

Régression linéaire multiple

Kaoutar IKKANE

FPO

Mathématiques appliquées pour la science de données

December 26, 2024

1. Introduction
2. Modèle
3. Estimation
4. Ecart au modèle : Résidus
5. Décomposition de la variabilité
6. Conclusion



Introduction

Le but de la régression simple (resp. multiple) est d'expliquer une variable Y à l'aide d'une variable X (resp. plusieurs variables X_1, \dots, X_q). La variable Y est appelée variable dépendante, ou variable à expliquer et les variables X_j ($j = 1, \dots, q$) sont appelées variables indépendantes, ou variables explicatives.



Modèle



Une variable quantitative Y est mise en relation avec p variables quantitatives X^1, \dots, X^p . Les données sont supposées provenir de l'observation d'un échantillon statistique de taille n ($n > p + 1$) de $\mathbb{R}^{(p+1)}$:

$$(x_i^1, \dots, x_i^j, \dots, x_i^p, y_i), \quad i = 1, \dots, n.$$

Le modèle s'écrit :

$$Y_i = \beta_0 + \beta_1 x_i^1 + \beta_2 x_i^2 + \cdots + \beta_p x_i^p + \epsilon_i \quad i = 1, 2, \dots, n.$$

Hypothèses :

$$L(\epsilon_i) = \mathcal{N}(0, \sigma^2) \quad \text{et} \quad \text{cov}(\epsilon_i, \epsilon_k) = 0 \quad \forall i, k$$



$$Y_1 = \beta_0 + \beta_1 x_{11} + \cdots + \beta_j x_{1j} + \cdots + \beta_p x_{1p}$$

$$\vdots$$

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_j x_{ij} + \cdots + \beta_p x_{ip}$$

$$\vdots$$

$$Y_n = \beta_0 + \beta_1 x_{n1} + \cdots + \beta_j x_{nj} + \cdots + \beta_p x_{np}$$

Ecriture matricielle du modèle : $Y = X\beta + \epsilon$

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_i \\ \vdots \\ Y_n \end{pmatrix} \quad \epsilon = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_i \\ \vdots \\ \epsilon_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_j \\ \vdots \\ \beta_p \end{pmatrix}$$



$$X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1j} & \dots & x_{1p} \\ 1 & x_{i1} & \dots & x_{ij} & \dots & x_{ip} \\ 1 & x_{n1} & \dots & x_{nj} & \dots & x_{np} \end{pmatrix}$$

Ecriture matricielle des hypothèses sur les résidus :

$$Y = X\beta + \epsilon$$

$$\epsilon = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_i \\ \vdots \\ \epsilon_n \end{pmatrix} \quad E(\epsilon) = \begin{pmatrix} E(\epsilon_1) \\ \vdots \\ E(\epsilon_i) \\ \vdots \\ E(\epsilon_n) \end{pmatrix} = 0$$



$$V(\varepsilon) = \begin{pmatrix} V(\varepsilon_1) \\ \text{cov}(\varepsilon_i, \varepsilon_1) & V(\varepsilon_i) \\ & & V(\varepsilon_n) \end{pmatrix} = \sigma^2 I_d$$



Estimation



$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \cdots + \hat{\beta}_j x_{ij} + \cdots + \hat{\beta}_p x_{ip} + e_i$$

Écriture matricielle : $Y = X\hat{\beta} + E$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_j \\ \vdots \\ \hat{\beta}_p \end{pmatrix} \quad E = \begin{pmatrix} e_1 \\ \vdots \\ e_i \\ \vdots \\ e_n \end{pmatrix}$$

Trouver $\hat{\beta}$ tel que : $\sum e_i^2$ minimum (MMC)

solution:

Equations normales : $X'X\hat{\beta} = X'Y$

1. $n = p+1$ et X inversible: $\hat{\beta} = X^{-1}Y$ et $E = 0$
2. $n > p+1$ et $X'X$ inversible: $\hat{\beta} = (X'X)^{-1}X'Y$
3. $n < p+1$: infinité de solution avec $E = 0$



Propriétés de l'estimateur $\hat{\beta} = (X'X)^{-1}X'Y$:

$$E(\hat{\beta}) = \beta \quad E(\epsilon) = 0$$

$$V(\hat{\beta}) = (X'X)^{-1}\sigma^2 \quad V(\epsilon) = \sigma^2 I_d$$



$$V(\hat{\beta}) = \begin{pmatrix} V(\hat{\beta}_1) \\ \text{cov}(\hat{\beta}_j, \hat{\beta}_1) & V(\hat{\beta}_j) \\ & & V(\hat{\beta}_p) \end{pmatrix}$$
$$\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}_1 + \cdots + \hat{\beta}_j \bar{x}_j + \cdots + \hat{\beta}_p \bar{x}_p$$

Cas des données centrées :

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \cdots + \hat{\beta}_j x_{ij} + \cdots + \hat{\beta}_p x_{ip} + e_i$$

$$y_i - \bar{y} = \hat{\beta}_1 (x_{i1} - \bar{x}_1) + \cdots + \hat{\beta}_j (x_{ij} - \bar{x}_j) + \cdots + \hat{\beta}_p (x_{ip} - \bar{x}_p) + e_i$$



Ecart au modèle : Résidus



$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \cdots + \hat{\beta}_j x_{ij} + \cdots + \hat{\beta}_p x_{ip}$$

$$e_i = y_i - \hat{y}_i \quad \hat{Y} = X\hat{\beta} \quad E = Y - \hat{Y}$$

Propriétés:

$$X'E = 0 \text{ (p+1 équations)} \quad \sum_i e_i = 0 \quad \sum_i e_i x_{ij} = 0 \quad \forall j$$

$$\implies \text{cov}[E, x_j] = 0$$

$$\hat{\sigma}^2 = \frac{\sum_i (y_i - \hat{y}_i)^2}{n - p - 1} = \frac{\sum_i (e_i)^2}{n - p - 1}$$

$$E[\hat{\sigma}^2] = \sigma^2$$

Ecart au modèle : Résidus

Exemple numérique $\hat{y}_i = 0 + 1x_{i1} + 2x_{i2}$

	x_1	x_2	y	\hat{y}	e
1	-1	-1	-4	-3	-1
2	-1	1	2	1	1
3	1	-1	0	-1	1
4	1	1	2	3	-1

$$\sum e_i^2 = 4$$

$$\hat{\sigma}^2 = \frac{4}{1}$$

$$\hat{\sigma} = 2$$



Décomposition de la variabilité



$$y_i - \bar{y} = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i)$$

$$\sum_i (y_i - \bar{y})^2 = \sum_i (\hat{y}_i - \bar{y})^2 + \sum_i (y_i - \hat{y}_i)^2$$

Variabilité	Totale	Due au modèle	Résiduelle
Notation	SCE_T	SCE_M	SCE_R
Degré de lib	$n - 1$	p	$n - p - 1$



avec:

$$\hat{y}_i - \bar{y} = \sum_{j=1}^{j=p} \hat{\beta}_j (x_{ij} - \bar{x}_j)$$

Décomposition de la variabilité



Exemple numérique

variabilité	SCE	ddl
totale	$SCE_T = 24$	$n - 1 = 3$
due au modèle	$SCE_M = 20$	$p = 2$
résiduelle	$SCE_R = 4$	$n - p - 1 = 1$

Coéfficient de détermination :

$$\frac{SCE_M}{SCE_T} = R^2$$

$$R^2 = [r(y, \hat{y})]^2$$



Conslusion

En conclusion, la régression linéaire multiple est un modèle efficace pour capturer les relations entre une variable cible et plusieurs variables explicatives. À travers l'estimation des coefficients, l'analyse des résidus et la décomposition de la variabilité, on peut évaluer et affiner la pertinence du modèle. Ce modèle reste donc un outil incontournable pour l'analyse de données et la prise de décisions informées.



Merci Pour votre attention