

Problème – Minimisation de l'erreur de prédiction

Dans ce problème, on s'intéresse à l'erreur de prédiction dans le modèle linéaire. On travaille ici avec un modèle linéaire, non supposé gaussien, et on note σ^2 la variance (constante) du bruit.

Première partie : cas à une variable explicative. On cherche à apprendre le modèle linéaire suivant pour prédire le réel Y_i en fonction d'une variable explicative réelle z_i :

$$\forall 1 \leq i \leq n, \quad Y_i = \beta_0 + \beta_1 z_i + \varepsilon_i.$$

On notera \bar{y} et \bar{z} les moyennes empiriques de $z = (z_1, \dots, z_n)^T$ et de $Y = (Y_1, \dots, Y_n)^T$ (attention, elles dépendent de n). On suppose que le modèle est identifiable et on note $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$ l'estimateur des moindres carrés de β .

Sous le même modèle, on observe une nouvelle valeur z_{n+1} de la variable explicative et on cherche à prédire la variable réponse Y_{n+1} avec l'estimateur

$$\hat{y}_{n+1} := \hat{\beta}_0 + \hat{\beta}_1 z_{n+1} = x_{n+1}^T \hat{\beta}, \quad \text{où } x_{n+1} := \begin{bmatrix} 1 \\ z_{n+1} \end{bmatrix}.$$

L'erreur de prédiction est définie par :

$$\text{err}(z_{n+1}) := \mathbb{E} [(Y_{n+1} - \hat{y}_{n+1})^2].$$

Notons que dans cette espérance, l'aléa vient du bruit ε_{n+1} dans Y_{n+1} ainsi que de l'aléa dans le $\hat{\beta}$.

1. Redémontrer que si S est un vecteur aléatoire de \mathbb{R}^m de matrice de covariance (finie) C , alors pour tout vecteur $u \in \mathbb{R}^m$, $\text{Var}(u^T S) = u^T C u$. *Solution.* On peut développer pour voir que

$$\begin{aligned} \text{Var}(u^T S) &= \text{Var} \left(\sum_{i=1}^m u_i S_i \right) \\ &= \sum_{1 \leq i, j \leq m} \text{Cov}(u_i S_i, u_j S_j) \\ &= \sum_{1 \leq i, j \leq m} u_i u_j C_{i,j} = u^T C u. \end{aligned}$$

2. Que vaut $\mathbb{E}[Y_{n+1} - \hat{y}_{n+1}]$? *Solution.* Elle est égale à $\mathbb{E} [x_{n+1}^T (\beta - \hat{\beta}) + \varepsilon_{n+1}] = x_{n+1}^T \mathbf{0} + 0 = 0$ par nullité du biais de $\hat{\beta}$ et d'après les hypothèses sur les ε_i .
3. Ecrire la matrice X du plan d'expérience dans ce modèle (on mettra l'intercept dans la première colonne), et montrer que

$$(X^T X)^{-1} = \frac{1}{nv(z)} \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n z_i^2 & -\bar{z} \\ -\bar{z} & 1 \end{bmatrix},$$

avec $v(z)$ la variance empirique de z . *Solution.* Dans notre cas on a $X = [\mathbf{1}_n \mid z] \in \mathbb{R}^{n \times 2}$, et $X^T X = n \begin{bmatrix} 1 & \bar{z} \\ \bar{z} & \frac{1}{n} \sum_{i=1}^n z_i^2 \end{bmatrix}$ de déterminant $nv(z)$. Il vient donc $(X^T X)^{-1} = \frac{1}{nv(z)} \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n z_i^2 & -\bar{z} \\ -\bar{z} & 1 \end{bmatrix}$.

4. En utilisant les résultats des questions 1, 2 et 3, montrer que

$$\text{err}(z_{n+1}) = \sigma^2 \left(1 + \frac{1}{n} + \frac{(z_{n+1} - \bar{z})^2}{\sum_{i=1}^n (z_i - \bar{z})^2} \right).$$

Solution. Comme la variable dont on prend le moment d'ordre deux est centré d'après la question 2, on calcule en fait une variance.

$$\begin{aligned} \text{err}(z_{n+1}) &= \text{Var}(\varepsilon_{n+1} + x_{n+1}^T(\beta - \hat{\beta})) \\ &= \sigma^2 + x_{n+1}^T(\sigma^2(X^T X)^{-1})x_{n+1}. \end{aligned}$$

D'après la question 3 on a $(X^T X)^{-1} = \frac{1}{nv(z)} \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n z_i^2 & -\bar{z} \\ -\bar{z} & 1 \end{bmatrix}$ avec $v(z)$ variance empirique de z . Cela donne

$$\begin{aligned} \text{err}(z_{n+1}) &= \sigma^2 + x_{n+1}^T(\sigma^2(X^T X)^{-1})x_{n+1} \\ &= \sigma^2 + \frac{\sigma^2}{nv(z)} \left(\frac{1}{n} \sum_{i=1}^n z_i^2 - 2\bar{z}z_{n+1} + z_{n+1}^2 \right) \\ &= \sigma^2 + \frac{\sigma^2}{nv(z)} (v(z) + \bar{z}^2 - 2\bar{z}z_{n+1} + z_{n+1}^2) \\ &= \sigma^2 \left(1 + \frac{1}{n} + \frac{(z_{n+1} - \bar{z})^2}{nv(z)} \right). \end{aligned}$$

5. Pour quelle(s) valeur(s) de z_{n+1} l'erreur de prédiction est-elle minimale ? Interpréter ce résultat. *Solution.* L'erreur est minimale en $z_{n+1} = \bar{z}$ et vaut $\sigma^2 + \sigma^2/n$. On est toujours meilleur pour prédire lorsqu'on observe le barycentre de ce qu'on a déjà observé. En effet, intuitivement, c'est le point le plus typique.

6. Quelle est la limite de l'erreur de prédiction minimale lorsque $n \rightarrow \infty$? Interpréter cette valeur limite. *Solution.* Cette erreur tend vers σ^2 lorsque n est grand, c'est l'erreur minimale de prédiction qu'on puisse faire, car la variable Y_{n+1} est bruitée avec un bruit indépendant et de variance σ^2 : on ne peut pas faire mieux.

Deuxième partie : cas général. Nous considérons cette fois-ci le cas général à plusieurs covariables. Ici, pour tout $1 \leq i \leq n$,

$$Y_i = \beta_0 + \beta_1 z_{1,i} + \dots + \beta_p z_{p,i} + \varepsilon_i.$$

On note

$$Z := \left[\begin{array}{c|c|c} z_1 & \dots & z_p \end{array} \right] \in \mathbb{R}^{n \times p},$$

avec pour tout $1 \leq \ell \leq p$, $z_\ell := [z_{\ell,1} \dots z_{\ell,n}]^T \in \mathbb{R}^n$. On note \bar{z} le vecteur des moyennes empiriques $\bar{z} := [\bar{z}_1 \dots \bar{z}_p]^T \in \mathbb{R}^p$ (dont les entrées dépendent de n).

7. En notant $\beta = [\beta_0 \ \beta_1 \ \dots \ \beta_p]^T$, écrire le modèle matriciellement sous la forme $Y = X\beta + \varepsilon$ en précisant X . *Solution.* On l'a déjà fait, X s'écrit

$$X = \left[\begin{array}{c|c|c|c} \mathbf{1}_n & z_1 & \dots & z_p \end{array} \right] \in \mathbb{R}^{n \times (p+1)},$$

8. Ecrire la matrice $X^T X$ sous forme de 4 blocs faisant intervenir Z , \bar{z} et n . [Solution.](#)
Déjà fait en séance d'exercices. On a que $X^T X = \begin{bmatrix} n & n\bar{z}^T \\ n\bar{z} & Z^T Z \end{bmatrix}$

Dans toute la suite, on suppose que le modèle est identifiable.

9. On donne la formule d'inversion matricielle par blocs : soit A une matrice inversible s'écrivant par blocs $A = \begin{bmatrix} T & U \\ V & W \end{bmatrix}$ avec T inversible. Alors $Q = W - VT^{-1}U$ est inversible et l'inverse de A est :

$$A^{-1} = \begin{bmatrix} T^{-1} + T^{-1}UQ^{-1}VT^{-1} & -T^{-1}UQ^{-1} \\ -Q^{-1}VT^{-1} & Q^{-1} \end{bmatrix}.$$

On note $\Gamma := \frac{1}{n}Z^T Z - \bar{z}\bar{z}^T$. Ecrire la matrice $(X^T X)^{-1}$ sous la forme d'une matrice par blocs en fonction de n , \bar{z} et Γ^{-1} . [Solution.](#) *La formule donne directement $Q = n\Gamma$ et*

$$(X^T X)^{-1} = \frac{1}{n} \begin{bmatrix} 1 + \bar{z}^T \Gamma^{-1} \bar{z} & -(\Gamma^{-1} \bar{z})^T \\ -\Gamma^{-1} \bar{z} & \Gamma^{-1} \end{bmatrix}.$$

10. Comme dans la première partie, on observe désormais un vecteur

$$x_{n+1} = (1, z_{n+1}) = (1, z_{1,n+1}, \dots, z_{\ell,n+1})$$

et l'on cherche à prédire la variable réponse Y_{n+1} avec l'estimateur $\hat{y}_{n+1} := x_{n+1}^T \hat{\beta}$. L'erreur de prédiction est définie comme précédemment.

Exprimer $\text{err}(z_{n+1})$ en fonction de n , de σ^2 , du vecteur $(z_{n+1} - \bar{z})$ et de la matrice Γ^{-1} . [Solution.](#) *On reprend les calculs de la question 3., et on a encore*

$$\begin{aligned} \text{err}(z_{n+1}) &= \sigma^2 + x_{n+1}^T (\sigma^2 (X^T X)^{-1}) x_{n+1} \\ &= \sigma^2 + \frac{\sigma^2}{n} (1 + z^T \Gamma^{-1} z - 2z_{n+1}^T \Gamma^{-1} \bar{z} + z_{n+1}^T \Gamma^{-1} z_{n+1}) \\ &= \sigma^2 \left(1 + \frac{1}{n} + (z_{n+1} - \bar{z})^T \Gamma^{-1} (z_{n+1} - \bar{z}) \right). \end{aligned}$$

11. On admet dans cette question que Γ est symétrique définie positive. Pour quelle(s) valeur(s) de z_{n+1} l'erreur de prédiction est-elle minimale ? [Solution.](#) *Ben du coup c'est évident, il faut dire que Γ^{-1} est encore symétrique définie positive, et le minimum est donc atteint en $z_{n+1} = \bar{z}$. Le reste du temps, c'est au-dessus !*
12. (★) Montrer que $\Gamma = \frac{1}{n}Z^T Z - \bar{z}\bar{z}^T$ est bien symétrique définie positive. On pourra chercher à l'écrire sous la forme $\frac{1}{n}WW^T$ avec W une matrice bien choisie. [Solution.](#) *Symétrie évidente. Ensuite, l'idée c'est de voir qu'elle ressemble fortement à une matrice de covariance. C'est même carrément une matrice de covariance empirique. Si on note \tilde{Z} la matrice*

$$\tilde{Z} := \begin{bmatrix} z_1 - \bar{z}_1 \mathbf{1}_n & \dots & z_p - \bar{z}_p \mathbf{1}_n \end{bmatrix} \in \mathbb{R}^{n \times p},$$

on a que $\Gamma = \frac{1}{n}\tilde{Z}^T \tilde{Z}$. On a donc la positivité qui en découle. On peut ensuite évoquer le résultat d'inversion par blocs qui implique que Γ est inversible, ou remarquer que \tilde{Z} est

injective. En effet, si $u \in \text{Ker} \tilde{Z}$, alors $u_1 z_1 + \dots + u_p z_p = (\sum_{\ell=1}^p u_\ell \bar{z}_\ell) \mathbf{1}_n$, et comme X est de rang plein par hypothèse, $u = 0$.