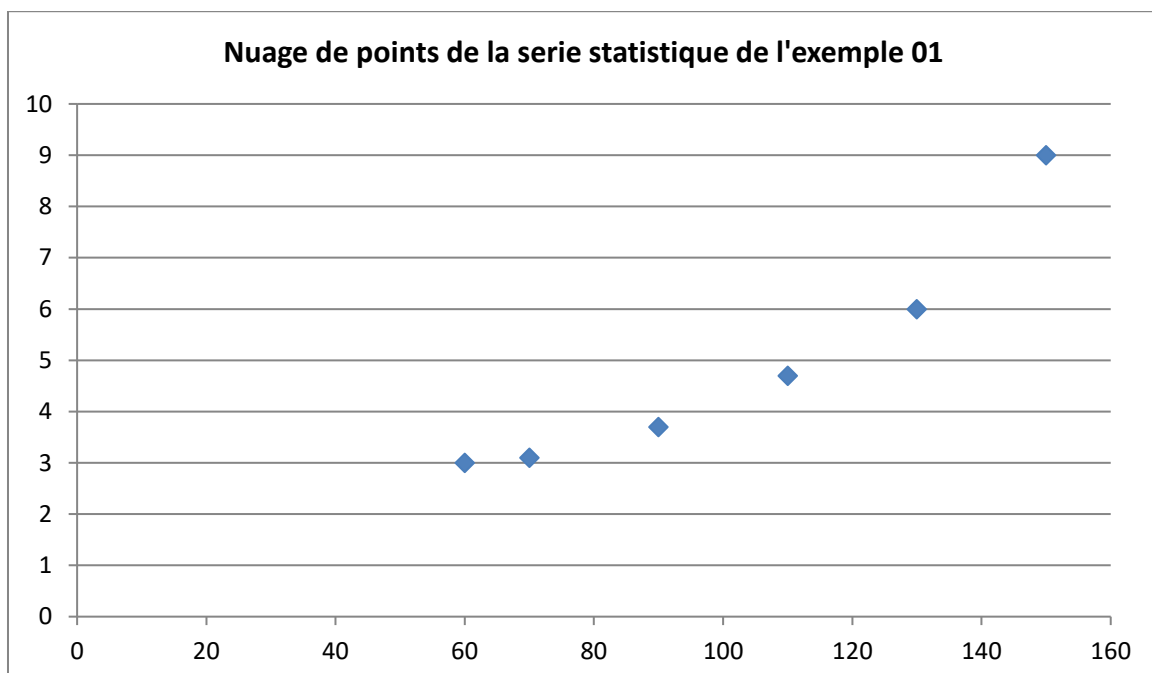
Définition : On appelle série statistique à deux variables (ou série statistique doubles) une série statistique à deux caractères sont étudiés simultanément.

Exemple 01:

On a relevé, pour un modèle de voiture, la consommation en carburant (en L/100 km) pour différentes vitesse (en km/h) sur le cinquième rapport :

Vitesse x_i (en km/h)	60	70	90	110	130	150
Consommation y_i (en L/100 km)	3	3.1	3.7	4.7	6	9

Définition : Dans un repère orthogonal, l'ensemble des points M_i de coordonnées (x_i, y_i) constitue le nuage de points associé à la série statistique à deux variables.



Point moyen

Définition : le Point moyen d'un nuage de points G de coordonnées (\bar{x}, \bar{y}) où :

\bar{x} représente la moyenne des x_i :

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{N} \sum_{i=1}^n x_i$$

\bar{y} représente la moyenne des y_i :

$$\bar{y} = \frac{y_1 + y_2 + \dots + y_n}{N} = \frac{1}{N} \sum_{i=1}^n y_i$$

Exemple : le Point moyen de l'exemple 01 est

$$\bar{x} = \frac{60 + 70 + 90 + 110 + 130 + 150}{6} = 101.66$$

$$\bar{y} = \frac{3 + 3.1 + 3.7 + 4.7 + 6 + 9}{6} = 4.91$$

Donc le point moyen $G(101.66, 4.91)$.

Ajustement affine par la méthode des moindres carrés

Définition : On appelle covariance de x et de y le nombre

$$\text{cov}(x, y) = \frac{1}{N} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \left(\frac{1}{N} \sum_{i=1}^n x_i y_i \right) - \bar{x} \bar{y}$$

Rappel : La variance de caractère x est :

$$V(x) = \frac{1}{N} \sum_{i=1}^n (x_i - \bar{x})^2 = \left(\frac{1}{N} \sum_{i=1}^n x_i^2 \right) - \bar{x}^2 = \text{cov}(x, x)$$

La variance de caractère y est :

$$V(y) = \frac{1}{N} \sum_{i=1}^n (y_i - \bar{y})^2 = \left(\frac{1}{N} \sum_{i=1}^n y_i^2 \right) - \bar{y}^2 = \text{cov}(y, y)$$

Elle est utilisée pour le calcul de l'écart type : $\sigma(x) = \sqrt{V(x)}$, $\sigma(y) = \sqrt{V(y)}$.

Exemple : Calculer dans l'exemple 01 $cov(x, y)$, $cov(x, x)$, $cov(y, y)$, $\sigma(x)$, $\sigma(y)$.

On a

							Somme
x_i	60	70	90	110	130	150	
y_i	3	3.1	3.7	4.7	6	9	
$x_i y_i$	180	217	333	517	780	1350	3377
x_i^2	3600	4900	8100	12100	16900	22500	68100
y_i^2	9	9.61	13.69	22.09	36	81	171.39

$$\bar{x} = 101.66, \quad \bar{y} = 4.91.$$

$$cov(x, y) = \left(\frac{1}{N} \sum_{i=1}^n x_i y_i \right) - \bar{x} \bar{y} = \frac{3377}{6} - 499.15 = 63.68.$$

$$V(x) = \left(\frac{1}{N} \sum_{i=1}^n x_i^2 \right) - \bar{x}^2 = \frac{68100}{6} - (101.66)^2 = 1015.2444.$$

$$V(y) = \left(\frac{1}{N} \sum_{i=1}^n y_i^2 \right) - \bar{y}^2 = \frac{171.39}{6} - (4.91)^2 = 4.4569.$$

$$\sigma(x) = \sqrt{V(x)} = \sqrt{1015.2444} = 31.86.$$

$$\sigma(y) = \sqrt{V(y)} = \sqrt{4.4569} = 2.11.$$

Théorème : Lors d'un ajustement affine par la méthode des moindres carrés.

1. La droite de régression D de Y en X a pour équation $D(Y/X) \ Y = aX + b$ où :

$$a = \frac{cov(x, y)}{V(x)}$$

Passe par le point moyen du nuage $G(\bar{x}, \bar{y})$ c'est-à-dire vérifié $\bar{Y} = a\bar{X} + b$, donc $b = \bar{Y} - a\bar{X}$.

2. La droite de régression D de X en Y a pour équation $D(X/Y) \ X = a'Y + b'$ où :

$$a' = \frac{cov(x, y)}{V(y)}$$

Passe par le point moyen du nuage $G(\bar{x}, \bar{y})$ c'est-à-dire vérifié $\bar{X} = a'\bar{Y} + b'$, donc $b' = \bar{X} - a'\bar{Y}$.

Exemple : Calculer dans l'exemple 01 La droite de régression D de Y en X

On a

$$\bar{x} = 101.66, \quad \bar{y} = 4.91, \quad cov(x, y) = 63.68, \quad V(x) = 1015.2444, \quad V(y) = 4.4569.$$

$$1. D(Y/X) \quad Y = aX + b$$

$$a = \frac{cov(x,y)}{V(x)} = 0.0627, \quad b = \bar{Y} - a\bar{X} = -1.46.$$

$$\text{Donc } D(Y/X) \quad Y = aX + b = 0.0627X - 1.46$$

$$2. D(X/Y) \quad X = a'Y + b'$$

$$a' = \frac{cov(x,y)}{V(y)} = 14.287, \quad b' = \bar{X} - a'\bar{Y} = 31.51$$

$$\text{Donc } D(X/Y) \quad X = a'Y + b' = 14.287Y + 31.51$$

Coefficient de corrélation linéaire

Définition : le Coefficient de corrélation linéaire d'une série statistique à deux variables x et y est le nombre r défini par :

$$r = \frac{cov(x,y)}{\sqrt{V(x)}\sqrt{V(y)}} = \frac{cov(x,y)}{\sigma(x)\sigma(y)}$$

Remarque :

1. $-1 \leq r \leq 1$.
2. Si $r = 1$ ou $r = -1$ alors il y a une corrélation positive ou négative parfaite entre X et Y et les points (x_i, y_i) sont tous sur la droite de régression.
 Une corrélation positive c'est-à-dire une augmentation de X entraîne une augmentation de Y .
 Une corrélation négative c'est-à-dire une augmentation de X entraîne une diminution de Y ou le contraire.
3. Si $r = 0$ alors il n'y a pas de corrélation entre X et Y et les points (x_i, y_i) sont dispersés au hasard.
4. Si $0 < r < 1$ alors il y a une corrélation positive faible, moyenne ou forte entre X et Y .
5. Si $-1 < r < 0$ alors il y a une corrélation négative faible, moyenne ou forte entre X et Y .

Exemple : Calculer dans l'exemple 01 le Coefficient de corrélation linéaire.

$$\text{On a } cov(x,y) = 63.68, \quad \sigma(x) = 31.86, \quad \sigma(y) = 2.11.$$

Donc

$$r = \frac{cov(x,y)}{\sigma(x)\sigma(y)} = \mathbf{0.947}$$

alors il y a une corrélation positive forte entre X et Y

Coefficient de détermination

Définition : le Coefficient de détermination d'une série statistique à deux variables x et y est le nombre R^2 défini par :

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{Y})^2}{\sum_{i=1}^n (y_i - \bar{Y})^2}$$

Tel que $\hat{y}_i = ax_i + b$

Exemple : Calculer dans l'exemple 01 le Coefficient de détermination.

							Somme
x_i	60	70	90	110	130	150	
y_i	3	3.1	3.7	4.7	6	9	
$(y_i - \bar{Y})^2$	3.65	3.27	1.46	0.044	1.188	16.72	26.332
$\hat{y}_i = ax_i + b$	2.3	2.92	4.18	5.43	6.69	7.94	
$(\hat{y}_i - \bar{Y})^2$	6.8	3.92	0.53	0.28	3.17	9.21	23.92

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{Y})^2}{\sum_{i=1}^n (y_i - \bar{Y})^2} = \frac{23.92}{26.332} = 0.9.$$

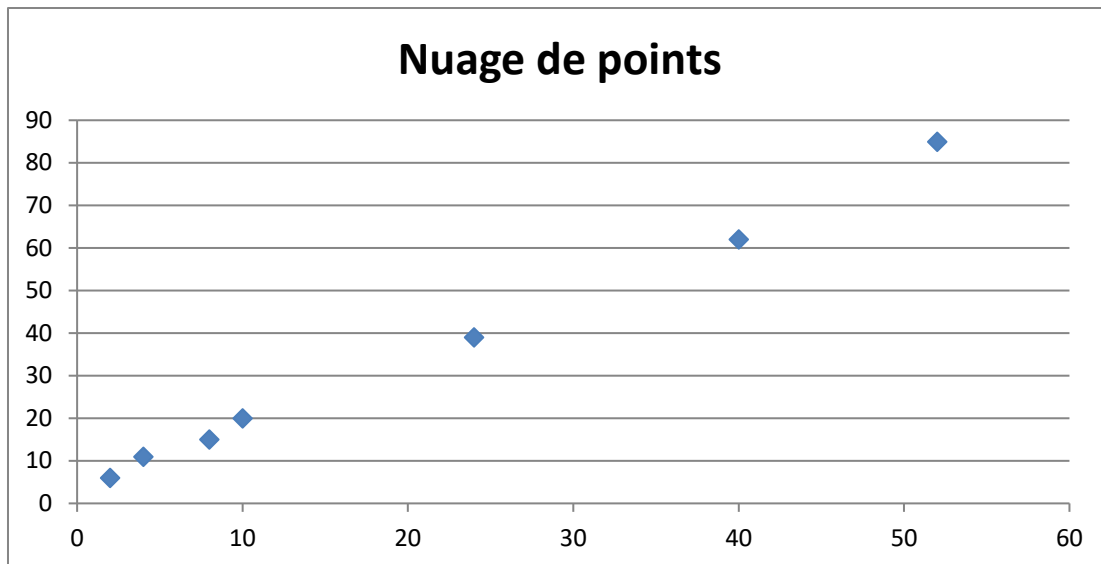
Exercice : Dans la série statistique suivante, X représente le nombre de jours d'exposition au soleil d'une feuille et Y le nombre de stomates aérifères au millimètre carré:

X	2	4	8	10	24	40	52
y	6	11	15	20	39	62	85

1. Tracer le nuage des points.
2. Calculer le coefficient de corrélation linéaire entre X et Y. Conclusion?
3. Déterminer l'équation de la droite de régression de Y en fonction de X.
4. Si on expose au soleil une feuille 15 jours; quel est le nombre de stomates aérifères peut-on prévoir ?

Corrigé type

1. Le nuage des points



2. Le coefficient de corrélation linéaire entre X et Y.

								Somme
x_i	2	4	8	10	24	40	52	140
y_i	6	11	15	20	39	62	85	238
$x_i y_i$	12	44	120	200	936	2480	4420	8212
x_i^2	4	16	64	100	576	1600	2704	5064
y_i^2	36	121	225	400	1521	3844	7225	13372

On a

$$\bar{x} = \frac{140}{7} = 20 \quad \bar{y} = \frac{238}{7} = 34.$$

$$\text{cov}(x, y) = \left(\frac{1}{N} \sum_{i=1}^n x_i y_i \right) - \bar{x} \bar{y} = \frac{8212}{7} - 680 = 493.14.$$

$$V(x) = \left(\frac{1}{N} \sum_{i=1}^n x_i^2 \right) - \bar{x}^2 = \frac{5064}{7} - (20)^2 = 323.43.$$

$$V(y) = \left(\frac{1}{N} \sum_{i=1}^n y_i^2 \right) - \bar{y}^2 = \frac{13372}{7} - (34)^2 = 754.28.$$

$$\sigma(x) = \sqrt{V(x)} = \sqrt{323.43} = 17.98.$$

$$\sigma(y) = \sqrt{V(y)} = \sqrt{754.28} = 27.46.$$

Donc le coefficient de corrélation linéaire entre X et Y.

$$r = \frac{\text{cov}(x, y)}{\sigma(x)\sigma(y)} = 0.99$$

alors il y a une corrélation positive forte entre X et Y

3. L'équation de la droite de régression de Y en fonction de X.

$$D(Y/X) \quad Y = aX + b$$

$$a = \frac{\text{cov}(x,y)}{V(x)} = 1.52, \quad b = \bar{Y} - a\bar{X} = 3.6.$$

$$\text{Alors } D(Y/X) \quad Y = 1.52 X + 3.6$$

4. le nombre de stomates aérifères

$$\text{On a } Y = 1.52 X + 3.6 = 1.52 \times 15 + 3.6 = 26.4$$

								Somme
x_i	2	4	8	10	24	40	52	140
y_i	6	11	15	20	39	62	85	238
$(y_i - \bar{Y})^2$	784	529	361	196	25	784	2601	
$\hat{y}_i = ax_i + b$	6.64	9.68	15.76	18.8	40.08	64.4	82.64	
$(\hat{y}_i - \bar{Y})^2$	748.57	591.46	332.69	231.04	36.96	924.16	2365.84	

$$Y = 1.52 X + 3.6$$