

# **MOOC Statistique pour ingénieur**

## **Thème 0 : statistique descriptive**

### **Vidéo 2 : Statistiques à deux variables**

F. Delacroix    M. Lecomte

Institut Mines-Télécom  
École Nationale Supérieure des Mines de Douai

# Sommaire

- 1 Distributions à deux caractères
- 2 Covariance
- 3 Coefficient de corrélation linéaire

# Un exemple

- Test en compression d'éprouvettes de béton
- $X$ =teneur en ciment ( $kg/m^3$ )
- $Y$ =résistance à la compression (MPa)



# Collecte des données

$n = 90$  mesures

$X \backslash Y$	60	80	100
300	15	4	1
350	10	20	10
400	5	10	15

$X \backslash Y$	$y_1$	$\dots$	$y_j$	$\dots$	$y_c$	Total
$x_1$			$n_{1j}$			
$\vdots$			$\vdots$			
$x_i$	$n_{i1}$	$\dots$	$n_{ij}$	$\dots$	$n_{ic}$	$n_{i.}$
$\vdots$			$\vdots$			
$x_r$			$n_{rj}$			
Total			$n_{.j}$			$n$

$n_{ij}$  = nombre d'observations avec  $X = x_i$  et  $Y = y_j$

# Collecte des données

$n = 90$  mesures

$X \backslash Y$	60	80	100
300	15	4	1
350	10	20	10
400	5	10	15

$X \backslash Y$	$y_1$	$\dots$	$y_j$	$\dots$	$y_c$	Total
$x_1$			$n_{1j}$			
$\vdots$			$\vdots$			
$x_i$	$n_{i1}$	$\dots$	$n_{ij}$	$\dots$	$n_{ic}$	$n_{i.}$
$\vdots$			$\vdots$			
$x_r$			$n_{rj}$			
Total			$n_{.j}$			$n$

$$n_{i.} = \text{effectif marginal de la } i^{\text{ème}} \text{ ligne} = \sum_{j=1}^c n_{ij}$$

# Collecte des données

$n = 90$  mesures

$X \backslash Y$	60	80	100
300	15	4	1
350	10	20	10
400	5	10	15

$X \backslash Y$	$y_1$	$\dots$	$y_j$	$\dots$	$y_c$	Total
$x_1$			$n_{1j}$			
$\vdots$			$\vdots$			
$x_i$	$n_{i1}$	$\dots$	$n_{ij}$	$\dots$	$n_{ic}$	$n_{i.}$
$\vdots$			$\vdots$			
$x_r$			$n_{rj}$			
Total			$n_{.j}$			$n$

$$n_{.j} = \text{effectif marginal de la } j^{\text{ème}} \text{ colonne} = \sum_{i=1}^r n_{ij}$$

# Collecte des données

$n = 90$  mesures

$X \backslash Y$	60	80	100
300	15	4	1
350	10	20	10
400	5	10	15

$X \backslash Y$	$y_1$	$\dots$	$y_j$	$\dots$	$y_c$	Total
$x_1$			$n_{1j}$			
$\vdots$			$\vdots$			
$x_i$	$n_{i1}$	$\dots$	$n_{ij}$	$\dots$	$n_{ic}$	$n_{i.}$
$\vdots$			$\vdots$			
$x_r$			$n_{rj}$			
Total			$n_{.j}$			$n$

$$n = \sum_{j=1}^c \sum_{i=1}^r n_{ij} = \sum_{i=1}^r n_{i.} = \sum_{j=1}^c n_{.j}$$

# Distribution conjointe

Fréquence de la cellule  $C_{ij}$  :  $f_{ij} = \frac{n_{ij}}{n}$

$X \backslash Y$	60	80	100	Total
300	15	4	1	$n_{1.}$ = 20
350	10	20	10	$n_{2.}$ = 40
400	5	10	15	$n_{3.}$ = 30
Total	$n_{.1}$ = 30	$n_{.2}$ = 34	$n_{.3}$ = 26	90

$X \backslash Y$	60	80	100	$f_{i.}$
300	16,7%	4,4%	1,1%	22,2%
350	11,1%	22,2%	11,1%	44,4%
400	5,6%	11,1%	16,7%	33,3%

$f_{i.} = \frac{n_{i.}}{n}$  distribution marginale en  $X$



# Distribution conjointe

Fréquence de la cellule  $C_{ij}$  :  $f_{ij} = \frac{n_{ij}}{n}$

$X \backslash Y$	60	80	100	Total
300	15	4	1	$n_{1.} = 20$
350	10	20	10	$n_{2.} = 40$
400	5	10	15	$n_{3.} = 30$
Total	$n_{.1} = 30$	$n_{.2} = 34$	$n_{.3} = 26$	90

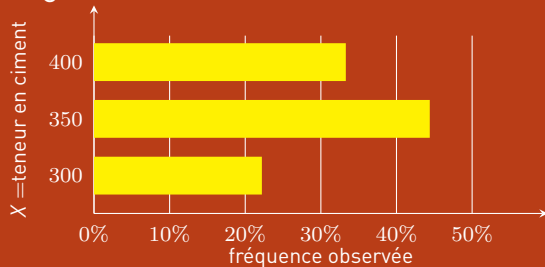
$X \backslash Y$	60	80	100	$f_{i.}$
300	16,7%	4,4%	1,1%	22,2%
350	11,1%	22,2%	11,1%	44,4%
400	5,6%	11,1%	16,7%	33,3%
$f_{.j}$	33,3%	37,8%	28,9%	

$f_{i.} = \frac{n_{i.}}{n}$  distribution marginale en  $X$

$f_{.j} = \frac{n_{.j}}{n}$  distribution marginale en  $Y$

# Distributions marginales

- diagramme en barres



- teneur moyenne en ciment des éprouvettes

$$\bar{x} = \frac{1}{n} \sum_{i=1}^r n_i \cdot x_i = \sum_{i=1}^r f_i \cdot x_i \simeq 355,5 \text{ kg/m}^3$$

# Distribution conditionnelles

$X \backslash Y$	60	80	100	Total
300	15	4	1	20
350	10	20	10	40
400	5	10	15	30
Total	30	34	26	90

$X \backslash Y$	$y_1$	$\dots$	$y_j$	$\dots$	$y_c$	Total
$y_j$ sachant que $X = x_i$	$n_{i1}$	$\dots$	$n_{ij}$	$\dots$	$n_{ic}$	$n_{i.}$

# Distribution conditionnelles

$X \backslash Y$	60	80	100	Total
300	15	4	1	20
$f_{j/i}$	75%	20%	5%	

$X \backslash Y$	$y_1$	$\dots$	$y_j$	$\dots$	$y_c$	Total
$y_j$ sachant que $X = x_i$	$n_{i1}$	$\dots$	$n_{ij}$	$\dots$	$n_{ic}$	$n_{i.}$

Fréquence de  $Y = y_j$  sachant que  $X = x_i$  :

$$f_{j/i} = \frac{n_{ij}}{n_{i.}} = \frac{f_{ij}}{f_{i.}}$$

# Indépendance

## Définition

$X$  et  $Y$  sont *indépendantes* si la distribution conditionnelle de  $Y$  sachant  $X = x_i$  ne dépend pas de  $i$  :

$$\forall i, j, \quad f_{j/i} = f_{.j}$$
$$f_{ij} = f_{i.} \times f_{.j}$$

$X \backslash Y$	60	80	100	$f_{i.}$
300	16,7%	4,4%	1,1%	22,2%
350	11,1%	22,2%	11,1%	44,4%
400	5,6%	11,1%	16,7%	33,3%
$f_{.j}$	33,3%	37,8%	28,9%	

$$0,333 \times 0,222 \neq 0,167$$

$X$  et  $Y$  ne sont *pas indépendantes*.

# Sommaire

- 1 Distributions à deux caractères
- 2 Covariance
- 3 Coefficient de corrélation linéaire

# Covariance : un exemple



$$\begin{aligned}\mathbb{C}ov(X,Y) &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) (y_i - \bar{y}) \\ &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}\end{aligned}$$

X engrais	20	24	28	22	32	28	32	36	41	41
Y rendement	16	18	23	24	28	29	26	31	32	34

# Covariance : un exemple

$$\bar{x} = \frac{1}{10} \sum_{i=1}^{10} x_i = 30,4$$

$$\bar{y} = \frac{1}{10} \sum_{i=1}^{10} y_i = 26,1$$

$$\frac{1}{n} \sum_{i=1}^n x_i y_i = 828,6$$

$$\begin{aligned} \text{Cov}(X, Y) &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y} \end{aligned}$$

$$\text{Cov}(X, Y) = 828,6 - 30,4 \times 26,1 = 35,16$$



# Propriétés de la covariance

## Proposition

- *Symétrie* :  $\mathbb{C}ov(X,Y) = \mathbb{C}ov(Y,X)$
- *lien avec la variance* :  $\mathbb{C}ov(X,X) = \mathbb{V}(X)$
- *transformation affine* :  $\mathbb{C}ov(aX + b, cY + d) = a c \mathbb{C}ov(X,Y)$
- *Si  $X$  et  $Y$  sont indépendantes alors  $\mathbb{C}ov(X,Y) = 0$ .*



La réciproque est **fausse** !

$\mathbb{C}ov(X,Y)=0$  **n'entraîne pas** que  $X$  et  $Y$  sont indépendantes

# Variance d'une somme

## Théorème

$$\mathbb{V}(X + Y) = \mathbb{V}(X) + 2\text{Cov}(X, Y) + \mathbb{V}(Y)$$

Cas de variables **décorrélées** :  $\mathbb{V}(X + Y) = \mathbb{V}(X) + \mathbb{V}(Y)$ .

# Inégalité de Cauchy-Schwarz

Pour  $t \in \mathbb{R}$  :

$$0 \leq \mathbb{V}(X + tY) = \mathbb{V}(X) + 2t \mathbb{Cov}(X, Y) + t^2 \mathbb{V}(Y)$$

$$\Delta = [2\mathbb{Cov}(X, Y)]^2 - 4\mathbb{V}(X)\mathbb{V}(Y) = 4[\mathbb{Cov}(X, Y)^2 - \mathbb{V}(X)\mathbb{V}(Y)] \leq 0$$

## **Théorème (Inégalité de Cauchy-Schwarz)**

$$|\mathbb{Cov}(X, Y)| \leq \sqrt{\mathbb{V}(X)\mathbb{V}(Y)} = \sigma(X)\sigma(Y)$$

# Sommaire

- 1 Distributions à deux caractères
- 2 Covariance
- 3 Coefficient de corrélation linéaire

# Coefficient de corrélation linéaire

## Définition

$$r(X,Y) = \frac{\text{Cov}(X,Y)}{\sigma(X) \sigma(Y)}$$

On a

$$-1 \leq r(X,Y) \leq 1$$

$r(X,Y) \simeq 0 \Rightarrow$  absence de relation linéaire (décorrélation)  
 $\nRightarrow$  indépendance

# Exemple

Corrélation entre rendement et quantité d'engrais d'une parcelle de blé

$X$ engrais	20	24	28	22	32	28	32	36	41	41
$Y$ rendement	16	18	23	24	28	29	26	31	32	34

$$\sigma(X) \simeq 7,40 \quad \sigma(Y) \simeq 5,91 \quad \text{Cov}(X,Y) \simeq 35,16$$

$$r(X,Y) \simeq \frac{35,16}{7,40 \times 5,91} \simeq 0,89$$

Il y a **corrélation linéaire forte**.

