

# SY09

## Régression linéaire

T. Dencœur

### 1 Principe

Nous avons vu dans le chapitre sur la théorie de la décision que la fonction de décision  $g$  minimisant le risque quadratique

$$R(g) = \mathbb{E}_{\mathbf{X}, Y}[(g(\mathbf{X}) - Y)^2]$$

est la fonction de régression :

$$g^*(\mathbf{x}) = \mathbb{E}(Y|\mathbf{x}).$$

Dans le modèle de régression linéaire, on suppose que cette fonction est une fonction affine de  $\mathbf{x}$  que l'on note

$$g^*(\mathbf{x}) = w_0^* + \sum_{j=1}^p w_j^* x_j = \mathbf{w}^{*'} \mathbf{x}$$

avec, par convention,  $\mathbf{x} = (1, x_1, \dots, x_p)'$ .

Le problème posé ici consiste à estimer le vecteur  $\mathbf{w}^*$  à partir d'un ensemble d'apprentissage  $\mathcal{L} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ . Pour cela, on remarque que  $\mathbf{w}^*$  s'obtient comme solution d'un problème d'optimisation :

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} R(\mathbf{w})$$

avec  $R(\mathbf{w}) = \mathbb{E}_{\mathbf{X}, Y}[\mathbf{w}'\mathbf{X} - Y]^2$ . On ne peut en pratique résoudre ce problème de manière exacte car la fonction  $R(\mathbf{w})$  est inconnue, mais on peut remplacer le risque théorique par le *risque empirique* défini par :

$$\hat{R}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}'\mathbf{x}_i - y_i)^2.$$

On notera  $\hat{\mathbf{w}}$  le vecteur de coefficients minimisant le risque empirique :

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \hat{R}(\mathbf{w}).$$

Intuitivement, lorsque  $n$  est assez grand,  $\hat{R}(\mathbf{w})$  sera "proche" de  $R(\mathbf{w})$ , et  $\hat{\mathbf{w}}$  sera donc "proche" de  $\mathbf{w}^*$ . Le vecteur  $\hat{\mathbf{w}}$  est appelé *estimateur des moindres carrés* de  $\mathbf{w}^*$ .

## 2 Méthode des moindres carrés

Notons

$$X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}$$

la matrice  $(n, p+1)$  contenant les valeurs des variables explicatives et  $\mathbf{y} = (y_1, \dots, y_n)'$  le vecteur des observations de la variable  $y$ . Le risque empirique  $\hat{R}(\mathbf{w})$  peut alors s'écrire matriciellement :

$$\begin{aligned} \hat{R}(\mathbf{w}) &= \frac{1}{n} (\mathbf{X}\mathbf{w} - \mathbf{y})' (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= \frac{1}{n} (\mathbf{w}' \mathbf{X}' \mathbf{X} \mathbf{w} - 2\mathbf{w}' \mathbf{X}' \mathbf{y} + \mathbf{y}' \mathbf{y}). \end{aligned}$$

On a

$$\frac{d\hat{R}(\mathbf{w})}{d\mathbf{w}} = \frac{1}{n} (2\mathbf{X}' \mathbf{X} \mathbf{w} - 2\mathbf{X}' \mathbf{y}),$$

d'où

$$\frac{d\hat{R}(\mathbf{w})}{d\mathbf{w}} = 0 \Leftrightarrow \mathbf{X}' \mathbf{X} \mathbf{w} = \mathbf{X}' \mathbf{y}. \quad (1)$$

Si la matrice  $\mathbf{X}' \mathbf{X}$  est inversible, le minimum du risque empirique est donc obtenu pour

$$\hat{\mathbf{w}} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y}.$$

On notera

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\mathbf{w}} = \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y}$$

le vecteur des prédictions obtenu en remplaçant le paramètre  $\mathbf{w}$  inconnu par son estimateur des moindres carrés  $\hat{\mathbf{w}}$ .

**Remarque 1** On a supposé  $\mathbf{X}' \mathbf{X}$  inversible, ce qui est le cas si la matrice  $\mathbf{X}$  est de rang  $p+1$ . Si ce n'est pas le cas, c'est qu'une variable (une colonne de  $\mathbf{X}$ ) s'exprime comme combinaison linéaire des autres. Il suffit alors de supprimer la ou les variables redondantes.

**Remarque 2** Si certaines variables sont très corrélées, la matrice  $\mathbf{X}' \mathbf{X}$  est mal conditionnée et les calculs numériques peuvent être très imprécis. Une solution (appelée ridge regression en anglais) consiste à ajouter un terme sur la diagonale de  $\mathbf{X}' \mathbf{X}$  :

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}' \mathbf{y}$$

où  $\lambda$  est une constante à déterminer. On montre que l'on améliore ainsi parfois les propriétés de l'estimateur.

**Remarque 3** On a

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\mathbf{w}} = \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} = \mathbf{P} \mathbf{y}$$

en notant  $\mathbf{P} = \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'$ . Cette matrice  $\mathbf{P}$  a des propriétés remarquables. En effet,  $\mathbf{P}$  est symétrique (évident), et de plus

$$\mathbf{P}^2 = \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' = \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' = \mathbf{P}.$$

La matrice  $P$  est donc idempotente (c'est un opérateur de projection orthogonale, comme nous le verrons par la suite). De même, on peut écrire

$$\hat{\varepsilon} = \mathbf{y} - \hat{\mathbf{y}} = (I_n - P)\mathbf{y} = R\mathbf{y}$$

avec  $R = I_n - P$ . On vérifie aisément que  $R$  a les mêmes propriétés que  $P$  (symétrie et idempotence) : c'est également un opérateur de projection orthogonale.

### 3 Analyse de la variance

#### 3.1 Point de vue géométrique

Plaçons nous dans  $\mathbb{R}^n$  et considérons les vecteurs

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \quad \mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \quad \mathbf{x}_j = \begin{pmatrix} x_{1j} \\ \vdots \\ x_{nj} \end{pmatrix}, \quad j = 1, p$$

La méthode des moindres carrés peut être interprétée comme la recherche de la meilleure approximation de  $\mathbf{y}$  dans le sous-espace  $\mathcal{L}$  de  $\mathbb{R}^n$  engendré par les  $p + 1$  vecteurs  $\mathbf{1}, \mathbf{x}_1, \dots, \mathbf{x}_p$ . On cherche en effet

$$\hat{\mathbf{y}} = \hat{w}_0 \mathbf{1} + \sum_{j=1}^p \hat{w}_j \mathbf{x}_j \in \mathcal{L}$$

tel que la distance euclidienne  $\|\mathbf{y} - \hat{\mathbf{y}}\|^2$  soit minimum. On sait que la solution consiste à définir  $\hat{\mathbf{y}}$  comme la projection orthogonale de  $\mathbf{y}$  sur  $\mathcal{L}$ . On a vu en effet que

$$\hat{\mathbf{y}} = P\mathbf{y},$$

$P$  étant un opérateur de projection orthogonale.

Cette représentation géométrique permet de retrouver sans calculs fastidieux plusieurs résultats intéressants. Tout d'abord,

$$\hat{\varepsilon} \perp \mathbf{1} \Rightarrow \sum_{i=1}^n \hat{\varepsilon}_i = 0,$$

d'où l'on déduit

$$\frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}.$$

Par ailleurs, la projection orthogonale de  $\mathbf{y}$  sur l'axe dirigé par  $\mathbf{1}$  a pour coordonnée

$$\frac{\langle \mathbf{y}, \mathbf{1} \rangle}{\|\mathbf{1}\|} = \bar{y}.$$

Il en est de même, d'après ce qui précède, pour la projection orthogonale de  $\hat{\mathbf{y}}$  sur  $\mathbf{1}$ . Enfin, on a de manière évidente :

$$\hat{\mathbf{y}} \perp \hat{\varepsilon}.$$

### 3.2 Equation d'analyse de la variance

Notons  $\bar{y} = \bar{y}1$ . En appliquant le théorème de Pythagore au triangle  $(y, \hat{y}, \bar{y})$ , on obtient finalement la relation très importante suivante, appelée *équation d'analyse de la variance* :

$$\|y - \bar{y}\|^2 = \|\hat{y} - \bar{y}\|^2 + \|\hat{\varepsilon}\|^2,$$

ce que l'on peut encore écrire, en divisant par  $n$  :

$$\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2$$

soit encore

$$S_{YY} = S_{reg} + S_{res}.$$

Cette équation est appelée *équation d'analyse de la variance*. Le terme de gauche ( $S_{YY}$ ) est la variance empirique des  $Y_i$ , il caractérise la dispersion des valeurs observées de la variable à expliquer. Le premier terme du membre de droite ( $S_{reg}$ ) est la variance empirique des  $\hat{Y}_i$ , que l'on appelle variance expliquée par la régression. Le second terme du membre de droite ( $S_{res}$ ) est la variance des résidus, ou variance résiduelle.

**Remarque 4** A chacun des termes de l'équation d'analyse de la variance est associé un nombre de degrés de liberté (d.d.l.), égal au nombre de combinaisons linéaires des  $Y_i$  utilisées dans le calcul :

- $S_{YY}$  dépend de  $n$  quantités  $y_1 - \bar{y}, \dots, y_n - \bar{y}$  liées par la relation

$$\sum_{i=1}^n (y_i - \bar{y}) = 0.$$

Ce terme a donc  $n - 1$  d.d.l.

- On a  $\hat{y}_i = \mathbf{x}_i' \hat{\mathbf{w}}$  et  $\bar{y} = \bar{\mathbf{x}}' \hat{\mathbf{w}}$ . Par conséquent, le terme  $S_{reg}$  est fonction des paramètres  $\hat{w}_1, \dots, \hat{w}_p$  (le terme  $\hat{w}_0$  s'annule dans chacune des différences  $\hat{y}_i - \bar{y}$ ). La variance expliquée a donc  $p$  d.d.l.
- Par conséquent, le nombre de d.d.l associé à la variance résiduelle est  $n - p - 1$ .

La plupart des logiciels statistiques présentent les résultats de la régression sous forme d'un tableau (appelé *tableau d'analyse de la variance*), où figurent les différents termes de l'équation d'analyse de la variance, et les nombres de d.d.l associés (cf. tableau 1).

### 3.3 Evaluation de la qualité de l'ajustement

On définit à partir de l'équation d'analyse de la variance le *coefficient de détermination*, égal à la proportion de la variance totale expliquée par la régression :

$$R^2 = \frac{S_{reg}}{S_{YY}} = 1 - \frac{S_{res}}{S_{YY}}.$$

Ce coefficient traduit la « qualité de l'ajustement », comme on le voit en considérant les deux situations extrêmes suivantes :

TABLE 1 – Tableau d'analyse de la variance (SS : *sum of squares*; MS : *mean square*).

source de variation	d.d.l.	SS	MS=SS/d.d.l.
régression	$p$	$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	$\frac{1}{p} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$
résiduelle	$n - p - 1$	$\sum_{i=1}^n (y_i - \hat{y}_i)^2$	$\frac{1}{n - p - 1} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \hat{\sigma}^2$
totale	$n - 1$	$\sum_{i=1}^n (y_i - \bar{y})^2$	$\frac{1}{n - 1} \sum_{i=1}^n (y_i - \bar{y})^2$

- Si les résidus sont nuls, on a  $S_{res} = 0$  et  $R^2 = 1$ . Les  $n$  points  $(\mathbf{x}_i, y_i) \in \mathbb{R}^{p+1}$  sont alors situés dans l'hyperplan d'équation

$$y = \hat{w}_0 + \hat{w}_1 x_1 + \dots + \hat{w}_p x_p.$$

Cela signifie que l'on peut retrouver sans erreur les  $y_i$  à partir des  $\mathbf{x}_i$ , c'est-à-dire que toute la variation des  $y_i$  est expliquée par les  $\mathbf{x}_i$ .

- Si les prédictions sont constantes ( $\hat{y}_i = \bar{y}, \forall i$ ), la variance expliquée est nulle et  $R^2 = 0$ . Dans ce cas, les  $\mathbf{x}_i$  n'expliquent pas du tout la variation des  $y_i$ .
- De manière générale, on a  $0 \leq R^2 \leq 1$ , et la valeur du  $R^2$  s'interprète comme un « degré de liaison » entre les variables explicatives et la variable à expliquer.

**Remarque 5** Géométriquement,  $R^2$  est égal au carré du cosinus de l'angle  $\theta$  entre les vecteurs  $\mathbf{y} - \bar{y}$  et  $\hat{\mathbf{y}} - \bar{y}$  : c'est donc le carré du coefficient de corrélation linéaire entre les  $y_i$  et les  $\hat{y}_i$ .