

Lien entre deux variables quantitatives:

individu (x_i, y_i)

	X	Y	\hat{Y}	ε
1	x_1	y_1		
i	x_i	y_i		
...		
n	x_n	y_n		

avec $\hat{Y} = bx_i + a$
 $\varepsilon = Y_i - \hat{Y}_i$

Les calculs:

$$\bar{x}; \bar{y} \quad s^2x, s^2y$$

covariance: indicateur entre deux variables

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

si on calcule la covariance entre x et x on trouve la variance des x .

$$s_{xy} = \frac{1}{n} \sum x_i y_i - \bar{x} \bar{y}$$

La covariance peut être négative

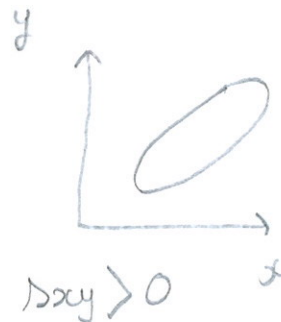
Le signe de la covariance indique la relation entre x et y .

graphique de dispersion:

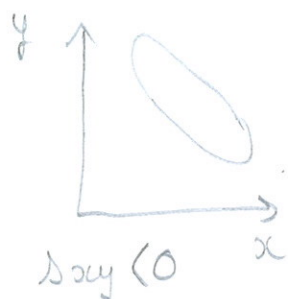
nuage de points.



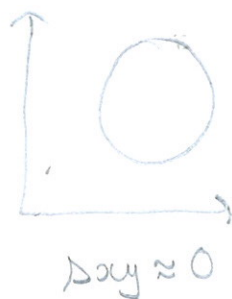
Le sens du nuage de point indique le lien entre x et y .



lorsque x et $y \uparrow$, liaison positive entre x et y .



quand $x \uparrow$ et $y \downarrow \rightarrow$ liaison négative entre x et y .

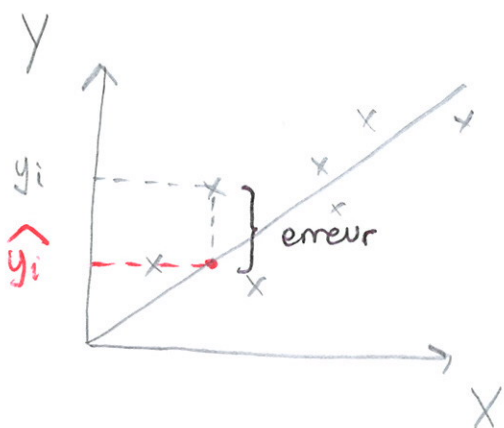


Si ce n'est pas linéaire il faut rechercher une équation (régression).

• régression linéaire:

Y variable à expliquer

X variable explicative : variable fixe régresse.



il faut chercher une fonction telle que $Y = f(X)$

Il n'y aura pas une fonction qui passe par tous les points : on cherche une approximation qui passe par le maximum de points.

$Y = bX + a$: on prend en compte les distances verticales (car X fixe) entre les points et la courbe. C'est l'erreur : $E_i = y_i - \hat{y}_i$

$Y = bX + a$ telle que $\sum E_i^2$ soit minimum (moindre carré ordinaire MKO)

$$S(b, a) = \sum (y_i - b x_i - a)^2$$

$$\begin{cases} \frac{\partial S}{\partial a} = 0 \\ \frac{\partial S}{\partial b} = 0 \end{cases}$$

$$\begin{aligned} b &= \frac{\Delta xy}{\Delta^2 x} \\ a &= \bar{y} - b \bar{x} \end{aligned}$$

\rightarrow ordonnée à l'origine

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

$$\hat{y}_i = b x_i + a$$

coefficient de corrélation:

$$r = \frac{\Delta xy}{b x \Delta y}$$

$r \in [-1, 1]$ donc $r \rightarrow 1, 0, -1$

r proche de 1: corrélation linéaire positive

_____ 0: \emptyset de corrélation linéaire

_____ -1: corrélation linéaire négative

Signe de r = signe de la covariance

• modélisation statistique:

$y_i = b x_i + a + \varepsilon_i$ le modèle

on appelle ε_i les résidus; $\varepsilon_i \sim N(0; \sigma_{\text{res}})$ et ils sont indépendants (on note $\varepsilon_i = \text{iid}$)

	\hat{y}	ε
moyenne		0
variance	$D^2 \hat{y}$	$D^2 \varepsilon$

on estime ensuite:

$$\sigma_{\text{res}} = \frac{n}{n-2} D^2 \varepsilon$$

décomposition de la variance et R^2 :

$V_{\text{totale}} = V_{\text{régression}} + V_{\text{résidus}}$

$$D^2 y = D^2 \hat{y} + D^2 \varepsilon$$

coefficient de détermination: $R^2 = \frac{D^2 \hat{y}}{D^2 y}$

ou a aussi $R^2 = r \times r$

R^2 est compris entre 0 et 1, car $r \in [0, 1]$

• analyse des résidus:

pour vérifier la normalité des résidus:

- histogrammes

- graphique quantil/quantil.

variance homogène: teste d'homoscédasticité.

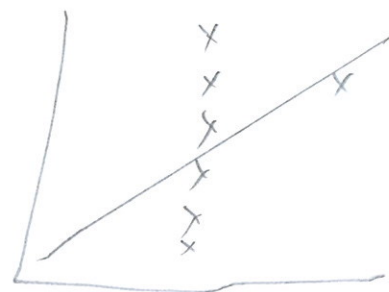
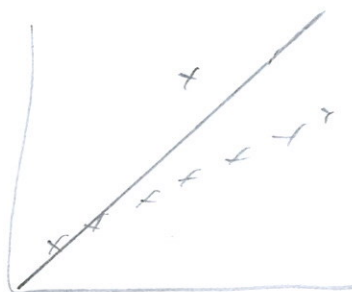
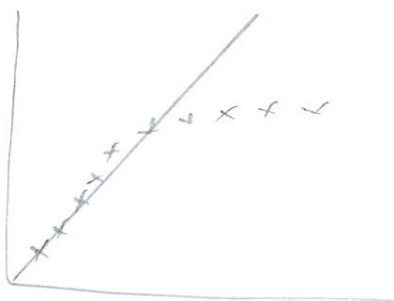
graphe des résidus:

$$\frac{\epsilon_i}{\sigma_{res}}$$

résidus standardisés



on va chercher les points externes aux limites
une corrélation n'est jamais la preuve d'un lien de causalité



coefficient de symétrie: γ_1 de Fisher

$$\gamma_1 = \frac{\mu_3}{\sigma_3}$$

- $< 0 \rightarrow$ étalement à gauche
- $= 0 \rightarrow$ symétrique
- $> 0 \rightarrow$ étalement à droite

coefficient de aplatissement: γ_2 de Fisher

$$\gamma_2 = \frac{\mu_4}{\sigma_4}$$

- < 0 + plate \rightarrow platycurtique
- $= 0$ normale \rightarrow mésocurtique
- > 0 pointue \rightarrow leptocurtique.