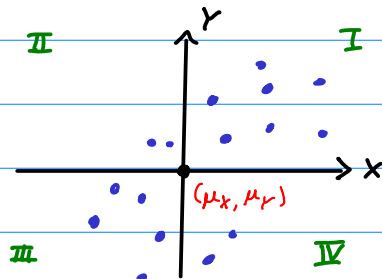


Régression linéaire

But : Établir un lien linéaire entre deux variables aléatoires X et Y .

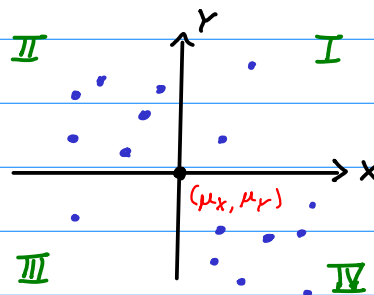
§ Retour sur la covariance.

$$\text{Cov}(X, Y) \stackrel{\text{def}}{=} \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y).$$



$\text{Cov}(X, Y) > 0$ tendance à se trouver dans les quadrants I et III.

$$\text{ou } \mu_X = \mathbb{E}(X), \mu_Y = \mathbb{E}(Y).$$



$\text{Cov}(X, Y) < 0$ tendance à se trouver dans les quadrants II et IV.

Estimateur de la covariance : soit $(X_1, Y_1), \dots, (X_n, Y_n)$ un échantillon iid d'une population de distribution (X, Y) . Il est naturel de prendre

$$\widehat{\text{Cov}}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}), \text{ où } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

Vérifier qu'il est sans biais. D'abord, on a

$$\begin{aligned} \mathbb{E}[(X_i - \bar{X})(Y_i - \bar{Y})] &= \mathbb{E}[(X_i - \mu_X + (\mu_X - \bar{X}))((Y_i - \mu_Y) + (\mu_Y - \bar{Y}))] \\ &= \mathbb{E}[(X_i - \mu_X)(Y_i - \mu_Y)] - \mathbb{E}[(X_i - \mu_X)(\bar{Y} - \mu_Y)] \\ &\quad - \mathbb{E}[(\bar{X} - \mu_X)(Y_i - \mu_Y)] + \mathbb{E}[(\bar{X} - \mu_X)(\bar{Y} - \mu_Y)] \\ &= \text{Cov}(X_i, Y_i) - \text{Cov}(X_i, \bar{Y}) - \text{Cov}(\bar{X}, Y_i) + \text{Cov}(\bar{X}, \bar{Y}) \end{aligned}$$

$$\textcircled{1} \text{Cov}\left(X_i, \frac{1}{n} \sum_{j=1}^n Y_j\right) = \frac{1}{n} \sum_{j=1}^n \text{Cov}(X_i, Y_j) = \frac{1}{n} (0 + \dots + 0 + \text{Cov}(X_i, Y_i) + 0 + \dots + 0)$$

*Cov(X_i, Y_j) = 0 pour i ≠ j
car elles sont indépendantes*

$$= \frac{1}{n} \text{Cov}(X_i, Y_i) = \frac{1}{n} \text{Cov}(X, Y)$$

② Comme le ①

$$\textcircled{3} \text{Cov}(\bar{X}, \bar{Y}) = \text{Cov}\left[\frac{1}{n} \sum_{j=1}^n X_j, \frac{1}{n} \sum_{m=1}^n Y_m\right] = \frac{1}{n^2} \sum_{j=1}^n \sum_{m=1}^n \text{Cov}(X_j, Y_m)$$

$$= \frac{1}{n^2} \sum_{j=1}^n \text{Cov}(X_j, Y_j) = \frac{1}{n} \text{Cov}(X, Y)$$

On a donc

$$\mathbb{E}[(X_i - \bar{X})(Y_i - \bar{Y})] = \text{Cov}(X, Y) - \frac{1}{n} \text{Cov}(X, Y) - \frac{1}{n} \text{Cov}(X, Y) + \frac{1}{n} \text{Cov}(X, Y)$$

$$= \frac{n-1}{n} \text{Cov}(X, Y)$$

Enfin, on a

$$\mathbb{E}[\hat{\text{cov}}(X, Y)] = \mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})\right]$$

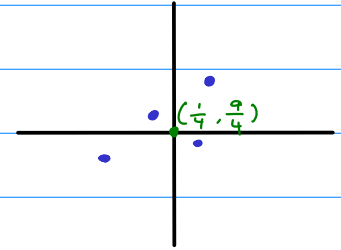
$$= \frac{1}{n-1} \sum_{i=1}^n \mathbb{E}[(X_i - \bar{X})(Y_i - \bar{Y})]$$

$$= \frac{1}{n-1} \sum_{i=1}^n \frac{n-1}{n} \text{Cov}(X, Y) = \text{Cov}(X, Y).$$

Pour calculer la covariance empirique d'un échantillon $(x_1, y_1), \dots, (x_n, y_n)$, on fera $\hat{\text{cov}}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$.

Ex. (1,3), (1,2), (-1,1), (0,3)

$$\bar{x} = \frac{1+1-1+0}{4} = \frac{1}{4} \quad \bar{y} = \frac{3+2+1+3}{4} = \frac{9}{4}$$



i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
1	0,75	0,75	0,5625
2	0,75	-0,25	-0,1875
3	-1,25	-1,25	1,5625
4	-0,25	0,75	-0,1875
			1,75 Somme

$$\widehat{\text{Cov}}(\bar{x}, \bar{y}) = \frac{1,75}{3} = 0,58\bar{3}$$

§ Corrélation

$$\rho = \text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

Rappel : $-1 \leq \text{Cor}(X, Y) \leq 1$ et $\text{Cor}(X, Y) = \pm 1$ si $\exists a, b, \mathbb{P}(Y = aX + b) = 1$

$\text{Cor}(X, Y)$ est une mesure ^{sans unités} de la force du lien linéaire entre X et Y .

Un estimateur naturel de la corrélation est

$$\hat{r} = \widehat{\text{Cor}}(X, Y) = \frac{\widehat{\text{Cov}}(X, Y)}{\hat{\sigma}_X \hat{\sigma}_Y} = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Les 1/n-1 se simplifient

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

mais celui-ci est biaisé.

Il est difficile de trouver un estimateur simple et sans biais. En pratique, on utilise plutôt

$$\hat{\rho} = r + \frac{r(1-r^2)}{2(n-4)}$$

qui est peu près sans biais lorsque $n \geq 18$ et la distribution de la population est normalement distribuée.

Pour le cours, on se contentera de \hat{r} , qui est plus simple à calculer.

Exemple de calcul Calculer la corrélation empirique de l'échantillon

(1,-1) (0,1) (-1/2, 3) (1,2) (1/2, 1/2)

pas besoin de
diviser par $\frac{1}{n-1}$

Il faut calculer : \bar{x} , \bar{y} , $\sum(x_i - \bar{x})^2$, $\sum(y_i - \bar{y})^2$, $\sum(x_i - \bar{x})(y_i - \bar{y})$

$$\bar{x} = \frac{1+0-\frac{1}{2}+1+\frac{1}{2}}{5} = \frac{2}{5} = 0,4 \quad \bar{y} = \frac{-1+1+3+2+\frac{1}{2}}{5} = \frac{5,5}{5} = 1,1$$

i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	$y_i - \bar{y}$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
1	0,6	0,36	-2,1	4,41	-1,26
2	-0,4	0,16	-0,1	0,01	0,04
3	-0,9	0,81	1,9	3,61	-1,71
4	0,6	0,36	0,9	0,81	0,54
5	0,1	0,01	-0,6	0,36	-0,06
Σ		1,7		9,2	-2,45

$$\hat{r} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{(\sum(x_i - \bar{x})^2)(\sum(y_i - \bar{y})^2)}} = \frac{-2,45}{\sqrt{1,7 \times 9,2}} = -0,6195$$

§ Loi normale multidimensionnelle

On fait seulement le cas en dimension 2.

On définit la matrice de covariance d'un vecteur aléatoire (X, Y) par

$$\Sigma = \begin{pmatrix} \text{Cov}(X, X) & \text{Cov}(X, Y) \\ \text{Cov}(Y, X) & \text{Cov}(Y, Y) \end{pmatrix} = \begin{pmatrix} \text{Var}(X) & \text{Cov}(X, Y) \\ \text{Cov}(X, Y) & \text{Var}(Y) \end{pmatrix}$$

C'est une matrice symétrique. Son déterminant est

$$\begin{aligned} \det \Sigma &= \text{Var}(X)\text{Var}(Y) - \text{Cov}(X,Y)^2 \\ &= \text{Var}(X)\text{Var}(Y) \left[1 - \frac{\text{Cov}(X,Y)^2}{\text{Var}(X)\text{Var}(Y)} \right] \\ &= \text{Var}(X)\text{Var}(Y) (1 - \text{Cor}(X,Y)^2) \end{aligned}$$

La loi normale multidimensionnelle $N(\vec{\mu}, \Sigma)$ est distribuée selon la fonction de densité

$$f(\vec{x}) = \frac{1}{2\pi \sqrt{\det \Sigma}} e^{-\frac{1}{2} \vec{x}^T \Sigma^{-1} \vec{x}} \quad \text{où } \vec{x} = \begin{pmatrix} x \\ y \end{pmatrix}$$

après simplification

$$f(x,y) = \frac{1}{2\pi \sigma_x \sigma_y \sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left(\frac{(x-\mu_x)^2}{\sigma_x^2} - \frac{2\rho(x-\mu_x)(y-\mu_y)}{\sigma_x \sigma_y} + \frac{(y-\mu_y)^2}{\sigma_y^2} \right)}$$

Dans le cas où $\mu_x = 0$ et $\mu_y = 0$, la densité devient

$$f(x,y) = \frac{1}{2\pi \sigma_x \sigma_y \sqrt{1-\rho^2}} \exp\left(-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}\right)$$

La corrélation apparaît de façon naturelle dans cette distribution importante.

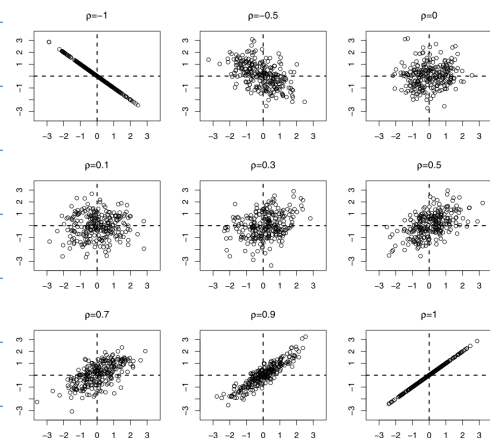
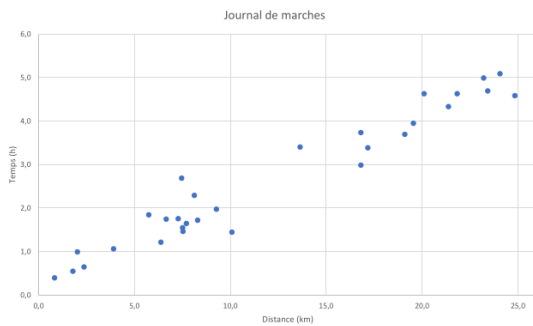


Diagramme de dispersion de $N(\vec{0}, \Sigma)$ avec différentes valeurs de corrélation ρ .

(Source: Rousson, V. Statistique appliquée aux sciences de la vie)

§ Régression linéaire.

Un diagramme de dispersion est un nuage de points dans le plan-XY où chaque point (x, y) est une donnée observée.



(Source : Note de cours
d'Élise Davignon)

Soit (X, Y) un vecteur aléatoire. On cherche une relation linéaire (affine) entre X et Y : $Y = aX + b$. On appelle X la variable explicative et Y , la variable de réponse.

On a vu que si $\text{Cor}(X, Y) = \pm 1$, alors $Y = aX + b$. En général, cette relation ne sera pas parfaite. On a plutôt

$$Y = aX + b + e, \text{ où } e \text{ est un terme d'erreur aléatoire appelé } \underline{\text{résidu}}$$

Le but est de trouver a et b (des constantes) qui minimisent l'erreur.

• Cas parfait : si $e = 0$, on a $\text{Cov}(X, Y) = a \text{Cov}(X, X) = b \text{Var}(X) \Rightarrow a = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$

$$\text{et } \mathbb{E}(Y) = a \mathbb{E}(X) + b \Rightarrow b = \mathbb{E}(Y) - a \mathbb{E}(X).$$

• Cas général : on s'attend à trouver la même sol que le cas parfait.

Cependant, les calculs sont plus compliqués, vu qu'on ne connaît pas e .

Soit $(X_1, Y_1), \dots, (X_n, Y_n)$ un échantillon iid. On $Y_i = aX_i + b + e_i$.
 Le but est de minimiser l'erreur. On utilise l'erreur quadratique :

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - aX_i - b)^2 \quad (*)$$

Si a et b minimisent l'erreur quadratique, alors en dérivant $(*)$, on obtient

$$(1) \quad \sum_{i=1}^n 2(Y_i - aX_i - b)(-X_i) = 0 \quad (\text{en dérivant par rapport à } a)$$

$$(2) \quad \sum_{i=1}^n 2(Y_i - aX_i - b)(-1) = 0 \quad (\text{en dérivant par rapport à } b)$$

$$(2) \Leftrightarrow 0 = \sum_{i=1}^n (Y_i - aX_i - b) = \sum_{i=1}^n Y_i - a \sum_{i=1}^n X_i - nb$$

$$\Leftrightarrow b = \frac{1}{n} \sum_{i=1}^n Y_i - a \frac{1}{n} \sum_{i=1}^n X_i = \bar{Y} - a\bar{X}$$

$$(1) \quad 0 = \sum_{i=1}^n (X_i Y_i - a X_i^2 - b X_i) = \sum_{i=1}^n X_i Y_i - a \sum_{i=1}^n X_i^2 - b \sum_{i=1}^n X_i$$

$$= \sum_{i=1}^n X_i Y_i - a \sum_{i=1}^n X_i^2 - (\bar{Y} - a\bar{X}) n\bar{X}$$

$$= \left(\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y} \right) - a \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2 \right)$$

$$= n \underbrace{\left(\frac{1}{n} \sum_{i=1}^n X_i Y_i - \bar{X}\bar{Y} \right)}_{\text{Cov empirique}} - a n \underbrace{\left(\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 \right)}_{\text{Variance empirique}}$$

$$\Rightarrow a = \frac{\widehat{\text{Cov}}(X, Y)}{\widehat{\sigma}_X^2} \quad \text{si } \widehat{\sigma}_X^2 \neq 0$$

$$\text{Conclusion: } a = \frac{\widehat{\text{Cov}}(X, Y)}{\widehat{\sigma}_x^2} = \frac{\sum_{i=1}^n X_i Y_i - n \bar{X} \bar{Y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$b = \bar{Y} - a \bar{X}$$

§ Distribution des estimateurs de a et b.

On peut voir a et b comme des paramètres et on peut utiliser les formules précédentes comme estimateurs :

$$\hat{a} = \frac{\widehat{\text{Cov}}(X, Y)}{\sigma_x^2} \quad \hat{b} = \bar{Y} - \hat{a} \bar{X}$$

On fait notre analyse dans le cas où on a observé les X_i , c'est-à-dire $X_i = x_i$, et l'on veut prédire Y_i à partir de x_i . On suppose également que e_1, \dots, e_n sont iid avec distribution $N(0, \sigma^2)$.

Comme $Y_i = ax_i + b + e_i$, il suit que $Y_i \sim N(ax_i + b, \sigma^2)$.

Distribution a:
$$\hat{a} = \frac{\widehat{\text{Cov}}(X, Y)}{\widehat{\sigma}_x^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{S_x^2}$$

$$\sum_{i=1}^n (x_i - \bar{x}) \bar{Y} = (n\bar{x} - n\bar{x}) \bar{Y} = 0$$

$$= \sum_{i=1}^n \frac{(x_i - \bar{x}) Y_i}{S_x^2}$$

où $S_x^2 = \sum_{i=1}^n (x_i - \bar{x})^2 = (n-1) \sigma_x^2$

C'est une somme de distributions normales indépendantes, donc \hat{a} suit une loi normale. On a

$Y = ax + b + e$
 $E(Y) = ax + b$

$$E(\hat{a}) = \sum_{i=1}^n \frac{(x_i - \bar{x})}{S_x^2} (E(Y_i) - E(\bar{Y})) = \sum_{i=1}^n \frac{(x_i - \bar{x})}{S_x^2} (ax_i + b - a\bar{x} - b)$$

$$= \frac{a}{S_x^2} \underbrace{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}_{S_x^2} = a$$

Donc \hat{a} est sans biais.

Pour la variance (= EQM), on a

$Y_i \sim N(ax_i + b, \sigma^2)$

$$\text{Var}(\hat{a}) = \frac{1}{\sigma_x^4} \sum_{i=1}^n (x_i - \bar{x})^2 \text{Var}(Y_i) = \frac{\sigma^2}{S_x^4} \underbrace{\sum_{i=1}^n (x_i - \bar{x})^2}_{= S_x^2} = \frac{\sigma^2}{S_x^2}$$

Conclusion : $\hat{a} \sim N(a, \frac{\sigma^2}{S_x^2})$.

Distribution de \hat{b} : $\hat{b} = \bar{Y} - \hat{a} \bar{x}$

\hat{b} est une combinaison linéaire des Y_i , donc \hat{b} suit une loi normale.

$$\begin{aligned} E(\hat{b}) &= \frac{1}{n} \sum_{i=1}^n E(Y_i) - E(\hat{a}) \bar{x} \\ &= b + \frac{a}{n} \sum_{i=1}^n x_i - a \bar{x} \\ &= b \end{aligned}$$

donc \hat{b} est sans biais

$$\begin{aligned} \text{Var}(\hat{b}) &= \text{Var}(\bar{Y}) + \bar{x}^2 \text{Var}(\hat{a}) \\ &= \frac{\sigma^2}{n} + \bar{x}^2 \frac{\sigma^2}{S_x^2} \\ &= \sigma^2 \left(\frac{S_x^2 + n \bar{x}^2}{n S_x^2} \right) \\ &= \frac{\sigma^2}{n S_x^2} \sum_{i=1}^n x_i^2 \end{aligned}$$

Donc $\hat{b} \sim N\left(b, \frac{\sigma^2}{n S_x^2} \sum_{i=1}^n x_i^2\right)$

De plus, on sait que $e_i = Y_i - aX_i - b$.

Estimation de la variance inconnue σ^2

On sait que les $e_i \sim N(0, \sigma^2)$. Ainsi, $\frac{e_i}{\sigma} \sim N(0, 1)$, donc

$$\sum_{i=1}^n \frac{e_i^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n e_i^2 \text{ suit une loi du } \chi^2. \text{ Il y aurait habituellement } n$$

degrés de liberté, mais nous allons utiliser la relation $e_i = Y_i - aX_i - b$, où a et b sont

inconnus. Si on remplace un paramètre par son estimateur, on perd un degré de liberté (c'est un principe que nous accepterons). Dans notre cas, avec deux paramètres, on perd deux degrés de liberté, donc $\frac{1}{\sigma^2} \sum_{i=1}^n e_i^2 \sim \chi_{n-2}^2$.

Il suit que $E\left(\frac{1}{\sigma^2} \sum_{i=1}^n e_i^2\right) = n-2 \Rightarrow E\left(\frac{1}{n-2} \sum_{i=1}^n e_i^2\right) = \sigma^2$.

Après des calculs fastidieux, on peut obtenir

$$\sum_{i=1}^n e_i^2 = \frac{S_x^2 S_y^2 - S_{xy}^2}{S_x^2} = (n-1) \frac{\hat{\sigma}_x^2 \hat{\sigma}_y^2 - \widehat{\text{Cov}}(X, Y)^2}{\hat{\sigma}_x^2}$$

où $S_x^2 = (n-1)\hat{\sigma}_x^2$, $S_y^2 = (n-1)\hat{\sigma}_y^2$ et $S_{xy} = (n-1)\widehat{\text{Cov}}(X, Y)$

Coefficient de détermination: $R^2 = \frac{S_y^2 - S_e^2}{S_y^2} = \frac{\hat{\sigma}_y^2 - \hat{\sigma}_e^2}{\hat{\sigma}_y^2}$, où $S_y^2 = (n-1)\hat{\sigma}_y^2$
 $S_e^2 = (n-1)\hat{\sigma}_e^2$

Par sa définition, on voit que

$$\sigma_y^2 = R^2 \sigma_y^2 + \sigma_e^2$$

Variance expliquée par X
autre variance

$$\sigma_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

Heuristique: $Y = aX + b + e$

$$\Rightarrow \text{Var}(Y) = a^2 \text{Var}(X) + \text{Var}(e)$$

il y a une partie de la variance de Y qui est expliquée par la variance de X et une partie qui est expliquée par le fait que le modèle n'est pas parfait (le terme d'erreur)

Ainsi, R^2 est la proportion de la variance de Y qui est expliquée par la variance de X.

Propriété de R^2 : • $R^2 \in [0, 1]$ (car $R^2 = 1 - \frac{\sigma_e^2}{\sigma_y^2}$ et $\sigma_y^2, \sigma_e^2 \geq 0$)

• Si $R^2 = 0$, la variance de Y n'est pas linéairement expliquée par celle de X.

• Si $R^2=1$, la variance de Y est entièrement expliquée par celle de X

• $R = \sqrt{R^2} = \text{corrélation} = \frac{\widehat{\text{Cov}}(X, Y)}{\hat{\sigma}_x \sigma_Y}$

$$\left[\text{En effet, on a } R^2 = \frac{\hat{\sigma}_Y^2 - \hat{\sigma}_e^2}{\hat{\sigma}_Y^2} = \frac{\hat{\sigma}_Y^2 - \frac{\widehat{\text{Cov}}(X, Y)^2}{\sigma_x^2}}{\hat{\sigma}_Y^2} \right.$$

$$= \frac{\widehat{\text{Cov}}(X, Y)^2}{\sigma_x^2 \sigma_Y^2} = \text{Cor}(X, Y)^2$$

où on a utilisé $\hat{\sigma}_e^2 = \frac{1}{(n-1)} \sum_{i=1}^n e_i^2 = (n-1) \frac{\hat{\sigma}_x^2 \hat{\sigma}_Y^2 - \widehat{\text{Cov}}(X, Y)^2}{\hat{\sigma}_x^2}$

Ex.



(pris dans les notes de cours d'Élise Daignon)

Dans cet exemple, on trouve $R^2 = 0,944$, donc 94% de la variance dans le temps de marche est expliquée par la distance parcourue.

La droite en orange est la droite de régression $Y = \hat{a}X + \hat{b}$, où \hat{a} et \hat{b} sont calculées à partir des données en utilisant les formules de leur estimateur.

Inférence statistique sur a et b

On peut faire des tests d'hypothèse sur les valeurs de a et b .

P.ex. $H_0: a = a_0$ $H_1: a \neq a_0$.

Dans ce cas, si l'erreur suit une loi normale, alors $\hat{a} \sim N(a_0, \frac{\sigma^2}{S_x^2})$

On ne connaît pas la variance σ^2 de e , donc on utilise $\hat{\sigma}_e^2$:

$$T_{\text{stat}} = \frac{\hat{a} - a_0}{\hat{\sigma}_e / (n-2) S_x} \sim \text{Student}(n-2).$$

Ex. On prend l'échantillon de données (distance (km), temps de marche (h)) = (x, y) :

(4, 1) (7, 1,5) (8, 1,6) (17, 3)

On s'intéresse à la régression $Y = aX + b$ et à l'hypothèse $H_0: a = 0,2$.

On utilise $T_{\text{stat}} = \frac{\hat{a} - 0,2}{\hat{\sigma}_e / (n-2) S_x}$. Pour calculer t_{stat} , il faut calculer

le \hat{a} empirique, $\hat{\sigma}_e^2$ et S_x^2 . Pour calculer $\hat{\sigma}_e^2$, il faut aussi calculer \hat{b} empirique.

On rappelle que $\hat{a} = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$ et $\hat{b} = \bar{Y} - \hat{a} \bar{X}$.

$$1. \bar{x} = \frac{4 + 7 + 8 + 17}{4} = \frac{36}{4} = 9 \quad \bar{y} = 1,775.$$

2. $\sigma_x^2, \sigma_y^2, \text{Cov}(X, Y)$:

i	$(x_i - \bar{x})$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
1	-5	25	-0,775	0,6	3,875
2	-2	4	-0,275	0,075	0,55
3	-1	1	-0,175	0,306	0,175
4	8	64	1,225	1,5	9,8
Σ		94		2,2075	14,4

$$\sigma_x^2 = \frac{94}{3} = 31,3 \quad \sigma_y^2 = \frac{2,2075}{3} = 0,735 \quad \text{Cov}(X, Y) = \frac{14,4}{3} = 4,8$$

$$\sigma_x = 5,6$$

$$\sigma_y = 0,857$$

$$3. \hat{a} = \frac{4,8}{5,6 \times 0,857} = 1 \quad \hat{b} = 1,775 - 1 \times 9 = -7,225$$

$$4. e_i = y_i - \hat{a}x_i - b$$

i	e_i	e_i^2
1	4,225	17,85
2	1,725	2,97
3	0,825	0,68
4	-6,775	45,9
Σ	\times	67,4

$$\hat{\sigma}_e^2 = \frac{67,4}{3} = 22,47$$

$$\hat{\sigma}_e = 4,74.$$

$$5. t_{\text{stat}} = \frac{\hat{a} - 0,2}{\hat{\sigma}_a / (n-2) s_x} = \frac{1 - 0,2}{4,74 / (2 \times 9,69)} = 3,27.$$

$$s_x^2 = (n-1) \hat{\sigma}_x^2$$

$$= 3 \times \hat{\sigma}_x^2 = 94$$

$$s_x = 9,69$$

$$dL = n - 2 = 4 - 2 = 2$$

$$6. T_{\text{stat}} \sim \text{Student}(2)$$

$$\text{valeur-p} = \mathbb{P}(|T_{\text{stat}}| \geq t_{\text{stat}} | H_0) = 0,082.$$

7. Avec un seuil $\alpha = 5\%$, on a $8,2\% \geq \alpha$, donc on ne rejette pas H_0 .