

Régression linéaire - Applications aux séries chronologiques

Présentation de la droite de régression linéaire par les moindres carrés sur un exemple

Loyer mensuel de locaux commerciaux (cf. en annexe jeu de données 1)

a) Les données

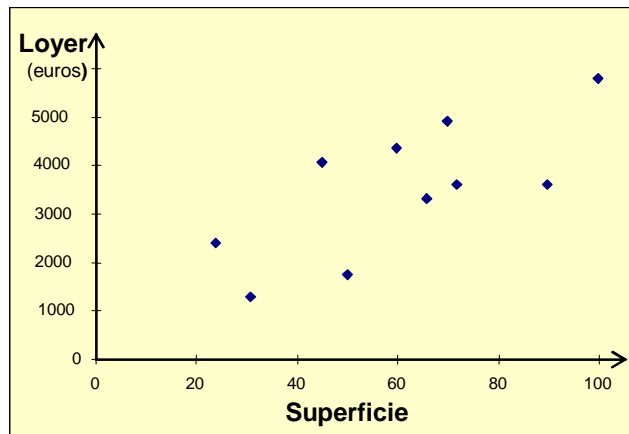
Population : 10 locaux commerciaux d'un quartier

Le loyer mensuel en euros : Y

La superficie en mètres carrés : X

Observations : $(X_i, Y_i) \quad i = 1, \dots, n$

Loyer Y (en Euros)	Superficie X (en m^2)
2 395	24
1 265	31
4 050	45
1 730	50
4 350	60
3 300	66
3 600	72
4 900	70
3 600	90
5 760	100



b) Les résumés numériques

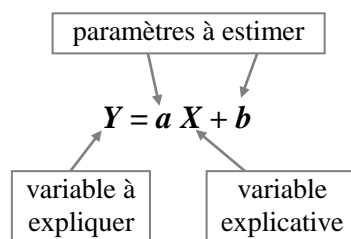
On obtient ici :

$$\bar{X} = 60.8 \text{ m}^2 ; \bar{Y} = 3495 \text{ €} ; \sigma(X) = 22.9 \text{ m}^2 \left(= \sqrt{V(X)} \right) ; \sigma(Y) = 1322.9 \text{ €} \left(= \sqrt{V(Y)} \right)$$

$$\text{COV}(X, Y) = 22148.5 \text{ m}^2 \times \text{€} ; \rho(X, Y) = 0.73$$

On a une assez bonne corrélation linéaire positive. Il est possible d'utiliser un modèle linéaire pour expliquer le loyer mensuel par la superficie du local.

c) Le modèle algébrique



Estimation des paramètres par la méthode des moindres carrés des erreurs

Introduction des erreurs E_i

$$Y_i = aX_i + b + E_i \quad i = 1, \dots, n$$

Les valeurs de a et de b qui rendent minimum la moyenne des carrés des erreurs :

$$f(a, b) = \frac{1}{n} \sum_{i=1}^n (Y_i - aX_i - b)^2$$

sont :

$$\hat{a} = \frac{\text{COV}(X, Y)}{V(X)} \quad \text{et} \quad \hat{b} = \bar{Y} - \hat{a}\bar{X}$$

La droite d'équation $y = \hat{a}x + \hat{b}$ est appelée droite de régression linéaire de Y en X .

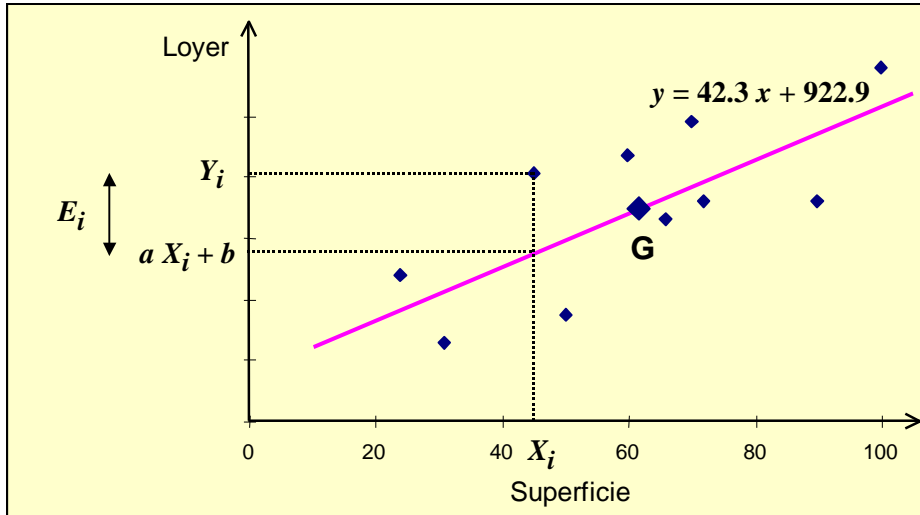
Ajustement du modèle et interprétation

La droite de régression linéaire de Y en X a pour équation :

$$y - 3495 = 42.3(x - 60.8) \text{ soit } y = 42.3x + 922.9$$

Le loyer mensuel moyen d'un local est de 3495 euros pour une superficie moyenne de 60.8 m² ;

- pour chaque m² en plus, le loyer augmente en moyenne de 42.3 euros
- pour chaque m² en moins, le loyer diminue en moyenne de 42.3 euros

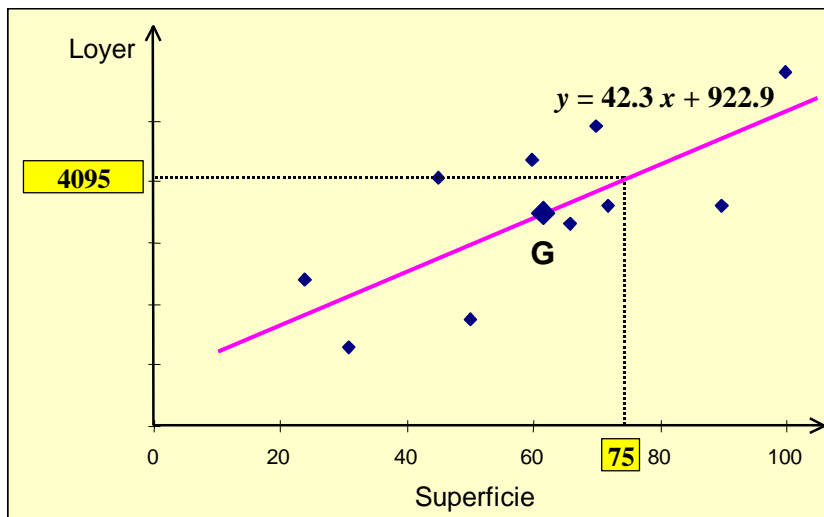


d) Prévision

Quel loyer prévoir pour un local de 75 m² ?

Le modèle ajusté est : $\hat{Y} = 42.3X + 922.9$

Pour $X = 75$ on a : $\hat{Y} = 42.3 \times 75 + 922.9 \cong 4095$ euros



e) Décomposition de la variance

On construit les deux variables :

$$\hat{Y} = \hat{a}X + \hat{b} \quad (\text{le modèle}) \text{ et}$$

$$\hat{E} = Y - \hat{Y} \quad (\text{l'erreur d'ajustement}).$$

	X (Surface)	Y (Prix)	$\hat{Y} = \hat{a}X + \hat{b}$	$\hat{E} = Y - \hat{Y}$
	24	2 395	1 938.2	456.8
	31	1 265	2 234.3	-969.3
	45	4 050	2 826.6	1 223.4
	50	1 730	3 038.1	-1 308.1
	60	4 350	3 461.2	888.8
	66	3 300	3 715.0	-415.0
	72	3 600	3 968.8	-368.8
	70	4 900	3 884.2	1 015.8
	90	3 600	4 730.3	-1 130.3
	100	5 760	5 153.3	606.7
Moyenne	61	3 495	3 495	0
Variance	523.56	1 750 150	936 962.435	813 187 .565

On vérifie $\overline{\hat{Y}} = \bar{Y}$, $\overline{\hat{E}} = 0$, $\text{COV}(\hat{Y}, \hat{E}) = 0$

et

$$\boxed{V(Y) = V(\hat{Y}) + V(\hat{E})}$$

$R^2 = \frac{V(\hat{Y})}{V(Y)}$ (= 0.535) est la part de variance expliquée par le modèle ;

$1 - R^2$ (= 0.465) est la part résiduelle.

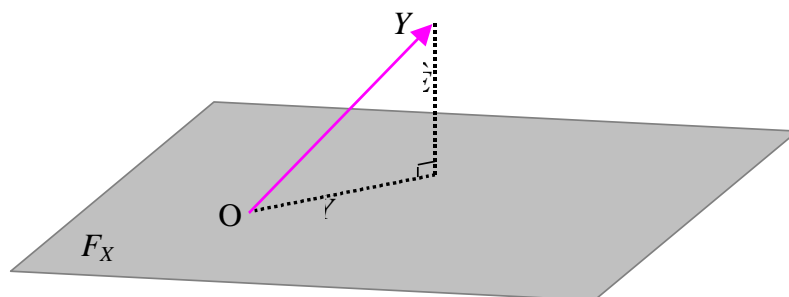
Interprétation géométrique

Dans \mathbf{R}^n muni du produit-scalaire : $\langle x, y \rangle_n = \frac{1}{n} \sum_{i=1}^n x_i y_i$

$$f(a, b) = \frac{1}{n} \sum_{i=1}^n (Y_i - aX_i - b)^2 = \|Y - aX - b\mathbf{1}_n\|_n^2$$

avec $Y = (Y_1, \dots, Y_n)$, $X = (X_1, \dots, X_n)$ et $\mathbf{1}_n = (1, \dots, 1)$.

Soit $F_X = \{aX + b\mathbf{1}_n \mid (a, b) \in \mathbf{R}^2\}$ le plan engendré par X et $\mathbf{1}_n$.



La recherche de a et b minimisant $f(a, b)$ est la recherche du vecteur de F_X le plus proche de Y .

La solution $\hat{Y} = \hat{a}X + \hat{b}\mathbf{1}_n$ est la projection orthogonale de Y sur F_X .

Généralisation

- Plusieurs variables explicatives (X_1, \dots, X_p) → Modèle de régression linéaire multiple
- Une variable explicative catégorielle à p modalités (X_1, \dots, X_p) les indicatrices des modalités → Analyse de variance à un facteur
- Plusieurs variables explicatives réelles ou catégorielles → Analyse de covariance

Coefficient de corrélation linéaire

$$\rho(X, Y) = \frac{\text{COV}(X, Y)}{\sigma(X)\sigma(Y)} \quad -1 \leq \rho \leq 1$$

$$|\rho(X, Y)| = 1 \Leftrightarrow \exists a \in \mathbf{R}^+ \text{ et } b \in \mathbf{R} \quad Y = aX + b$$

On a ici :
$$\rho(X, Y) = 0.73 \quad (= \sqrt{R^2})$$

Origine historique du mot "régression"

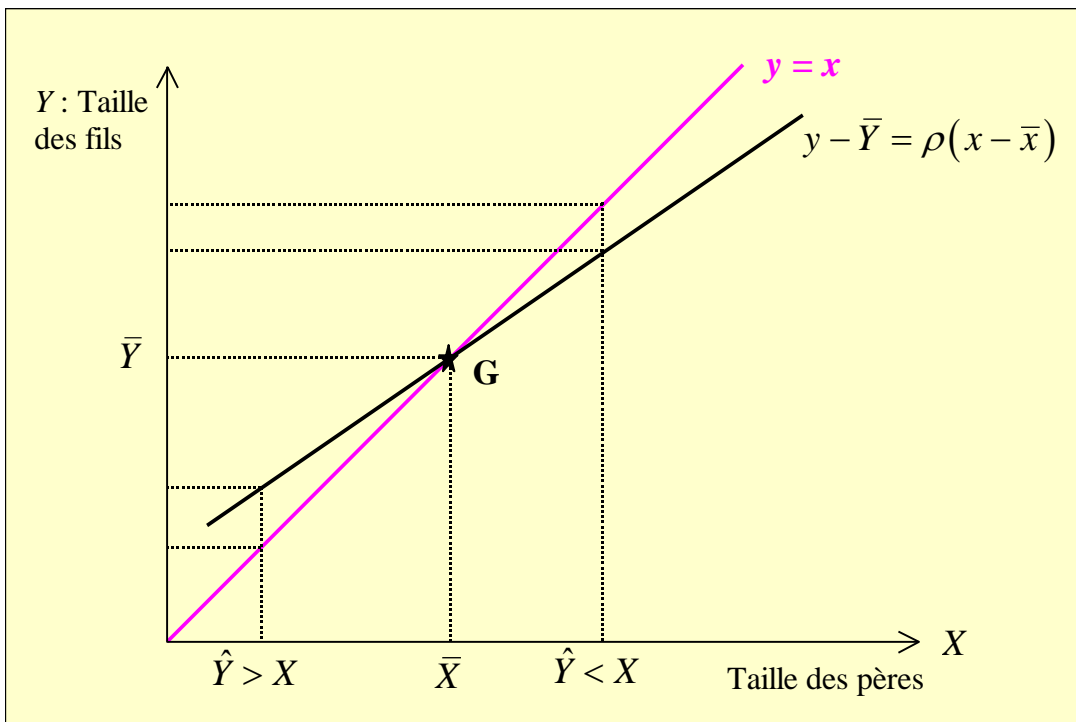
$$Y - \bar{Y} = \frac{\text{COV}(X, Y)}{V(X)}(X - \bar{X}) + E$$

Galton (1822, 1911) étudie la taille des fils (Y) en fonction de la taille des pères (X).

Bien que les moyennes (resp. les variances) des tailles des fils et des pères soient égales, il observe une *régression vers la moyenne*.

En effet, l'équation de la droite de régression linéaire de Y en X s'écrit alors :

$$y - \bar{Y} = \rho(x - \bar{X}) \text{ et } \rho < 1$$



Applications aux séries chronologiques

Chiffre d'affaires trimestriel d'une entreprise (cf. annexe jeu de données 2)

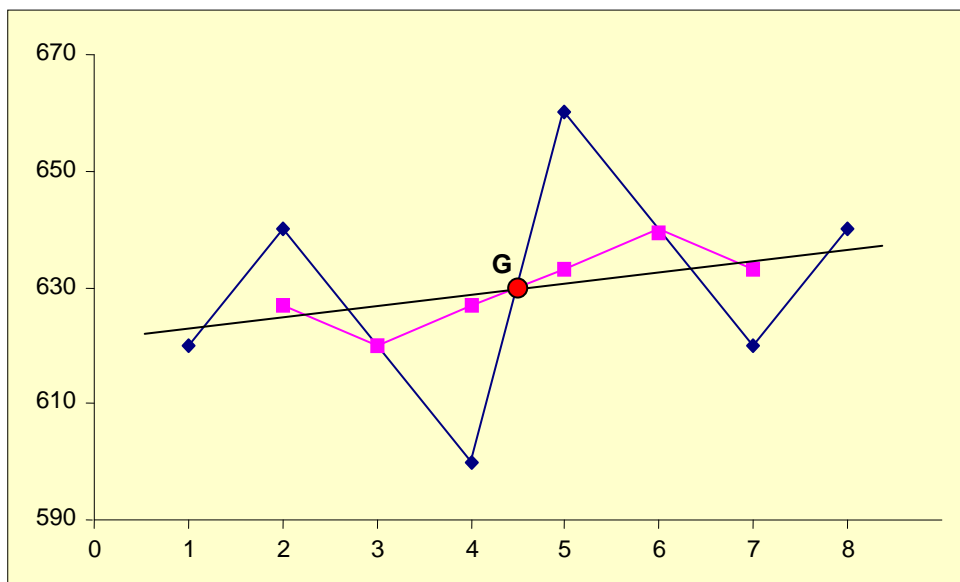
a) et b)

Le chiffre d'affaire (CA) trimestriel en KE d'une entreprise est observé sur 8 trimestres consécutifs. On a a priori une seule variable réelle, le CA, observée sur n individus ordonnés (les trimestres) mais pour étudier l'évolution de la variable en fonction du temps, on préfère considérer le temps comme une variable réelle explicative.

Soit Y le CA trimestriel en KE et T la variable temps (réindexé de 1 à 8). On représente les points (y_i, t) pour $t = 1, \dots, 8$ dans un repère orthogonal.

Les moyennes mobiles d'ordre 3 sont données dans le tableau suivant et représentées sur le graphique.

Temps	1	2	3	4	5	6	7	8
CA	620	640	620	600	660	640	620	640
Moy. Mob.	/	627	620	627	633	640	633	/



Les moyennes mobiles sont une technique de lissage permettant de repérer la tendance d'une série chronologique ayant de fortes perturbations. Nous observons ici une légère croissance linéaire.

c) Soit $y = at + b$, $t = 1, \dots, 8$, l'équation de la droite de régression du CA sur le temps.

Si T désigne la variable temps, on a $a = \frac{\text{COV}(Y, T)}{\text{V}(T)}$ et $b = \bar{Y} - a\bar{T}$.

La droite de régression linéaire du CA par rapport au temps passe par le point moyen (appelé aussi centre de gravité du nuage de points pondérés par les effectifs) $G(\bar{T}, \bar{Y})$. La droite est représentée sur le graphique.

Les calculs sont présentés dans le tableau suivant ; ils sont réalisés à partir du changement de variable

$$U = \frac{Y - 620}{20}.$$

t	y_t	u_t	t^2	u_t^2	$t u_t$
1	620	0	1	0	0
2	640	1	4	1	2
3	620	0	9	0	0
4	600	-1	16	1	-4
5	660	2	25	4	10
6	640	1	36	1	6
7	620	0	49	0	0
8	640	1	64	1	8
Σ	36	4	204	8	22
$(1/8)\Sigma$	4.5	0.5	25.5	1	2.75

On a : $\bar{T} = 4.5$, $V(T) = 25.5 - (4.5)^2 = 5.25$, $\bar{U} = 0.5$, $V(U) = 1 - (0.5)^2 = 0.75$, $\sigma(U) = 0.87$

$COV(U, T) = 2.75 - (4.5)(0.5) = 0.5$

On en déduit : $\bar{Y} = 20\bar{U} + 620 = 630$, $V(Y) = (20)^2 V(U) = 300$, $\sigma(Y) = 17.32$

$COV(Y, T) = 20COV(U, T) = 10$.

Enfin on a : $a = \frac{COV(Y, T)}{V(T)} = \frac{10}{5.25} = 1.9$ et $b = \bar{Y} - a\bar{T} = 630 - 1.9 \times 4.5 = 621.43$

L'équation de la droite de régression linéaire du CA par rapport au temps est donc $y = 1.9t + 621.43$. La droite passe par le point moyen G de coordonnées (4.5, 630) et le coefficient directeur correspond à l'accroissement du CA pour une unité de temps, soit un accroissement de 1.9 KE en moyenne par trimestre.

- d) La prévision du CA trimestriel pour le 3^{ème} et le 4^{ème} trimestre 2003 est obtenue à partir de la droite de régression linéaire du CA sur le temps, en posant $t = 9$ et $t = 10$.

Prévisions de Y: 3^{ème} trimestre 03 $\hat{Y} = 1.9 \times 9 + 621.43 = 638.53$

4^{ème} trimestre 03 $\hat{Y} = 1.9 \times 10 + 621.43 = 640.43$

Ces prévisions sont faites sans tenir compte des fortes fluctuations d'un trimestre à l'autre. L'augmentation de 1.9 KE en moyenne par trimestre est faible par rapport à ces fluctuations. Le coefficient de corrélation linéaire entre le CA et le temps est de 0.25 ce qui montre une corrélation positive mais bien faible. Les prévisions sont donc à prendre avec précaution.

Ventes trimestrielles d'un produit (cf. annexe jeu de données 3)

t	$Y(t)$	$Y(t) - Y(t-1)$	$Y(t) / Y(t-1)$
1	255	/	/
2	330	75	1.29
3	435	105	1.32
4	570	130	1.31
5	740	170	1.30
6	960	220	1.30

b/ Le modèle exponentiel est préférable au modèle linéaire car les rapports $Y(t) / Y(t-1)$ sont à peu près constants.

On pose donc : $Y(t) = b a^t$

On pose $X(t) = \ln Y(t)$, $A = \ln a$, $B = \ln b$

On a alors : $X(t) = At + B$ avec $A = \frac{COV(X, T)}{V(T)}$ et $B = \bar{X} - A\bar{T}$

t	$X(t)$	$\bar{X}'(t) = \ln X(t)$	$t X'(t)$
1	255	5.64	5.54
2	330	5.80	11.60
3	435	6.08	18.24
4	570	6.35	25.40
5	740	6.61	33.05
6	960	6.87	41.22
	Σ	37.25	135.05
	$(1/6) \Sigma$	6.21	22.51

On a :

$$\bar{T} = \frac{6+1}{2} = 3.5 \quad V(T) = \frac{6^2 - 1}{12} = 2.92$$

$$\bar{X} = 6.21$$

$$COV(X, T) = 22.51 - (3.5)(6.21) = 0.775$$

$$A = \frac{0.775}{2.92} = 0.265 \text{ d'où } a = e^A = 1.3$$

$$B = 6.21 - (0.265)(3.5) = 5.28 \text{ d'où } b = e^B = 197$$

Il vient finalement : $Y(t) = 197(1.3)^t$

Les prévisions pour les trimestres 7 et 8 sont :

$$Y(7) = 197(1.3)^7 = 1236 \text{ unités vendues et } Y(8) = 197(1.3)^8 = 1607 \text{ unités vendues.}$$

Quelques propriétés

Définitions de la covariance et du coefficient de corrélation linéaire.
Rappel des propriétés usuelles de ces indices.

- **Deuxième formule de la covariance**

Données brutes

$$COV(X, Y) = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X})(Y_k - \bar{Y}) = \frac{1}{n} \sum_{k=1}^n (X_k Y_k - \bar{X} Y_k - X_k \bar{Y} + \bar{X} \bar{Y}) = \frac{1}{n} \sum_{k=1}^n X_k Y_k - \bar{X} \bar{Y}$$

Distribution de fréquences

$$COV(X, Y) = \sum_{i=1}^l f_i (x_i - \bar{X})(y_i - \bar{Y}) = \sum_{i=1}^l f_i (x_i y_i - \bar{X} y_i - x_i \bar{Y} + \bar{X} \bar{Y}) = \sum_{i=1}^l f_i x_i y_i - \bar{X} \bar{Y}$$

- **Propriétés de la covariance**

$$COV(X, Y) = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X})(Y_k - \bar{Y}) = COV(Y, X), \text{ la covariance est symétrique}$$

$$COV(aX + bY, Z) = \frac{1}{n} \sum_{k=1}^n (aX_k + bY_k - a\bar{X} - b\bar{Y})(Z_k - \bar{Z})$$

$$= a \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X})(Z_k - \bar{Z}) + b \frac{1}{n} \sum_{k=1}^n (Y_k - \bar{Y})(Z_k - \bar{Z}) = a COV(X, Z) + b COV(Y, Z)$$

la covariance est linéaire par rapport au premier argument ; étant symétrique la covariance est linéaire par rapport au deuxième argument et donc bilinéaire.

$COV(X, X) = V(X) \geq 0$ la covariance est positive et $V(X) = 0$ si et seulement si X est constante.

- **Propriétés du coefficient de corrélation linéaire**

$$1) \text{ On a } \rho(X, Y) = \frac{COV(X, Y)}{\sigma(X)\sigma(Y)} = \frac{COV(Y, X)}{\sigma(Y)\sigma(X)} = \rho(Y, X).$$

Le coefficient de corrélation linéaire est symétrique.

2) Pour tout réel a on a vu que l'on avait :

$$V(Y - aX) = a^2 V(X) - 2a \text{COV}(X, Y) + V(Y) \quad (= f(a))$$

Comme une variance est positive, pour tout a , le trinôme du second degré $f(a)$ est positif.

On en déduit que le discriminant $\text{COV}^2(X, Y) - V(X)V(Y)$ est négatif ou nul,

$$\text{c-à-d } \text{COV}^2(X, Y) \leq V(X)V(Y), \quad \frac{\text{COV}^2(X, Y)}{V(X)V(Y)} \leq 1, \quad \left| \frac{\text{COV}(X, Y)}{\sigma(X)\sigma(Y)} \right| \leq 1, \quad |\rho(X, Y)| \leq 1,$$

soit $-1 \leq \rho(X, Y) \leq 1$

Par ailleurs $f(a)$ atteint son minimum pour $a = \frac{\text{COV}(X, Y)}{V(X)}$ et, pour cette valeur de a ,

on obtient : $f(a) = \frac{V(X)V(Y) - (\text{COV}(X, Y))^2}{V(X)}$; on en déduit :

$$V(Y - aX) = 0 \Leftrightarrow \exists b \in \mathbf{R}, Y = aX + b \Leftrightarrow V(X)V(Y) = [\text{COV}(X, Y)]^2 \Leftrightarrow |\rho(X, Y)| = 1$$

3) On peut préciser :

$$\rho(X, Y) = 1 \Leftrightarrow Y = aX + b, \quad a > 0$$

$$\rho(X, Y) = -1 \Leftrightarrow Y = aX + b, \quad a < 0$$

- **Coefficients de la droite de régression linéaire de Y en X**

Montrer que la fonction f de deux variables réelles a et b définie ci-après est minimum pour

$$a = \frac{\text{COV}(X, Y)}{V(X)} \quad \text{et} \quad b = \bar{Y} - a\bar{X}$$

$$\begin{aligned} f(a, b) &= \frac{1}{n} \sum_{k=1}^n (Y_k - aX_k - b)^2 = \frac{1}{n} \sum_{k=1}^n (Y_k - \bar{Y} - aX_k + a\bar{X} + \bar{Y} - a\bar{X} - b)^2 \\ &= \frac{1}{n} \sum_{k=1}^n (Y_k - \bar{Y} - aX_k + a\bar{X})^2 + \frac{2}{n} \sum_{k=1}^n (Y_k - \bar{Y} - aX_k + a\bar{X})(\bar{Y} - a\bar{X} - b) + \frac{1}{n} \sum_{k=1}^n (\bar{Y} - a\bar{X} - b)^2 \\ &= \frac{1}{n} \sum_{k=1}^n (Y_k - \bar{Y} - aX_k + a\bar{X})^2 + (\bar{Y} - a\bar{X} - b)^2 \quad (\text{car le double produit est nul}). \end{aligned}$$

C'est la somme de 2 termes positifs dont le second est minimum (car nul) pour $b = \bar{Y} - a\bar{X}$

Il s'agit donc de trouver a minimisant le premier terme ; il s'écrit :

$$\frac{1}{n} \sum_{k=1}^n [(Y_k - \bar{Y})^2 - 2a(Y_k - \bar{Y})(X_k - \bar{X}) + a^2(X_k - \bar{X})^2]$$

c'est-à-dire : $a^2 V(X) - 2a \text{COV}(X, Y) + V(Y)$

Ce trinôme du second degré en a admet un minimum pour $a = \frac{\text{COV}(X, Y)}{V(X)}$ d'où la solution.

Annexe : jeu de données

1) Loyer mensuel de locaux commerciaux

Le loyer mensuel en euros et la superficie en m² de 10 locaux commerciaux d'un quartier sont donnés dans le tableau suivant :

Loyer mensuel (en €)	Superficie (en m ²)
2 395	24
1 265	31
4 050	45
1 730	50
4 350	60
3 300	66
3 600	72
4 900	70
3 600	90
5 760	100

- Décrire les données et représenter les dans un repère cartésien, la superficie en abscisse, le loyer en ordonnée.
- Calculer moyenne et écart-type de chacune des deux variables, covariance et coefficient de corrélation linéaire des deux variables. Commenter.
- Calculer les coefficients a et b de la droite de régression linéaire, d'équation $y = a x + b$, du loyer Y par rapport à la superficie X . Représenter sur le graphique le point moyen et cette droite de régression linéaire. Interpréter le coefficient directeur de la droite.
- Donner une prévision du loyer mensuel net en euro d'un local commercial du quartier de 75 m².
- On note $\hat{Y} = aX + b$ le modèle et $\hat{E} = X - \hat{Y}$ l'erreur.
Calculer les valeurs de ces variables sur les 10 observations.
Calculer moyenne, variance, écart-type de chacune de ces deux variables, covariance et coefficient de corrélation linéaire des deux variables.
Écrire le carré du coefficient de corrélation linéaire de X et Y en fonction des variances de \hat{Y} et de Y . Commenter.
(On renvoie à la correction pour des compléments sur la régression linéaire).

2) Chiffre d'affaires (CA) trimestriel d'une entreprise

Le CA en milliers d'euros d'une entreprise est observé sur deux années consécutives.

3° Trim 2001	4° Trim 2001	1° Trim 2002	2° Trim 2002	3° Trim 2002	4° Trim 2002	1° Trim 2003	2° Trim 2003
620	640	620	600	660	640	620	640

- Décrire les données et représenter les dans un repère cartésien, le temps (noté T) en abscisse et le CA (noté Y) en ordonnée. Joignez les points consécutifs.
- Calculer les moyennes mobiles d'ordre 3. Représenter les points correspondants sur le graphique. Commenter.
- Calculer les coefficients a et b de la droite de régression linéaire, d'équation $y = a t + b$, du CA Y par rapport au temps T (on indexera de 1 à 8 les valeurs de T). Représenter sur le graphique le point moyen et cette droite de régression linéaire. Interpréter le coefficient directeur de la droite.
- Donner une prévision du CA pour les 3^{ème} et 4^{ème} trimestres de l'année 2003. Commenter.

3) Ventes trimestrielles d'un produit

L'effectif des ventes d'un produit est donné dans le tableau suivant.

Trim.	1	2	3	4	5	6
Quantité vendue	255	330	435	570	740	960

- Décrire les données et représenter les dans un repère cartésien, le temps (noté T) en abscisse et le nombre de produits vendus (noté Y) en ordonnée. Joignez les points consécutifs.
- Pour $t = 2$ à 6 , calculer $Y(t) - Y(t-1)$ et $Y(t)/Y(t-1)$. Entre le modèle linéaire $Y(t) = at + b$ et le modèle exponentiel $Y(t) = ba^t$, lequel pensez-vous le mieux adapté à cette série chronologique ?
- Ajuster le modèle choisi en b) et prévoir les trimestres 7 et 8.