

Régression linéaire multiple

Frédéric Bertrand et Myriam Maumy¹

¹IRMA, Université Louis Pasteur
Strasbourg, France

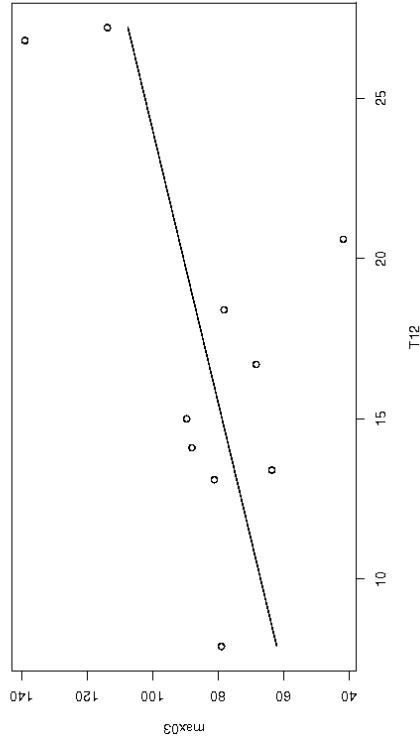
Master 1ère Année 06-02-2008

- **Problème** : Étude de la concentration d'ozone dans l'air.
- **Modèle** : La température (v.a. X) et la concentration d'ozone (v.a. Y) sont liées de manière linéaire :

$$Y = \beta_0 + \beta_1 X + \varepsilon.$$

- **Observations** : $n = 10$ mesures de la température et de la concentration d'ozone.
- **But** : Estimer β_0 et β_1 afin de prédire la concentration d'ozone connaissant la température.

Frédéric Bertrand et Myriam Maumy	Régression linéaire multiple
Introduction	Régression linéaire multiple
Présentation du modèle	Régression linéaire simple
Méthode des moindres carrés	Exemple
Propriétés des moindres carrés	Affiner le modèle
Tests	



Souvent la régression linéaire est trop simpliste. Il faut alors utiliser d'autres modèles plus réalistes mais parfois plus complexes :

- Utiliser d'autres fonctions que les fonctions affines comme les fonctions polynômiales, exponentielles, logarithmiques...
- Considérer plusieurs variables explicatives.

Exemple : La température **et** la vitesse du vent

- **Problème :** On cherche $\hat{\beta}$ qui annule cette dérivée. Donc on doit résoudre l'équation suivante :

$${}^t\mathbf{XX}\hat{\beta} = {}^t\mathbf{Xy}.$$

- **Solution :** On trouve après avoir inversé la matrice ${}^t\mathbf{XX}$ (il faut naturellement vérifier que ${}^t\mathbf{XX}$ est carrée et inversible c'est-à-dire qu'aucune des colonnes qui compose cette matrice ne soit proportionnelle aux autres colonnes)

$$\hat{\beta} = ({}^t\mathbf{XX})^{-1}{}^t\mathbf{Xy}.$$

$$\begin{aligned} \|\varepsilon\|^2 &= {}^t(\mathbf{y} - \mathbf{X}\hat{\beta})(\mathbf{y} - \mathbf{X}\hat{\beta}) \\ &= {}^t\mathbf{yy} - {}^t\hat{\beta}{}^t\mathbf{Xy} - {}^t\mathbf{yX}\hat{\beta} + {}^t\hat{\beta}{}^t\mathbf{XX}\hat{\beta} \\ &= {}^t\mathbf{yy} - 2{}^t\hat{\beta}{}^t\mathbf{Xy} + {}^t\hat{\beta}{}^t\mathbf{XX}\hat{\beta} \end{aligned}$$

car ${}^t\hat{\beta}{}^t\mathbf{Xy}$ est un scalaire. Donc il est égal à sa transposée.

La dérivée par rapport à $\hat{\beta}$ est alors égale à :

$$-2{}^t\mathbf{Xy} + 2{}^t\mathbf{XX}\hat{\beta}.$$

Retrouvons les résultats de la régression linéaire simple ($p = 2$)

$${}^t\mathbf{XX} = \begin{pmatrix} n & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i y_i \end{pmatrix}; \quad {}^t\mathbf{Xy} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}.$$

Donc :

$$\begin{aligned} ({}^t\mathbf{XX})^{-1} &= \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix} \\ &= \frac{1}{\sum (x_i - \bar{x})^2} \begin{pmatrix} \sum x_i^2/n & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix}. \end{aligned}$$

Finalement on retrouve bien :

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \begin{pmatrix} \frac{\bar{y} \sum x_i^2 - \bar{x} \sum x_i y_i}{\sum (x_i - \bar{x})^2} \\ \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum (x_i - \bar{x})^2} \end{pmatrix}$$

ce qui correspond aux estimateurs de la régression linéaire simple que nous avons déjà rencontrés dans le cours 1.

Résultats préliminaires :

```
> a <- lm(max03 ~ T12 + VX)
> summary(a)

Call :
lm(formula = max03 ~ T12 + VX)

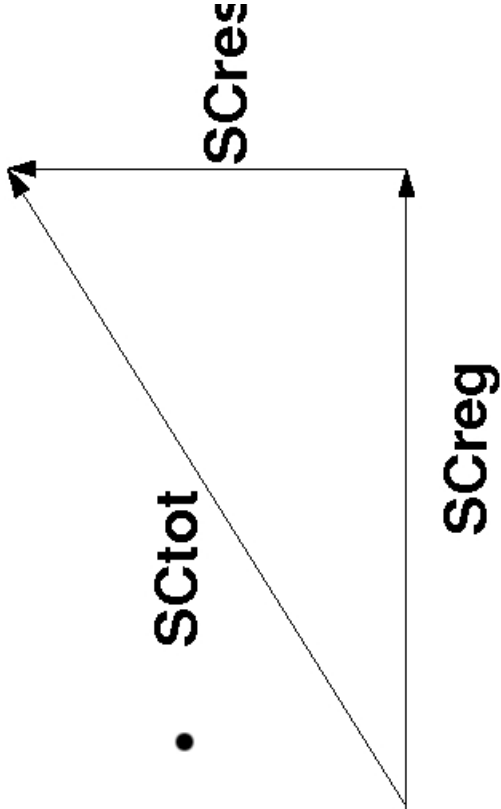
Residuals :
Min 1Q Median 3Q Max
-47.860 -10.561  5.119 10.645 26.506

Coefficients :
Estimate Std. Error t value Pr(>|t|)
(Intercept) 36.6520 26.5324  1.381 0.210
T12 2.6623 1.4202  1.875 0.103
VX 0.5431 0.7775  0.699 0.507
Residual standard error: 24.78 on 7 degrees of freedom
Multiple R-Squared: 0.3351, Adjusted R-squared: 0.1452
F-statistic: 1.764 on 2 and 7 DF, p-value: 0.2396
```

- $\sum \hat{y}_i^2 = \sum \hat{y}_i y_i$ ou (forme matricielle) ${}^t\hat{y}\hat{y} = {}^t y\hat{y}$
- $\sum \hat{y}_i = \sum y_i$

Propriété des moindres carrés :

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$
$$SC_{tot} = SC_{reg} + SC_{res}$$



Le coefficient de détermination est défini par :

$$R^2 = \frac{SC_{reg}}{SC_{tot}}.$$

Intuitivement ce coefficient de détermination quantifie la capacité du modèle à expliquer les variations de Y .

- Si R^2 est proche de 1 alors le modèle est proche de la réalité.
- Si R^2 est proche de 0 alors le modèle explique très mal la réalité. Il faut alors trouver un meilleur modèle.

<div> <div> Hypothèses pour les tests Hypothèses pour les tests Estimation de σ^2 Tests d'hypothèses </div> <div> Introduction Présentation du modèle Méthode des moindres carrés Propriétés des moindres carrés Tests </div> </div>	<div> <div> Hypothèses pour les tests Hypothèses pour les tests Estimation de σ^2 Tests d'hypothèses </div> <div> Introduction Présentation du modèle Méthode des moindres carrés Propriétés des moindres carrés Tests </div> </div>
<div> <div> Méthode : <ul style="list-style-type: none"> Calculer la statistique </div> <div> $t_{obs} = \frac{\hat{\beta}_j - b_j}{s(\hat{\beta}_j)}$ </div> <div> <ul style="list-style-type: none"> où $s^2(\hat{\beta}_j)$ est l'élément diagonal d'indice j de $s^2({}^t\mathbf{XX})^{-1}$. Si l'hypothèse nulle \mathcal{H}_0 est vraie, alors t_{obs} suit une loi de Student avec $(n - p)$ degrés de liberté. </div> </div>	<div> <div> <ul style="list-style-type: none"> Valeur critique : $t_{(\alpha/2, n-p)}$ le $(1 - \alpha/2)$-quantile d'une loi de Student avec $(n - p)$ degrés de liberté (cf table de la loi de Student). On rejette l'hypothèse nulle \mathcal{H}_0 si $t_{obs} \geq t_{(\alpha/2, n-p)}$. </div> <div> <p>Cas particulier : Tester si « $\beta_j = 0$ » pour un certain j.</p> <p>Si l'hypothèse nulle \mathcal{H}_0 : « $\beta_j = 0$ » est acceptable alors la variable X_j n'est pas significative au sein du modèle. On peut simplifier le modèle, ...et recommencer !</p> </div> </div>
<div> <div> Frédéric Bertrand et Myriam Maumy </div> <div> Régression linéaire multiple </div> </div>	<div> <div> Frédéric Bertrand et Myriam Maumy </div> <div> Régression linéaire multiple </div> </div>