

Analyse des données prostate: les moindres carrés ordinaires

julien.chiquet@genopole.cnrs.fr

Module MPR – option modélisation, 12 novembre 2009

1 Le modèle linéaire

Le modèle linéaire est défini par

$$y = \beta_0 + \sum_{i=1}^p X_i \beta_i + \varepsilon,$$

où y est la réponse à expliquer par le vecteur de variables d'entrées $X = (X_1, \dots, X_p)$. Le scalaire β_0 (appelé biais ou *intercept*) et les $(\beta_i)_{i=1, \dots, p}$ sont les paramètres à estimer. On suppose ici que les résidus ε sont gaussiens, c'est-à-dire que $\varepsilon \sim \mathcal{N}(0, \sigma^2)$.

On dispose de n données relatives aux variables y et X . Notons \mathbf{y} le vecteur de taille n contenant les n observations de y , et \mathbf{X} la matrice de taille $n \times (p+1)$ dont la $(i+1)^{\text{e}}$ colonne contient les n observations relatives à la variable X_i ; la première colonne de \mathbf{X} est remplie de 1, ce qui permet de gérer l'intercept. Ainsi, le modèle linéaire s'écrit pour ce n -échantillon

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

où $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^\top$ est le vecteur des paramètres à estimer.

2 L'estimateur des moindres carrés

L'estimateur des moindres carrés ordinaires est défini comme le vecteur minimisant la somme des carrés résiduels :

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{p+1}} \text{RSS}(\boldsymbol{\beta}), \quad \text{RSS}(\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2. \quad (1)$$

Le théorème suivant donne la solution de ce problème et les grandeurs caractéristiques standards associées à un estimateur (biais, variance).

Théorème. Lorsque $\mathbf{X}^\top \mathbf{X}$ est inversible¹, la solution $\hat{\boldsymbol{\beta}}$ des équations normales,

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad (2)$$

¹c'est-à-dire, lorsque $n > p + 1$: moralement, cela signifie qu'on a suffisamment de données pour estimer nos $p + 1$ paramètres.

est unique et est appelée estimateur de Gauss-Markov. C'est un estimateur sans biais de β et de matrice de variance-covariance

$$\text{Var}(\hat{\beta}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}. \quad (3)$$

Démonstration. La fonction RSS est convexe : son minimum est atteint lorsque son gradient $\nabla_{\beta} \text{RSS}(\beta)$ s'annule. Ainsi

$$\nabla_{\beta} \text{RSS}(\beta) = \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta) = 0 \Leftrightarrow \mathbf{X}^\top \mathbf{X}\beta = \mathbf{X}^\top \mathbf{y},$$

d'où la solution (2) pour $\mathbf{X}^\top \mathbf{X}$ inversible.

Comme² $\mathbb{E}(\mathbf{y}) = \mathbb{E}(\mathbf{X}\beta + \varepsilon) = \mathbf{X}\beta$, on a

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}(\mathbf{y}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\beta = \beta,$$

et $\hat{\beta}$ est donc sans biais. Concernant sa matrice de variance-covariance, on a³

$$\text{Var}(\hat{\beta}) = \text{Var}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Var}(\mathbf{y}) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}.$$

Or, $\text{Var}(\mathbf{y}) = \text{Var}(\varepsilon) = \sigma^2 \mathbf{I}_n$, avec \mathbf{I}_n la matrice identité de taille $n \times n$; d'où l'expression (3). \square

2.1 Risque quadratique, erreur de prédiction

Définition (Risque quadratique). *On appelle risque quadratique (Mean Squared Error) de l'estimateur $\hat{\theta}$ du paramètre θ la grandeur définie par :*

$$\text{MSE}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta)^2] = \text{biais}^2(\hat{\theta}) + \text{Var}(\hat{\theta}). \quad (4)$$

Évidemment, le risque quadratique de $\hat{\beta}$ se ramène simplement à sa variance puisqu'il est sans biais.

Définition (Erreur de prédiction). *On appelle erreur de prédiction du modèle de régression*

$$y = f(X) + \varepsilon$$

la grandeur définie par⁴

$$\text{EPE}(f) = \mathbb{E}[(y - f(X))^2]. \quad (5)$$

Dans le cas du modèle linéaire, l'estimateur de Gauss-Markov prévoit pour \mathbf{y} le vecteur de valeurs $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$. Ainsi, nous avons

$$\text{EPE}(X\hat{\beta}) = \mathbb{E}[(y - X\hat{\beta})^2] = \mathbb{E}(y - X\hat{\beta})^2 + \text{Var}(y - X\hat{\beta}) = \sigma^2 + \text{MSE}(X\hat{\beta}).$$

²Rappel : $\mathbb{E}(AX) = A\mathbb{E}(X)$ pour A une matrice déterministe et X un vecteur aléatoire

³Rappel : $\text{Var}(AX) = A\text{Var}(X)A^\top$ pour A une matrice déterministe et X un vecteur aléatoire

⁴EPE pour *Expected Prediction Error*

Ensemble de test, ensemble d'apprentissage. Pour évaluer les performances d'un modèle, il est classique lorsque l'on dispose de suffisamment de données de les séparer en un ensemble d'*apprentissage*, qui permet de calculer les paramètres du modèle (ici $\hat{\beta}$), et un ensemble de *test*, qui permet d'évaluer la qualité du modèle. Sans faire cette procédure, on peut montrer que l'on sous-estime l'erreur de prédiction.

Estimation du risque quadratique. Le risque quadratique dépend de la variance de $\hat{\beta}$: il faudrait connaître la valeur σ^2 pour le calculer. Si l'on note $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$, les valeurs prédites par le modèle pour \mathbf{y} , on peut estimer le paramètre σ^2 par la variance empirique corrigée des résidus empiriques $\mathbf{y} - \hat{\mathbf{y}}$, c'est-à-dire

$$\hat{\sigma}^2 = \frac{1}{n-p-1} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2. \quad (6)$$

Ce calcul se fait sur les données de l'ensemble d'apprentissage, puisque ce sont celles qui ont permis de déterminer $\hat{\beta}$.

Estimation de l'erreur de prédiction. L'erreur de prédiction s'estime par exemple par

$$\widehat{\text{EPE}} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \left(\hat{\mathbf{y}}_{\text{test}} - \mathbf{X}_{\text{test}} \hat{\beta} \right)^2. \quad (7)$$

Nous verrons plus tard d'autres manières d'estimer cette erreur, en particulier par validation croisée.

2.2 Le Z -score : une première méthode de sélection de variables

Sous l'hypothèse de normalité des résidus ε , le vecteur $\hat{\beta}$ suit évidemment une loi normale multivariée donc on connaît l'espérance et la variance. Ainsi,

$$\hat{\beta} \sim \mathcal{N}(0, \sigma^2 \mathbf{X}^\top \mathbf{X}).$$

De même,

$$(n-p-1)\hat{\sigma}^2 \sim \sigma^2 \chi^2(n-p-1),$$

puisque $\hat{\sigma}^2$ est définie comme la somme du carré des résidus (qui sont gaussiens). On peut donc construire le test d'hypothèse suivant sur la nullité des composants de $\hat{\beta}$:

$$\begin{cases} H_0 : \hat{\beta}_i = 0, \\ H_1 : \hat{\beta}_i \neq 0. \end{cases}$$

Le Z -score, défini par $z_i = \hat{\beta}_i / \hat{\sigma} \sqrt{\text{diag}(\mathbf{X}^\top \mathbf{X})}$, suit une loi de Student à $n-p-1$ degrés de liberté sous H_0 , puisque défini comme le rapport d'une loi normale et d'un χ^2 .

Ainsi, on rejette H_0 au niveau de confiance α lorsque $|z_i| > t_{1-\alpha/2; n-p-1}$, le fractile d'ordre d'une loi de Student à $n-p-1$ degrés de liberté.

3 Commande R commentée

3.1 Préparation des données

On charge les données prostate, qui contiennent trois variables : un vecteur `x`, un vecteur `y` et un vecteur de booléens `set` indiquant les données de l'ensemble d'apprentissage.

```
> load("prostate.rda")
```

La manipulations des matrices (objet `matrix`) est plus simple que celle des tableaux de données (objet `data.frame`). Transformons donc `x` en matrice à l'aide de la fonction `as.matrix` :

```
> x <- as.matrix(x)
```

Pour retrouver les résultats du livre [ESLII], page 50, nous devons commencer par centrer et réduire les données⁵ :

```
> x <- scale(x)
```

Ensuite, nous construisons les données test (vecteurs `x.test` et `y.test`) et les données d'apprentissage (vecteurs `x` et `y`).

```
> test <- !set
> n <- sum(set)
> n.test <- sum(test)
> x.test <- x[test, ]
> y.test <- y[test]
> x <- x[set, ]
> y <- y[set]
```

Concaténons un vecteur $\mathbf{1}_n$ aux matrices \mathbf{X} et \mathbf{X}_{test} afin d'inclure l'intercept dans la procédure d'estimation des paramètres :

```
> p <- ncol(x)
> x <- cbind(rep(1, n), x)
> x.test <- cbind(rep(1, n.test), x.test)
```

3.2 Estimation

Calculons l'estimateur du vecteur des paramètres (notez l'utilisation de `solve` : pour inverser une matrice A , on résout $AX = I$) :

```
> beta.ols <- solve(t(x) %*% x, t(x) %*% y)
```

⁵Théoriquement, on devrait faire l'opération de centrage/réduction *après* avoir construit les ensembles de test et d'apprentissage – tenons-nous en à l'approche du livre, afin de retrouver exactement les mêmes valeurs

Vient ensuite l'estimation de σ^2 , de la variance de $\hat{\beta}$ et de l'erreur de prédiction.

```
> sigma.hat <- sqrt(1/(n - p - 1) * sum((y - x %*% beta.ols)^2))
> var.beta.ols <- solve(t(x) %*% x) * sigma.hat^2
> std.error <- sqrt(diag(var.beta.ols))
> epe <- 1/n.test * sum((y.test - x.test %*% beta.ols)^2)
```

On teste au niveau $\alpha = 5\%$ la nullité des $\hat{\beta}_i$ et on stocke dans un vecteur `rejet` les résultats (1 en cas de rejet de H_0).

```
> Z.score <- beta.ols/std.error
> rejet <- rep(FALSE, p + 1)
> alpha <- 0.05
> rejet[which(abs(Z.score) > qt(1 - alpha/2, df = n + p + 1))] <- TRUE
```

3.3 Résultats

On stocke tout les résultats dans une matrice `results` en y nommant les colonnes de manière appropriée.

```
> results <- cbind(beta.ols, std.error, Z.score, rejet)
> colnames(results) <- c("beta.hat", "std.error", "Z.score", "H1 à 5%")
```

Voici pour appeler l'affichage des résultats :

```
> cat("\n Modèle linéaire sur données prostate\n\n")
> print(results)
> cat("\n Valeur estimée de l'erreur de prédiction:", epe, "\n")
```

Ce qui donne à l'écran :

```
Modèle linéaire sur données prostate

      beta.hat  std.error  Z.score H1 à 5%
      2.46493292 0.08931498 27.5982031    1
lcavol  0.67952814 0.12662903  5.3662905    1
lweight 0.26305307 0.09562821  2.7507894    1
age     -0.14146483 0.10134245 -1.3959090    0
lbph    0.21014656 0.10221904  2.0558456    1
svi     0.30520060 0.12360027  2.4692552    1
lcp     -0.28849277 0.15452934 -1.8669126    0
gleason -0.02130504 0.14524723 -0.1466812    0
pgg45   0.26695576 0.15361357  1.7378397    0
```

Valeur estimée de l'erreur de prédiction: 0.521274

Références

[ESLII] Hastie, T., Tibshirani, R. and Friedman, J., *The Elements of Statistical Learning*, second edition, Springer, 2008.