

Tester l'association linéaire entre deux variables quantitatives : corrélation, régression linéaire simple



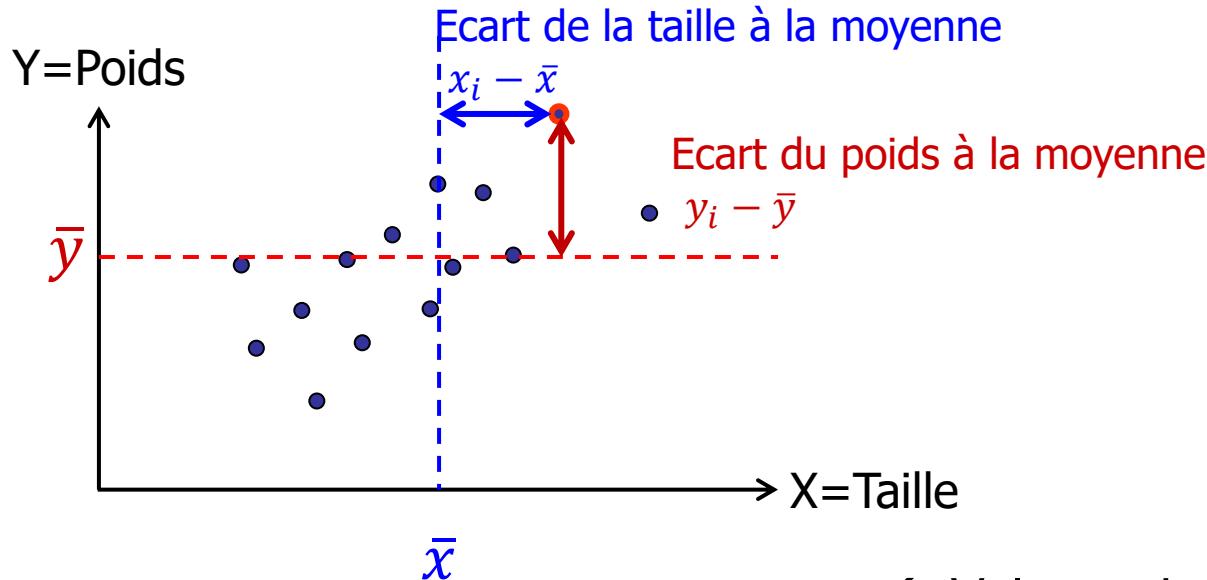
Introduction

- Etudier la relation bivariée, entre deux variables quantitatives X et Y, mesurées sur les mêmes individus
- Exemples :
 - Phénomènes différents mesurées chez les mêmes individus :
poids et taille
 - Mesures successives d'un même phénomène :
pression artérielle le matin et le soir
 - Etude de la variabilité inter-opérateur :
mesure de pression artérielle par l'infirmière X et l'infirmière Y
- Questions :
 - Sont-elles indépendantes, en relation linéaire, ou +/- associées ?
 - Varient-elles dans le même sens ou non ?

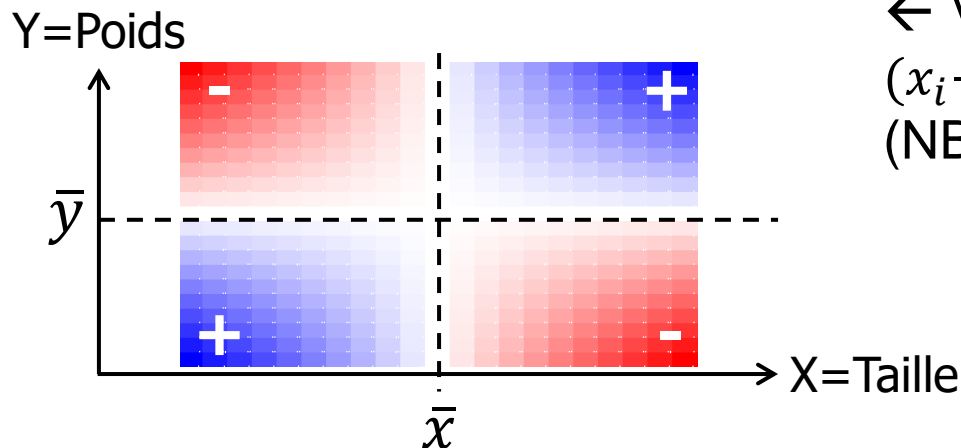


Rappel sur la Covariance

Exemples : poids et taille



1 point =
1 individu $i \{x_i, y_i\}$



← Valeurs du produit
 $(x_i - \bar{x}) \cdot (y_i - \bar{y})$ pour chaque individu i
(NB : nul sur les lignes des moyennes)

=> Dans ce cas, la moyenne de ce produit devrait être nettement positive.



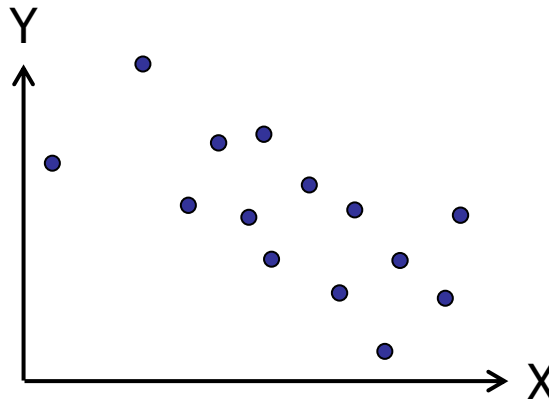
Rappel sur la Covariance

Autres exemples

Variables :

Y = montants des aides sociales perçues (€)
X = revenus du travail (€)

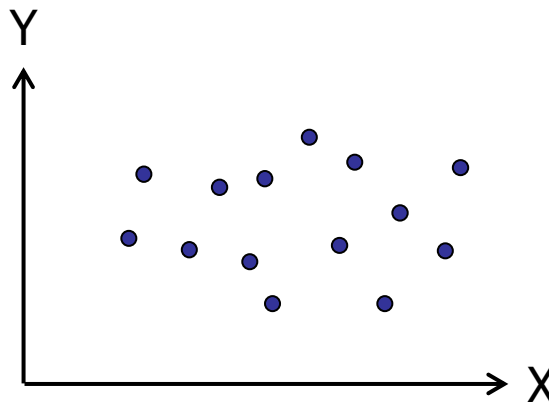
Graphique :



Valeur de
moyenne $[(x_i - \bar{x}) \cdot (y_i - \bar{y})]$

Fortement
négative

Y = distance travail (km)
X = taille (cm)



Nulle

Valeur de
moyenne $[(x_i - \bar{x}) \cdot (y_i - \bar{y})]$

Rappel sur la Covariance

Calcul dans un échantillon

- Covariance sur l'échantillon (empirique), avec une ligne pour chacun des n individus i :

$$\begin{aligned} cov_{ech}(X, Y) &= \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{n} \\ &= \left(\sum_{i=1}^n x_i \cdot y_i \right) / n - \bar{x} \cdot \bar{y} \end{aligned}$$

- Calcul analogue à l'estimateur biaisé de la variance :

$$cov_{ech}(X, X) = var_{ech}(X) = s^2_{ech}(X)$$



Rappel sur le Coefficient de corrélation linéaire de Pearson

Définition

- Coefficient insensible aux ordres de grandeur contrairement à la covariance
- Coefficient théorique de corrélation linéaire dans la population (rhô) : ρ
- Coefficient empirique de corrélation linéaire dans un échantillon :

$$r = \frac{\text{COV}_{ech}(X, Y)}{S_{echX} \cdot S_{echY}}$$

Attention : utiliser les estimateurs biaisés des écarts types

- Varie entre -1 et 1



Rappel sur le Coefficient de corrélation linéaire de Pearson

Interprétation

- En population :
 - Absence de relation linéaire entre X et Y $\Leftrightarrow \rho=0$
 - X et Y indépendantes $\Rightarrow \rho=0$
 - $Y=aX+b$ (relation strictement linéaire $a \neq 0$) $\Leftrightarrow \rho=1$ si $a>0$ ou -1 si $a<0$
 - X et Y corrélées linéairement
 - dans le même sens $\Leftrightarrow \rho \in]0 ; 1]$
 - en sens contraire $\Leftrightarrow \rho \in [-1 ; 0[$
 - fortement $\Leftrightarrow |\rho| \in [0,5 ; 1]$
 - faiblement $\Leftrightarrow |\rho| \in]0 ; 0,5[$

- Sur l'échantillon :
 - idem, en tenant compte des fluctuations liées à l'échantillonnage

- ATTENTION : $\rho=0$ signifie absence de relation linéaire, mais $\rho=0$ ne signifie pas toujours indépendance !

Exemples de nuages de points

avec $\rho=0$ mais X et Y non indépendantes :



Tester l'absence de relation linéaire entre deux variables quantitatives

- Dans un échantillon
 - Taille n
 - Observation de X et Y
 - Calcul de r , par exemple non-nul
- Extrapolation en population :
 - X et Y sont-elles en relation linéaire ?
 - Le coefficient ρ est-il non-nul ?
Test de nullité du coefficient de corrélation
 - Soit la relation $Y=aX+b+\varepsilon$, le coefficient a est-il non-nul ?
Test de nullité de la pente de la droite de régression
 - Calcul différent, mais résultat identique
 - ATTENTION : ne teste pas l'indépendance stricto sensu... => importance de l'exploration graphique !



Test de nullité du coeff. de corrélation

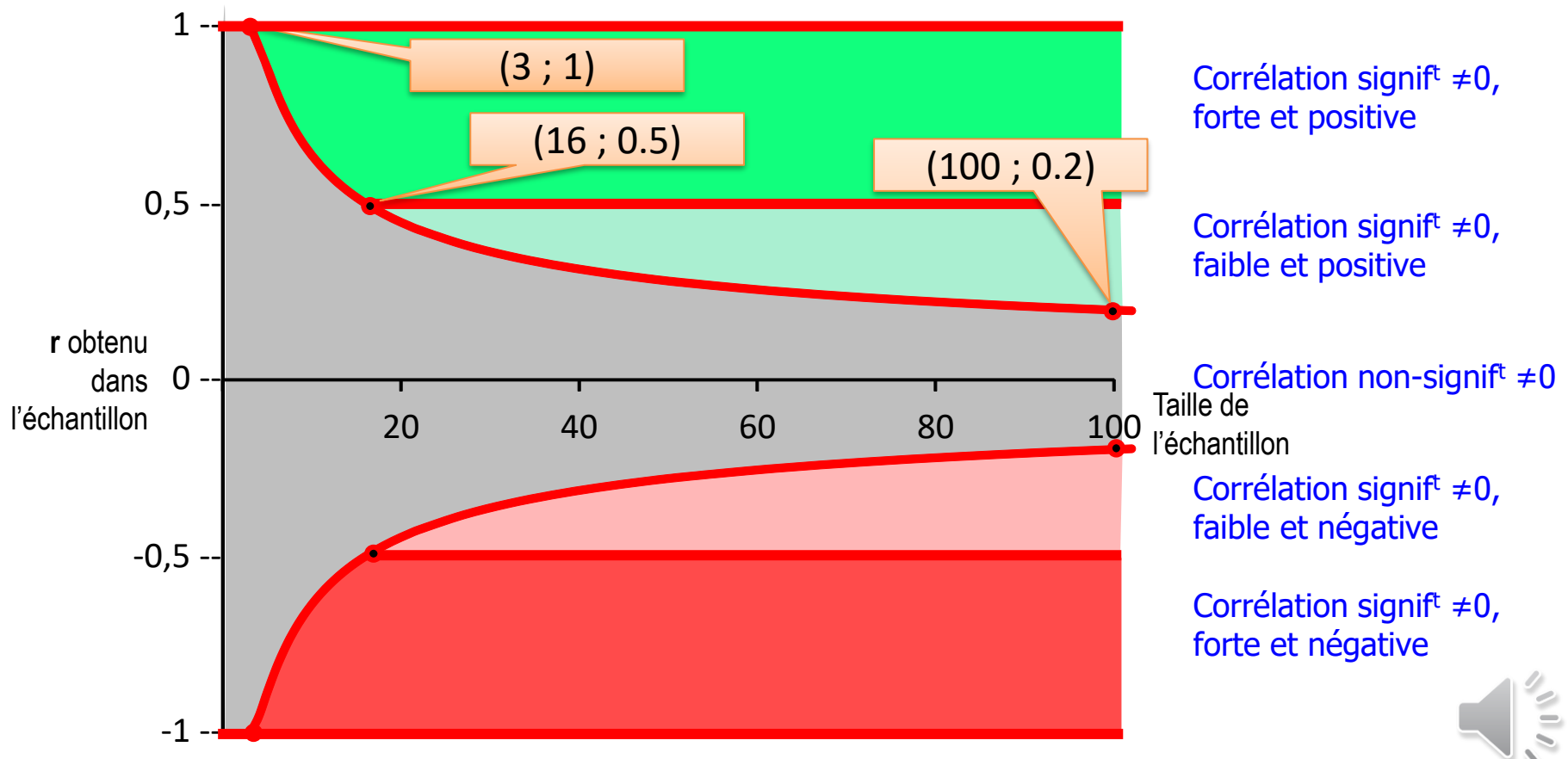
- Test bilatéral
- Hypothèse H_0 :
 - $\rho=0$
 - Autrement dit, X et Y ne sont pas corrélées linéairement
- Condition de validité :
 - (...) acquises si X et Y suivent une distribution normale
 - (ce sont celles de la régression linéaire simple, cf. post)
- Réalisation (autre approche existe) :
La quantité t suit une loi de Student à $v=n-2$ ddl
 - Si $|t| > t_{\alpha v}$ alors rejet de H_0 :
X et Y sont linéairement corrélées
 - Sinon, non-rejet de H_0 : on n'observe pas de corrélation linéaire significative entre les variables en population

$$t = r \cdot \sqrt{\frac{n-2}{1-r^2}}$$



Synthèse : interprétation de r

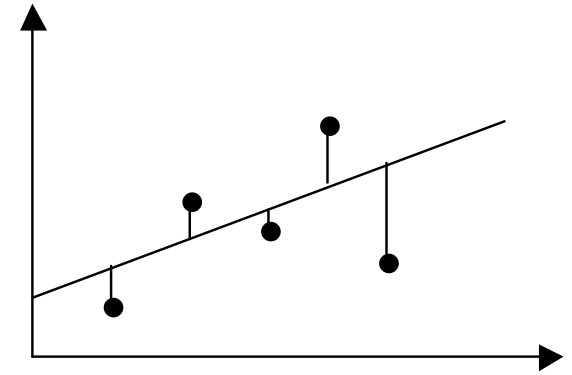
Interprétation en fonction de r (en Y) et de l'effectif n (en X)
(graphe obtenu d'après la table de Student, test bilatéral, $\alpha=5\%$):



Régression linéaire simple de Y sur X

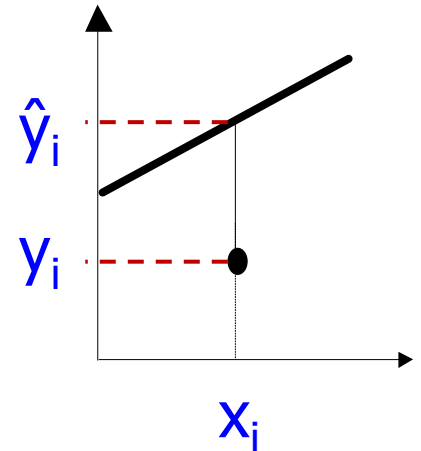
- On souhaite « modéliser » la relation entre X et Y par une équation $Y=aX+b+\varepsilon$ (ε : fluctuation aléatoire)
- Prédit Y et cherche à minimiser l'erreur « verticale » (c'est le plus fréquent)
- On calcule a et b sur l'échantillon :

$$a = r \cdot \frac{S_y}{S_x} = r \cdot \frac{S_{ech Y}}{S_{ech X}} = \frac{COV_{ech X,Y}}{S^2_{ech X}}$$
$$b = \bar{y} - a \cdot \bar{x}$$



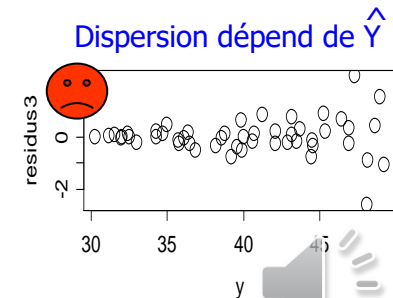
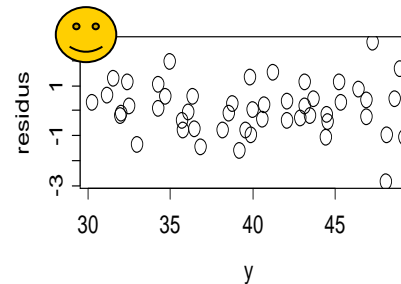
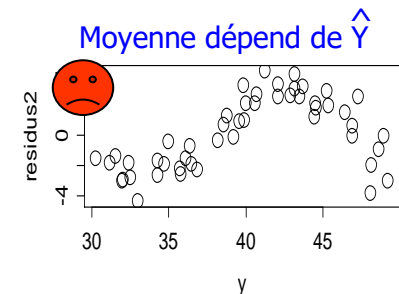
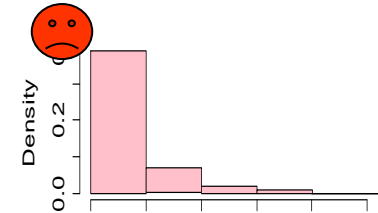
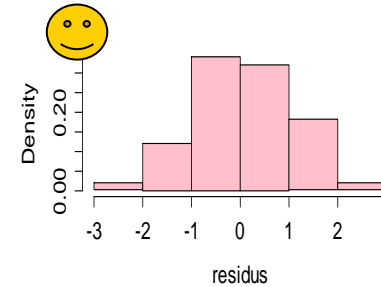
Régression de Y sur X

- On obtiendra ainsi sur l'échantillon :
 - a, b coefficients calculés
 - $y=ax+b$ équation de la droite de régression
 - $Y=aX+b+\varepsilon$ modèle de régression
- Et pour chaque individu i , avec x_i :
 - y_i valeur observée
 - $\hat{y}_i=a.x_i+b$ valeur prédite par la régression (sur la droite)
 - \hat{y}_i-y_i erreur (verticale) de prédiction ou « résidu »
- Condition de validité : à valider a posteriori
 - Analyse des résidus (relation linéaire)
 - Recherche d'individus influents



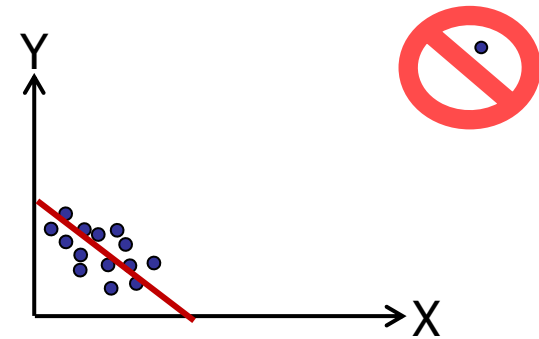
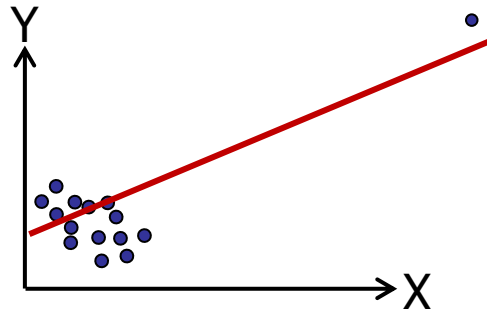
Analyse des résidus

- Rappel : une valeur par individu, $a \cdot x_i + b - y_i$
- Distribution univariée des résidus
 - Toujours : moyenne=0 et écart type déterminé (cf. post)
 - Objectif : distribution normale (d'emblée acquis si $Y \sim$ normale) \Rightarrow tracer un histogramme
- Résidus en fonction de \hat{Y}
 - Moyenne ne dépend pas de \hat{Y}
 - Variance ne dépend pas de \hat{Y} (homoscédasticité)



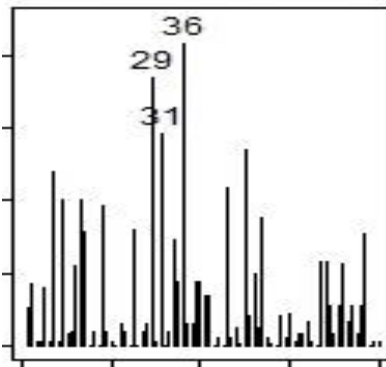
Recherche d'individus influents

■ Méthode graphique :

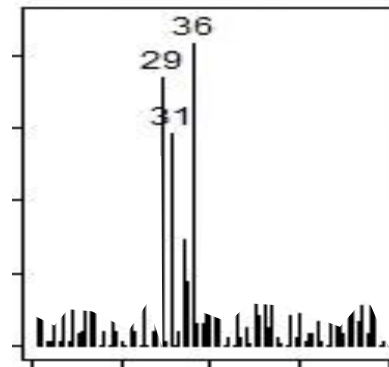


■ Distance de Cook (conseillée) :

Aucun individu ne se détache, OK



Deux ou trois individus influents...



Prédiction

- Pour une nouvelle valeur x_0 (avec y_0 inconnu), on peut tenter de prédire y_0 par \hat{y}_0
- Prédiction probabiliste : un individu n'aura pas exactement la valeur prédite, mais :
- L'erreur de prédiction suit une loi normale centrée et d'écart type s_{yx} :

$$s_{yx} = s_y \cdot \sqrt{1 - r^2}$$

$|r|$ élevée \Leftrightarrow prédiction précise

- IC de la prédiction à 95% :

$$IC_{y_0; 1-\alpha} = \hat{y}_0 \pm u_\alpha \cdot s_{yx}$$



Commentaire sur l'erreur de prédiction

$$s_{yx} = s_y \cdot \sqrt{1 - r^2}$$

- Si $r=0$

$$s_{yx} = s_y$$

=> aucune réduction d'incertitude sur la valeur de y

- Si $|r|=1$

$$s_{yx} = 0$$

=> y est sur la droite, aucune incertitude

- Cas général : l'existence d'une relation linéaire réduit l'incertitude sur y , sachant x .



Epilogue

- Test de nullité de la pente : équivalent au test de nullité du coefficient de corrélation de Pearson
- Test de nullité de l'ordonnée à l'origine : peu utilisé en pratique
- Possibilité de forcer un modèle $y=ax + 0$
- Régression de X sur Y :
 - Minimise l'erreur de prédiction de X : régression « horizontale »
 - Formules différentes
 - Rarement utilisée
- Régression orthogonale :
 - Minimise la distance (orthogonale) entre les points et la droite
 - Formules différentes
 - Rarement utilisée

