



---

**TP Analyse bivariable**  
**Croisement Qualitatif-Quantitatif**  
*Corrigé*

---

**Exercice 1**

*Décomposition de la variance*

Dans une population  $\Omega$  de taille  $n$ , on observe deux variables :

- une qualitative,  $x = \{x_k\}_{k=1, \dots, n}$ , à  $p$  modalités notées,  $m_1, \dots, m_p$
- une quantitative continue  $y = \{y_k\}_{k=1, \dots, n}$  de moyenne  $\bar{y}$  et de variances  $s_y^2$ .

On suppose que les modalités de la série  $x$  définissent des sous-populations

$$\Omega = \Omega_1 \cup \dots \cup \Omega_p \text{ où } \Omega_i \cap \Omega_j = \emptyset,$$

de tailles respectives  $n_1, \dots, n_p$ .

On peut alors considérer les restrictions de la caractéristique  $y$  sur chacune des sous-populations et calculer les indicateurs numériques usuels pour chaque modalités de  $x$ ,

- moyennes :  $\bar{y}_i, i=1, \dots, p$
- variances :  $s_i^2, i=1, \dots, p$

Montrer que

$$\bar{y} = \frac{1}{n} \sum_{i=1}^p n_i \bar{y}_i$$
$$s_y^2 = \frac{1}{n} \sum_{i=1}^p n_i (\bar{y}_i - \bar{y})^2 + \frac{1}{n} \sum_{i=1}^p n_i s_i^2 = s_E^2 + s_R^2$$

A quoi correspondent les termes  $s_E^2$  et  $s_R^2$  ?

On définit un indice de liaison entre les deux caractéristiques  $x$  et  $y$  par le rapport de corrélation

$$s_{y/x} = \sqrt{\frac{s_E^2}{s_y^2}}$$

Donner un encadrement de  $s_{y/x}$ . A quoi correspondent les cas  $s_{y/x}=0$  et  $s_{y/x}=1$  ?

**Exercice 2**

*Données : SalairesData.xls*

Le fichier présente les salariés d'une entreprise ayant 3 sites (A, B et C). On y indique leur sexe, leur salaire annuel, leur catégorie (CS : cadre supérieur, CM : cadre moyen, OE : ouvrier employé), leur âge et leur site.

- 1) Construire un tableau répertoriant le salaire des individus par catégorie, puis par site.
- 2) Sur un même graphique, représenter les boîtes de Tuckey pour chaque catégorie. Commenter
- 3) Faire le même graphique mais par site. Sur quel site vaut-il mieux travailler à votre avis ?
- 4) Pouvez-vous justifier votre réponse à l'aide d'un indicateur numérique ?

### Exercice 3

*Données : EnsSuperieurData.xls*

Le fichier EnsSuperieur.xls comptabilise le nombre d'étudiants par sexe dans l'enseignement supérieur de premier et deuxième cycles. Il s'agit chiffres relevés par Eurostat en 2008.

Illustrer et commenter ces chiffres en travaillant dans un premier temps sur le nombre d'étudiants et ensuite sur le taux d'étudiants pour 1000 habitants.

Exercice 1

Notons  $y_{k,l}$ ,  $k=1, \dots, n_l$ , la  $k^{\text{ème}}$  observation de la caractéristique Y appartenant à la sous-population  $\Omega_l$ .

- Décomposition de la moyenne

$$\bar{y}_1 = \frac{1}{n_1} \sum_{k=1}^{n_1} y_{k,1} \Rightarrow n_1 \bar{y}_1 = \sum_{k=1}^{n_1} y_{k,1} \Rightarrow \frac{1}{n} \sum_{l=1}^p n_l \bar{y}_l = \frac{1}{n} \sum_{l=1}^p \sum_{k=1}^{n_l} y_{k,l} = \frac{1}{n} \sum_{k=1}^n y_k = \bar{y}$$

- Décomposition de la variance (Calculs très longs ne pas laisser chercher)

$$\begin{aligned} s_E^2 &= \frac{1}{n} \sum_{l=1}^p n_l (\bar{y}_l - \bar{y})^2 = \frac{1}{n} \sum_{l=1}^p n_l (\bar{y}_l^2 - 2\bar{y}_l \bar{y} + \bar{y}^2) = \frac{1}{n} \sum_{l=1}^p n_l \bar{y}_l^2 - \frac{2}{n} \bar{y} \sum_{l=1}^p n_l \bar{y}_l + \frac{1}{n} \bar{y}^2 \sum_{l=1}^p n_l \\ &= \frac{1}{n} \sum_{l=1}^p n_l \bar{y}_l^2 - 2\bar{y}^2 + \bar{y}^2 = \frac{1}{n} \sum_{l=1}^p n_l \bar{y}_l^2 - \bar{y}^2 \end{aligned}$$

Nous avons

$$s_1^2 = \frac{1}{n_1} \sum_{k=1}^{n_1} (y_{k,1} - \bar{y}_1)^2,$$

d'où

$$\begin{aligned} s_R^2 &= \frac{1}{n} \sum_{l=1}^p n_l s_l^2 = \frac{1}{n} \sum_{l=1}^p n_l \frac{1}{n_l} \sum_{k=1}^{n_l} (y_{k,l} - \bar{y}_l)^2 = \frac{1}{n} \sum_{l=1}^p \sum_{k=1}^{n_l} (y_{k,l}^2 - 2y_{k,l} \bar{y}_l + \bar{y}_l^2) \\ &= \frac{1}{n} \sum_{l=1}^p \sum_{k=1}^{n_l} y_{k,l}^2 - \frac{2}{n} \sum_{l=1}^p \bar{y}_l \sum_{k=1}^{n_l} y_{k,l} + \frac{1}{n} \sum_{l=1}^p \bar{y}_l^2 \sum_{k=1}^{n_l} 1 \\ &= \frac{1}{n} \sum_{l=1}^p \sum_{k=1}^{n_l} y_{k,l}^2 - \frac{2}{n} \sum_{l=1}^p \bar{y}_l n_l \bar{y}_l + \frac{1}{n} \sum_{l=1}^p \bar{y}_l^2 n_l = \frac{1}{n} \sum_{l=1}^p \sum_{k=1}^{n_l} y_{k,l}^2 - \frac{1}{n} \sum_{l=1}^p n_l \bar{y}_l^2 \end{aligned}$$

Donc

$$\begin{aligned} s_E^2 + s_R^2 &= \frac{1}{n} \sum_{l=1}^p n_l \bar{y}_l^2 - \bar{y}^2 + \frac{1}{n} \sum_{l=1}^p \sum_{k=1}^{n_l} y_{k,l}^2 - \frac{1}{n} \sum_{l=1}^p n_l \bar{y}_l^2 \\ &= \frac{1}{n} \sum_{l=1}^p \sum_{k=1}^{n_l} y_{k,l}^2 - \bar{y}^2 = \frac{1}{n} \sum_{k=1}^n y_k^2 - \bar{y}^2 = s_y^2 \end{aligned}$$

Le premier terme correspond à la *variance expliquée* par la partition de la série observée X et le deuxième terme est un reste appelé *variance résiduelle*.

- Rapport de corrélation

$$s_y^2 = s_E^2 + s_R^2 \Rightarrow 0 \leq s_E^2 \leq s_y^2 \Rightarrow 0 \leq \frac{s_E^2}{s_y^2} \leq 1 \Rightarrow 0 \leq s_{y/x} \leq 1$$

Cela signifie que  $s_E^2$  représente le pourcentage de variabilité de y expliquée par x.

- Si  $s_{y/x}=0$  alors la variance expliquée est nulle, il n'y a donc aucun lien entre y et x
- Si  $s_{y/x}=1$  alors la variance expliquée est égale à la variance de y donc y est entièrement expliquée par x.

Attention :  $s_{y/x} \neq s_{x/y}$ . Cela signifie que si  $s_{y/x}$  est proche de 1 alors x permet d'expliquer y mais pas réciproquement.







