

# Chapitre 3. Corrélation et ajustement affine

## R2.13 – Techniques quantitatives et représentations 2

D. Fourer



Département TC-Juvisy

Semestre 2 – 2023/2024

# Contexte et notations

## VARIABLES STATISTIQUES

- Variables  $X, Y$  continues **non regroupées en classes** sur une **même population** (taille  $N$ ).
- ▶ Définitions **brutes** des quantités à partir des individus.

# Contexte et notations

## VARIABLES STATISTIQUES

- Variables  $X, Y$  continues **non regroupées en classes** sur une **même population** (taille  $N$ ).
- ▶ Définitions **brutes** des quantités à partir des individus.

## PARAMÈTRES DE POSITION ET DE DISPERSION

- Moyennes **marginales** :  $\bar{x}$  et  $\bar{y}$ .
- Variances et écarts-types **marginiaux** :  $\text{var } X, \text{var } Y$  et  $\sigma_X, \sigma_Y$ .

1 Notion de corrélation

2 Droite de régression

1 Notion de corrélation

2 Droite de régression

# Covariance – Construction

## IDÉES DIRECTRICES

- Obtenir un **indicateur numérique** (plutôt que le « oui/non » de l'indépendance statistique).
- Obtenir une notion **plus facilement réalisée** que l'indépendance statistique.

# Covariance – Construction

## IDÉES DIRECTRICES

- Obtenir un **indicateur numérique** (plutôt que le « oui/non » de l'indépendance statistique).
- Obtenir une notion **plus facilement réalisée** que l'indépendance statistique.

### Définition

**Covariance**  $\text{cov}(X, Y)$  : Moyenne statistique des produits écarts des valeurs à leur moyenne marginale respective.

# Covariance – Construction

## IDÉES DIRECTRICES

- Obtenir un **indicateur numérique** (plutôt que le « oui/non » de l'indépendance statistique).
- Obtenir une notion **plus facilement réalisée** que l'indépendance statistique.

## Définition

**Covariance**  $\text{cov}(X, Y)$  : Moyenne statistique des produits écarts des valeurs à leur moyenne marginale respective.

## FORMULES

$$\frac{1}{N} \sum_{i=1}^N (X(i) - \bar{x})(Y(i) - \bar{y}), \quad \sum_{i,j} f_{ij}(x_i - \bar{x})(y_j - \bar{y}), \quad \dots$$

# Covariance – Propriétés

## LIEN AVEC LA VARIANCE

Extension et **généralisation** de la notion de variance :

$$\text{cov}(X, X) = \text{var } X.$$

# Covariance – Propriétés

## LIEN AVEC LA VARIANCE

Extension et **généralisation** de la notion de variance :

$$\text{cov}(X, X) = \text{var } X.$$

## CALCUL EFFECTIF

- Différence entre la moyenne statistique des produits des valeurs et le produit des moyennes marginales.
- Formules :

$$\left( \frac{1}{N} \sum_{i=1}^N X(i) Y(i) \right) - \bar{X}\bar{Y}, \quad \left( \sum_{i,j} f_{ij} X_i Y_j \right) - \bar{X}\bar{Y}, \quad \dots$$

# Corrélation et lien avec l'indépendance statistique

## Définition

On dit que  $X$  et  $Y$  sont **corrélées** lorsque  $\text{cov}(X, Y) \neq 0$  et **non corrélées** lorsque  $\text{cov}(X, Y) = 0$ .

# Corrélation et lien avec l'indépendance statistique

## Définition

On dit que  $X$  et  $Y$  sont **corrélées** lorsque  $\text{cov}(X, Y) \neq 0$  et **non corrélées** lorsque  $\text{cov}(X, Y) = 0$ .

## Proposition

*Si  $X$  et  $Y$  sont indépendantes, alors elles sont non corrélées.  
La réciproque est fausse.*

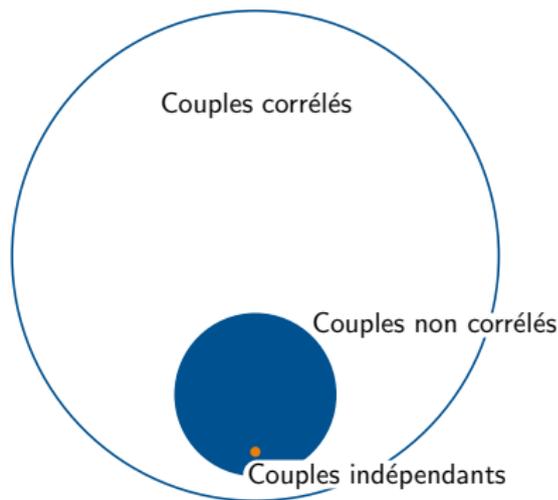
# Corrélation et lien avec l'indépendance statistique

## Définition

On dit que  $X$  et  $Y$  sont **corrélées** lorsque  $\text{cov}(X, Y) \neq 0$  et **non corrélées** lorsque  $\text{cov}(X, Y) = 0$ .

## Proposition

*Si  $X$  et  $Y$  sont indépendantes, alors elles sont non corrélées.*  
*La réciproque est fausse.*



## ATTENTION !

Les couples non corrélés sont « **rares** ». (Cf. site Spurious correlations.)

# Coefficient de corrélation

## Définition

Coefficient de corrélation (ou de régression) de  $X$  et  $Y$  :

$$r(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}.$$

# Coefficient de corrélation

## Définition

**Coefficient de corrélation** (ou de régression) de  $X$  et  $Y$  :

$$r(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}.$$

- Il discrimine la non corrélation :  $r(X, Y) \neq 0 \iff \text{cov}(X, Y) \neq 0$ .
- Il est **borné** :  $-1 \leq r(X, Y) \leq 1$ . (Inégalité de Cauchy-Schwarz.)

# Coefficient de corrélation

## Définition

**Coefficient de corrélation** (ou de régression) de  $X$  et  $Y$  :

$$r(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}.$$

- Il discrimine la non corrélation :  $r(X, Y) \neq 0 \iff \text{cov}(X, Y) \neq 0$ .
- Il est **borné** :  $-1 \leq r(X, Y) \leq 1$ . (Inégalité de Cauchy-Schwarz.)

## VALEURS REMARQUABLES

- $r(X, Y) = 1$  : points du nuage alignés sur une **droite croissante**.
- $r(X, Y) = 0$  :  $X$  et  $Y$  non corrélées, nuage « **loin** » de toute droite.
- $r(X, Y) = -1$  : points du nuage alignés sur une **droite décroissante**.

# Qu'est-ce que l'ajustement affine ?

## Définition

**Ajustement affine** : Méthode de construction d'une droite modélisant (au mieux) le nuage de points.

# Qu'est-ce que l'ajustement affine ?

## Définition

**Ajustement affine** : Méthode de construction d'une droite modélisant (au mieux) le nuage de points.

## POINTS ESSENTIELS

- Comment évaluer la **fiabilité** de la modélisation ?
- Existe-t-il une droite **meilleure** que les autres ?
  - ▶ *Le cas échéant, en quel sens est-elle meilleure ?*

# Qu'est-ce que l'ajustement affine ?

## Définition

**Ajustement affine** : Méthode de construction d'une droite modélisant (au mieux) le nuage de points.

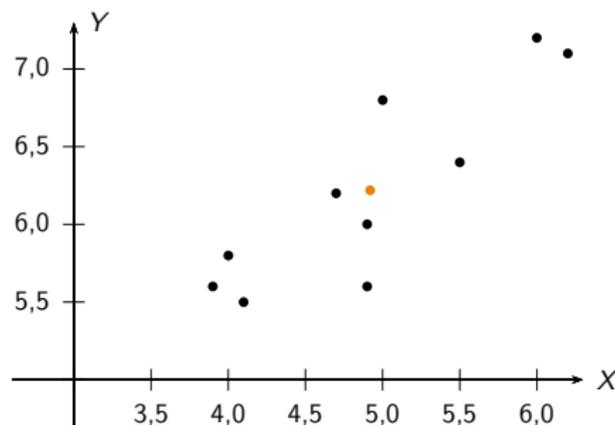
## POINTS ESSENTIELS

- Comment évaluer la **fiabilité** de la modélisation ?
- Existe-t-il une droite **meilleure** que les autres ?
  - ▶ *Le cas échéant, en quel sens est-elle meilleure ?*

## Définition (utile)

**Point moyen**  $G$  : Point de coordonnées  $(\bar{x}, \bar{y})$  dans le nuage.

# Exemple fondamental (TD 3 – Exercice 1)



COORDONNÉES DU POINT MOYEN  $G$

$$\bar{x} = 4,92 \quad \text{et} \quad \bar{y} = 6,22.$$

1 Notion de corrélation

2 Droite de régression

# Construction

## PRINCIPE

- Variable statistique des **écarts verticaux** :  $E_{a,b} = Y - (aX + b)$ .
  - ▶ Choisir  $a$  et  $b$  pour que  $E_{a,b}$  soit **proche de la variable nulle**.
- Minimisation de la **variance** de  $E_{a,b}$  (moindres carrés) :

$$\text{var } E_{a,b} = a^2 \text{var } X - 2a \text{cov}(X, Y) + \text{var } Y.$$

- Annulation de la **moyenne** de  $E_{a,b}$  :  $\bar{e}_{a,b} = \bar{y} - (a\bar{x} + b)$ .

# Construction

## PRINCIPE

- Variable statistique des **écarts verticaux** :  $E_{a,b} = Y - (aX + b)$ .
  - ▶ Choisir  $a$  et  $b$  pour que  $E_{a,b}$  soit **proche de la variable nulle**.
- Minimisation de la **variance** de  $E_{a,b}$  (moindres carrés) :

$$\text{var } E_{a,b} = a^2 \text{var } X - 2acov(X, Y) + \text{var } Y.$$

- Annulation de la **moyenne** de  $E_{a,b}$  :  $\bar{e}_{a,b} = \bar{y} - (a\bar{x} + b)$ .

## Définition

**Droite de régression** de  $Y$  en  $X$  : équation réduite  $y = \tilde{a}x + \tilde{b}$  avec :

$$\tilde{a} = \frac{\text{cov}(X, Y)}{\text{var } X} = \frac{\text{cov}(X, Y)}{\sigma_X^2} = r(X, Y) \frac{\sigma_Y}{\sigma_X} \quad \text{et} \quad \tilde{b} = \bar{y} - \tilde{a}\bar{x}.$$

## Exemple (TD 3 – Exercice 1)

## PARAMÈTRES STATISTIQUES

- Point moyen  $G(4,92; 6,22)$ .
- Variances et écart-types :

$$\text{var } X = 0,5756 \quad \text{et} \quad \sigma_X = 0,7587 \text{ à } 10^{-4} \text{ près,}$$

$$\text{var } Y = 0,3616 \quad \text{et} \quad \sigma_Y = 0,6013 \text{ à } 10^{-4} \text{ près.}$$

- Covariance :  $\text{cov}(X, Y) = 0,3966$ .

## Exemple (TD 3 – Exercice 1)

## PARAMÈTRES STATISTIQUES

- Point moyen  $G(4,92; 6,22)$ .
- Variances et écart-types :

$$\text{var } X = 0,5756 \quad \text{et} \quad \sigma_X = 0,7587 \text{ à } 10^{-4} \text{ près,}$$

$$\text{var } Y = 0,3616 \quad \text{et} \quad \sigma_Y = 0,6013 \text{ à } 10^{-4} \text{ près.}$$

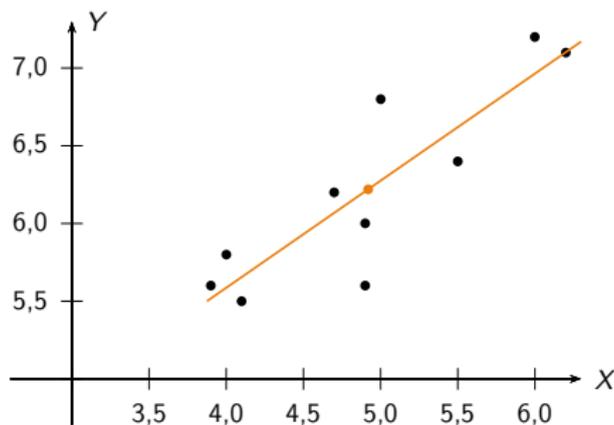
- Covariance :  $\text{cov}(X, Y) = 0,3966$ .

## DROITE DE RÉGRESSION

- Paramètres :  $\tilde{a} = 0,6890$  et  $\tilde{b} = 2,8301$  à  $10^{-4}$  près.
- Coefficient de corrélation :  $r(X, Y) = 0,8693$  à  $10^{-4}$  près.

## Exemple (TD 3 – Exercice 1)

2/2



## AVANTAGES

- **Meilleure** droite possible (au sens des écarts verticaux).
- Les projections sont **fiables** lorsque :

$$|r(X, Y)| \geq 0,9 \iff (r(X, Y) \geq 0,9 \text{ ou } r(X, Y) \leq -0,9).$$

## FIN DU CHAPITRE 3

---