

Statistiques Bivariées

Croisement Quanti - Quanti

• Statistiques Bivariées = étude de deux caractères (variables) simultanément chez les individus.

But: Etudier s'il ya un lien entre ces deux variables.

Variables:

Ex: On prend deux variables

X = taux d'alcoolisme

Y = taux de suicides

et on étudie s'il ya un lien entre ces 2 variables

Type de croisement: ces deux variables peuvent être :

① quantitatif \times quantitatif

② qualitatif \times quantitatif

③ qualitatif \times qualitatif

Pour chaque type de croisement, on a une étude par le Quanti \times Quanti

1. Nuage des points:

- X et Y deux variables quantitatives sur un échantillon de taille (n); les objectifs:
 - ① Déterminer si il ya un lien (ou une corrélation) entre les deux variables.
 - ② Construire un modèle qui permet d'expliquer Y par X (ou l'inverse) si il ya un lien.

2. Droite de Régression:

- * Soit $\{x_i\}_{i=1, \dots, n}$ les observations de X et $\{y_i\}_{i=1, \dots, n}$ les observations de Y.
- Objectif: trouver une fonction f tel que:
 $y_i = f(x_i) + \varepsilon_i$ / $\varepsilon = \text{erreur}$.
 On se restreint aux fonctions affines: $f(x) = ax + b$
- On cherche les coefficients a et b qui minimisent l'erreur quadratique moyenne.

$$\begin{aligned} EQ(a, b) &= \frac{1}{n} \sum_i \varepsilon_i^2 \\ &= \frac{1}{n} \sum_i (y_i - (ax_i + b))^2 \end{aligned}$$

②

* Droite de Régression de Y en X est:

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

$$\hat{b} = \bar{y} - \hat{a}\bar{x}$$

tel que: $\hat{a} = \frac{C_{xy}}{S_x^2}$

. $C_{xy} = \frac{1}{n} \sum_i (x_i - \bar{x})(y_i - \bar{y})$

c'est la covariante entre X et Y.

. $S_x^2 = \text{variance de } X$

Ce droite traduit les variations de Y qui peuvent être expliquées par X.

* Droite de Régression de X en Y.

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

$$\hat{b} = \bar{x} - \hat{a}\bar{y}$$

tel que: $\hat{a} = \frac{C_{xy}}{S_y^2}$

Exp: lien \tilde{x} et \tilde{y} chez les enfants de 6 ans.

. Équation de Y en X;

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

$$\hat{a} = \frac{C_{xy}}{S_x^2} = \frac{16,27}{38,52} = 0,42$$

$$\hat{b} = \bar{y} - \hat{a}\bar{x} = (20,30) - (0,42)(113,20) = 27,38$$

donc

$$y = 0,42x - 27,38$$

Exp: Si taille d'un enfant est $x = 115$, estimer le poids ??

$$y = 0,42x - 27,38 = 0,42(115) - 27,38$$

$$\textcircled{3} = 20,92 \text{ Kg.}$$

Équation de la droite X en Y :

$$X = \hat{a}Y + \hat{b}$$

$$\hat{a} = \frac{C_{XY}}{S_y^2} = \frac{16,27}{8,46} = 1,92$$

$$\hat{b} = \bar{x} - \hat{a}\bar{y} = 113,2 - (1,92)(20,30) \\ = +74,15$$

donc

$$X = 1,92Y + 74,15$$

taille

poids

Ex: Si le poids d'un enfant est 23 kg, estimer son taille:

$$X = 1,92Y + 74,15$$

$$= 1,92(23) + 74,15 = 118,31 \text{ cm}$$

* Covariance: C'est un indicateur numérique du lien entre X et Y : plus il est loin de 0, plus les variables sont liées. [$C_{XY} = 0$, pas de lien entre les 2 variables alors elles sont indépendantes].

⚠ La covariance n'est pas normée, donc il prendra n'importe quel valeur sur R. Pour cela pour étudier le lien entre X et Y , on définit le coefficient de Corrélation.

Coefficient de Corrélation (Coefficient de Pearson) :
pour étudier le lien entre deux variables quantitatives
 X et Y : prend des valeurs entre $[-1, 1]$

$$r_{xy} = \frac{C_{xy}}{S_x \cdot S_y}$$

on a alors :

$$\hat{a} = \frac{C_{xy}}{S_x^2} \Rightarrow \hat{a} = \frac{r_{xy} \cdot S_y}{S_x}$$

Interpretation: ① $|r|$ proche de 1, alors X et Y sont très liés entre eux par une droite affine

② $r < 0$: X et Y varie en sens inverse

Ex: $\begin{cases} X = \text{dépense} \\ Y = \text{Epargne} \end{cases}$

Si X augmente donc Y diminue.

③ $r > 0$: X et Y varie de même sens.

Ex: $\begin{cases} X = \text{taux nicotine dans sang} \\ Y = \text{risque avoir maladie poumon} \end{cases}$

Si X augmente alors Y aussi

④ $|r| = 0$: On ne peut rien dire sur un lien éventuel entre X et Y .

Exemple: pour poids et âge chez les enfants de 6 mois:

$$r_{xy} = \frac{C_{xy}}{S_x \cdot S_y} = \frac{16,27}{\sqrt{38,6^2} \times \sqrt{8,46}} \approx 0,9$$

- $r_{xy} \approx 1 \rightarrow$ l'éq de droite est pleinement justifiée.
- $r_{xy} > 0 \rightarrow$ plus la variable taille est grande
- $r_{xy} > 0 \rightarrow$ plus la variable poids est important et vice versa.

⑤

Prévisions: on appelle (prévisions), les valeurs données par la droite de régression.

Pour chaque x_i on peut calculer \hat{y}_i :

\hat{y}_i = c'est la valeur approchée (estimée) de la vraie valeur y_i .

$$\hat{y} = \bar{y}$$

mais

pas en variance.

$$S_{\hat{y}}^2 = r^2_{xy} \cdot S_y^2$$

Les \hat{y}_i éoliné
ont en moyenne
que les vraies valeurs
que y_i .

c'est l'écart entre y_i et \hat{y}_i :

$$e_i = y_i - \hat{y}_i$$

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

Erreurs globales: $EQ(\hat{a}, \hat{b}) = \frac{1}{n} \sum_i e_i^2 = S_y^2 (1 - r^2_{xy})$

A) Variance $\gamma \neq$ variance $\hat{\gamma}$.

$$S_y^2 = S_{\hat{y}}^2 + S_e^2$$

Variance totale Variance expliquée Variance des résidus.

B) Coefficient de détermination:

$$R^2 = \frac{S_{\hat{y}}^2}{S_y^2} = \frac{\sum xy}{\sum x^2} \quad (\text{entre } 0 \text{ et } 1).$$

Donne proportion des observations expliquées par la droite de régression (plus il est élevé plus c'est mieux)

Ex: On donne:

$$\begin{cases} S_{\hat{y}}^2 = 6,17 \\ S_e^2 = 1,44 \\ S_y^2 = 7,61 \end{cases}$$

$$\text{donc } R^2 = \frac{6,17}{7,61} = 0,81$$

$\Rightarrow 81\%$ de la variation des poids observés est expliquée par la droite de régression.

$$\underline{\underline{S_y^2}} = \underline{\underline{S_{\hat{y}}^2}} + \underline{\underline{S_e^2}}$$

Variance totale
 Variance explicative
 Variance résiduelle

- $S_{\hat{y}}^2 = r^2_{xy} \cdot S_y^2$ | r_{xy} = coeff de corrélat.
 - $S_e^2 = S_y^2 (1 - r^2_{xy})$
 - $S_y^2 = \frac{1}{n} \sum_{k=1}^n (y_k - \bar{y})^2$
- $r_{xy} = \frac{C_{xy}}{S_x \cdot S_y}$.

⚠️ Outliers et (valeur aberrantes)

Comment faire ??



+ Si on a une (plusieurs) valeur aberrante(s) :

- ① On retire la $i^{\text{ème}}$ observation
- ② On refait un nouveau modèle sans le $i^{\text{ème}}$
- ③ On calcule $\hat{y}_{(-i)}$ (la prévision) avec le nouveau modèle
- ④ On calcule le résidu :

$$e_{(-i)} = y_i - \hat{y}_{(-i)}$$

pour le \checkmark
modèle sans
 (i) .

⚠️ un résidu à valeur importante signale une observation abnégante.

On:

$$e_{(i)} = \frac{e_i}{1-h_{ii}}$$

calcul avec
résidu sans (i)

calcul avec
modèle initial
avec $i^{\text{ème}}$ don.

ou

$$h_{ii} = \frac{1}{n} + \frac{1}{n-1} \frac{(x_i - \bar{x})^2}{s_x^2}$$

h_{ii} = Levier, prend ses valeurs entre: $\frac{1}{n} \leq h_{ii} \leq 1$
Si h_{ii} proche de 1 alors on a une observation influente

(9)

- Les résidus standardisés doivent être toujours entre $[-2, +2]$ [I.D.C à voir en ING2].

[ce c'est l'hypothèse de normalité des résidus)

$$S_i = \frac{e_i}{\hat{\sigma} \cdot \sqrt{1-h_{ii}}}$$

Résidu Standardisé