

Cours - Informatique - Régression linéaire

Gustave Cortal

école —————
normale —————
supérieure —————
paris-saclay —————

université
PARIS-SACLAY

Exemple introductif

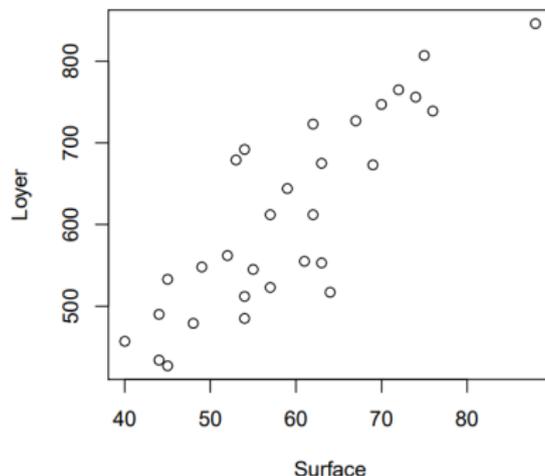


Figure: Prix du loyer en fonction de la surface de l'appartement.

Le prix d'un appartement est une variable aléatoire Y dont l'espérance dépend de la surface x . Y est la *variable à expliquer* et x est la *variable explicative*, supposée connue et non aléatoire.

Problèmes abordés

Les problèmes sont multiples :

- ▶ spécifier le modèle ;
- ▶ estimer les paramètres du modèle ;
- ▶ vérifier qu'il y a bien une relation entre les deux variables ;
- ▶ vérifier la validité du modèle retenu ;
- ▶ prédire le prix d'un nouvel appartement en fonction de sa surface.

Modèle de la régression linéaire simple

L'espérance de la variable aléatoire à expliquer Y est une fonction linéaire de la variable x :

$$E(Y_i) = a + bx_i \quad (1)$$

Nous avons :

$$Y_i = a + bx_i + \epsilon_i \quad i = 1, \dots, n. \quad (2)$$

avec $E(\epsilon_i) = 0$. On suppose que les ϵ_i sont **indépendants** et suivent **la même loi normale** $N(0, \sigma^2)$. Le modèle est alors :

$$Y_i \sim N(a + bx_i, \sigma^2), \quad i = 1, \dots, n. \quad (3)$$

Hypothèses fondamentales du modèle

- ▶ **Linéarité** : $E(Y_i)$ est une forme linéaire des paramètres :
 $E(Y_i) = a + bx_i$
- ▶ **Normalité des erreurs** : les erreurs suivent une loi normale
- ▶ **Homoscédasticité** : les erreurs partagent la même variance σ^2
- ▶ **Indépendance** : les erreurs sont indépendantes

Estimation des paramètres du modèle (1)

On estime les paramètres du modèle avec la méthode du maximum de vraisemblance. Soit la fonction de densité de Y_i :

$$f(Y_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(Y_i - a - bx_i)^2\right) \quad i = 1, \dots, n. \quad (4)$$

Les variables Y_i sont indépendantes. On en déduit la fonction de vraisemblance de l'échantillon :

$$L(a, b, \sigma^2; y_1, \dots, y_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(Y_i - a - bx_i)^2}{2\sigma^2}\right) \quad (5)$$

et la fonction de log-vraisemblance :

$$l(a, b, \sigma^2; y_1, \dots, y_n) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - a - bx_i)^2 \quad (6)$$

Estimation des paramètres du modèle (2)

Pour maximiser la vraisemblance par rapport à a et b , il faut minimiser :

$$Q(a, b) = \sum_{i=1}^n (Y_i - a - bx_i)^2 \quad (7)$$

$Q(a, b)$ est le *critère des moindres carrés*. Pour le minimiser, il faut résoudre le système à deux équations :

$$\frac{\partial Q}{\partial a} = 0 \quad \frac{\partial Q}{\partial b} = 0 \quad (8)$$

On trouve les *estimateurs des moindres carrés* \hat{a} et \hat{b} :

$$\hat{b} = \frac{S_{xY}}{s_x^2} \quad \hat{a} = \bar{Y} - \frac{S_{xY}}{s_x^2} \bar{x} \quad (9)$$

S_{xY} est la covariance empirique des x_i et des Y_i , et s_x^2 la variance empirique des x_i . \bar{Y} est la moyenne empirique des Y_i et \bar{x} la moyenne empirique des x_i .

Estimation des paramètres du modèle (3)

Nous obtenons la droite des moindres carrés de Y en x :

$$\hat{Y}_i = \hat{a} + \hat{b}x_i \quad \text{et} \quad \hat{\epsilon}_i = Y_i - \hat{Y}_i \quad i = 1, \dots, n \quad (10)$$

$\hat{\epsilon}_i$ sont appelées *résidus*.

L'estimateur du maximum de vraisemblance de σ^2 est :

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{a} - \hat{b}x_i)^2 = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 \quad (11)$$

C'est la moyenne des carrés des résidus.

Analyse de la variance

Équation d'analyse de la variance :

$$\frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2 = \frac{1}{n} \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (12)$$

que l'on peut noter :

$$S_Y^2 = S_{reg} + S_{res} \quad (13)$$

La variance totale est égale à la variance expliquée par la régression et la variance résiduelle. Les *estimateurs des moindres carrés* minimisent la variance résiduelle, donc maximise la variance expliquée par la régression.

Mesure de l'ajustement

On mesure la qualité de l'ajustement de la régression par le *coefficient de détermination* :

$$R^2 = \frac{S_{reg}}{S_Y^2} \quad (14)$$

C'est une quantité variant entre 0 et 1. Une valeur proche de 1 indique que la proportion de variance expliquée par la régression est grande, la régression est donc très explicative.

Prédiction sur de nouvelles valeurs

Soit une nouvelle valeur x_0 (e.g., une nouvelle surface), on peut prédire grâce au modèle de régression Y_0 (e.g., le prix de l'appartement correspondant) :

$$\hat{Y}_0 = \hat{a} + \hat{b}x_0 \quad (15)$$

Aller plus loin

- ▶ Propriétés des estimateurs \hat{a} , \hat{b} et $\hat{\sigma}^2$ (e.g., quelles lois suivent-ils?)
- ▶ Estimateur sans biais de σ^2
- ▶ Intervalles de confiance sur \hat{a} et \hat{b}
- ▶ Tester la non-nullité de la pente \hat{b}
- ▶ Vérification des hypothèses du modèle