

# Un exemple de régression linéaire

## Exercice 1 (Calculs pour une régression linéaire)

Soient  $p \in ]0; 1[$ , et  $X, Y \hookrightarrow \mathcal{G}(p)$  indépendantes.

On s'intéresse au couples de variables aléatoires :  $\triangleright S = \max(X, Y)$

$$\triangleright I = \min(X, Y)$$

1. Justifier les relations *a priori* :
- $\triangleright I + S = X + Y.$
  - $\triangleright I \cdot S = X \cdot Y.$
  - $\triangleright I^2 + S^2 = X^2 + Y^2.$

2. **Rappels sur la loi géométrique** On a choisi  $X \hookrightarrow \mathcal{G}(p)$ .

a) Rappeler  $X(\Omega)$ , et pour  $n \in X(\Omega)$ , la probabilité  $\mathbb{P}(X = n)$ .

$$(\cdot)^{1-u} b \cdot d = (u = X) \mathbb{P} \quad \text{, } 1 \leq u \text{ et pour } n \in \mathbb{N} : \text{Réponse : } X(\Omega) = \mathbb{N}^*$$

b) Rappeler l'expression de  $\mathbb{E}[X]$  et de  $\text{Var}(X)$ .

Par la formule de Kœnig-Huygens, en déduire :  $\mathbb{E}[X^2] = \frac{1+q}{p^2}.$

$$(\cdot)^{\frac{d}{b}} = \frac{\mathbb{E}[X^2]}{\mathbb{E}[X]^2} = \frac{\text{Var}(X) + \mathbb{E}[X]^2}{\mathbb{E}[X]^2} : \text{Réponse : } \frac{d}{b} = \frac{\mathbb{E}[X^2]}{\mathbb{E}[X]^2}$$

c) Rappeler, pour  $n \in \mathbb{N}$ , la probabilité :  $\mathbb{P}(X > n).$  (fonction d'*antirépartition* de  $X$ .)

$$(\cdot)^{1-u} b \cdot d = (u = X) \mathbb{P} \quad \text{, } 1 \leq u \text{ et pour } n \in \mathbb{N} : \text{Réponse : } \mathbb{P}(X > n) = \sum_{k=n+1}^{\infty} p \cdot q^{k-1}$$

3. **Lois marginales de  $I, S$**

- a) Montrer l'égalité d'événements :
- $\triangleright [S \leq k] = [X \leq k] \cap [Y \leq k],$
  - $\triangleright [I > k] = [X > k] \cap [Y > k].$

- b) En déduire :
- $\triangleright \mathbb{P}(S \leq n) = 1 - 2q^n + (q^2)^n,$
  - $\triangleright \mathbb{P}(I > n) = (q^2)^n.$

- c) Déduire enfin :
- $\triangleright \mathbb{P}(S = n) = 2 \cdot p \cdot q^{n-1} - (1 - q^2) \cdot (q^2)^{n-1}, \text{ pour } n \geq 1,$
  - $\triangleright I \hookrightarrow \mathcal{G}(1 - q^2).$

4. **Calcul de la covariance**

a) Justifier que l'on a :  $\mathbb{E}[I] = \frac{1}{1-q^2}.$

b) Montrer :  $\mathbb{E}[S] = \mathbb{E}[X] + \mathbb{E}[Y] - \mathbb{E}[I].$

(On utilisera l'une des relations de 1.)

En déduire :  $\mathbb{E}[S] = \frac{2p-1}{p^2}.$

c) Montrer :  $\mathbb{E}[I \cdot S] = \mathbb{E}[X] \cdot \mathbb{E}[Y].$

(On utilisera 1., et l'indépendance de  $X, Y$ .)

En déduire :  $\mathbb{E}[I \cdot S] = \frac{1}{p^2}.$

d) Obtenir enfin :  $\text{Cov}(I, S) = \frac{1}{p^2} - \frac{2p-1}{p^3 \cdot (2-p)} = \frac{1}{p^3 \cdot (2-p)} [p \cdot (2-p) - (2p-1)] = \frac{1-p^2}{p^3 \cdot (2-p)}.$

**Exercice 2 (Calculs pour la régression linéaire)**

Soient  $X$  et  $Y$  deux variables aléatoires admettant un moment d'ordre 2.

Toutes les espérances, variances, covariances apparaissant convergent donc bien.

On fait la régression linéaire de  $Y$  par  $X$ , en approximant  $Y$  par  $\hat{Y} = a \cdot X + b$ .

On cherche pour quelles valeurs de  $a, b \in \mathbb{R}$ , l'erreur quadratique  $r_{a,b} = \mathbb{E}[(Y - aX - b)^2]$ .

**1. Détermination des coefficients  $a, b$  optimaux**

a) Par la formule de Koenig-Huygens, montrer :  $r_{a,b} = \text{Var}(Y - aX) + (\mathbb{E}[Y] - a \cdot \mathbb{E}[X] - b)^2$ .

b) Montrer que :  $\text{Var}(Y - aX) = \left[ a \cdot \sigma_X - \frac{\text{Cov}(X,Y)}{\sigma_X} \right]^2 + \text{Var}(Y) - \frac{1}{\text{Var}(X)} \cdot [\text{Cov}(X,Y)]^2$   
 (On écrira :  $\text{Var}(Y - aX) = \text{Var}(Y) - 2a \cdot \text{Cov}(X,Y) + a^2 \cdot \sigma_X^2$ )

c) En déduire que  $r_{a,b}$  est minimisée pour :  
 ▶  $a = \frac{\text{Cov}(X,Y)}{\text{Var}(X)}$ ,  
 ▶  $b = \mathbb{E}[Y] - a \cdot \mathbb{E}[X]$ .

2. On choisit ces valeurs de  $a, b$ . Le résidu est noté :  $\epsilon = Y - a \cdot X - b$ .

a) Montrer que  $\mathbb{E}[\epsilon] = 0$ . En déduire :  $\mathbb{E}[a \cdot X + b]$ .

b) Montrer que :  $\text{Var}(\epsilon) = \text{Var}(Y) - \frac{1}{\text{Var}(X)} \cdot [\text{Cov}(X,Y)]^2$ .

c) Montrer que :  $\text{Cov}(X, \epsilon) = 0$ .

En déduire que :  $\text{Var}(a \cdot X + b) + \text{Var}(\epsilon) = \text{Var}(Y)$

d) Montrer :  
 ▶  $\text{Var}(a \cdot X + b) = \rho^2(X,Y) \cdot \text{Var}(Y)$ ,  
 ▶  $\text{Var}(\epsilon) = [1 - \rho^2(X,Y)] \cdot \text{Var}(Y)$ .

**Détermination pratique en Scilab**

La commande `reglin` permet de faire la régression linéaire d'un échantillon  $x$  par  $y$ .

**▶ Syntaxe**

```
1 // x,y sont deux échantillons de même longueur
2 [a,b] = reglin(x,y) // la régression linéaire s'écrit alors : yr = a*x + b
```

**▶ Écart-type du résidu**

La syntaxe `[a,b,sig] = reglin(x,y)` retourne de plus `sig` l'écart-type du résidu.

**▶ Part expliquée, part inexpliquée de la variance**

▶ On obtient la part expliquée  $\text{Var}(a \cdot X + b) = \rho^2(X,Y) \cdot \text{Var}(Y)$  comme suit :

```
1 partExpliquee = variance(a*x+b)
```

▶ On obtient la part inexpliquée  $\text{Var}(\epsilon) = [1 - \rho^2(X,Y)] \cdot \text{Var}(Y)$  comme suit :

```
1 partInexpliquee = sig^2
2 coefDetermination = 1 - sig^2 / variance(Y)
3 // coefDetermination proche de 1 <=> régression linéaire de qualité
4 // coefDetermination proche de 0 <=> piètre régression linéaire
```