

Chapitre 5

Modèle linéaire mixte

On appelle modèle mixte un modèle statistique dans lequel on considère à la fois des facteurs à effets fixes (qui interviennent sur la moyenne dans différents groupes du modèle) et des facteurs à effets aléatoires (qui interviennent sur la variance du modèle). Un modèle est dit mixte lorsqu'il y a au moins un facteur de chaque nature. Dans le cadre de ce cours, nous ne considérons que des modèles linéaires gaussiens mixtes, mais la notion de modèle mixte se rencontre également dans d'autres contextes, notamment dans le modèle linéaire généralisé.

5.1 Écriture du modèle

Modèle

Un modèle linéaire gaussien mixte à n observations s'écrit sous la forme matricielle suivante :

$$\begin{aligned}\mathbf{Y} &= \mathbf{X}\beta + \sum_{k=1}^K \mathbf{Z}_k \mathbf{A}_k + U \\ &= \mathbf{X}\beta + \mathbf{Z}\mathbf{A} + U\end{aligned}$$

où :

- \mathbf{Y} est le vecteur aléatoire réponse de \mathbb{R}^n .
- \mathbf{X} est la matrice $n \times p$ relative aux effets fixes du modèle, où p est le nombre total d'effets fixes pris en compte dans le modèle.
- β est le vecteur des p effets fixes $\beta_j, j = 1, \dots, p$ à estimer.
- \mathbf{Z}_k est la matrice des indicatrices (disposées en colonnes) des niveaux du k ème facteur à effets aléatoires ($k = 1, \dots, K$). On note q_k le nombre de niveaux de ce facteur. \mathbf{Z}_k est alors de dimension $n \times q_k$.
- On note A_{kl} la v.a.r. associée au l ème niveau du k ème facteur à effets aléatoires avec $l = 1, \dots, q_k$. Pour tout l lié au facteur k , on suppose $A_{kl} \sim \mathcal{N}(0, \sigma_k^2)$.
- Pour un facteur k donnée, on note $\mathbf{A}_k = (A_{k1}, \dots, A_{kq_k})'$ le vecteur colonne des A_{kl} . On suppose que $\mathbf{A}_k \sim \mathcal{N}(0, \sigma_k^2 \mathbf{I}_{q_k})$.
- Enfin, U est le vecteur aléatoire des erreurs du modèle qui vérifie $U \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$.

Le modèle peut alors être écrit sous la forme :

$$\begin{array}{c}
 \begin{array}{c} \uparrow n \\ \downarrow \end{array} \left(\begin{array}{c} Y \\ \leftarrow 1 \end{array} \right) = \left(\begin{array}{c} X \\ \leftarrow p \end{array} \right) \left(\begin{array}{c} \beta \\ \leftarrow 1 \end{array} \right) + \left(\begin{array}{c} Z_1 \dots Z_K \\ \leftarrow q_1 \dots q_K \end{array} \right) \left(\begin{array}{c} A_1 \\ \vdots \\ A_K \\ \leftarrow 1 \end{array} \right) + \left(\begin{array}{c} U \\ \leftarrow 1 \end{array} \right)
 \end{array}$$

réponse *effets fixes* *vecteur des effets fixes* *Z_k est l'indicateur des q_k niveaux du k ème facteur α* *effets aléatoires* *vecteur des effets aléatoires* *Bruit*

Moments du modèle

Il est évident de montrer que $\mathbb{E}(\mathbf{Y}) = \mathbf{X}\beta$ d'après les hypothèses sur le modèle. On note aussi \mathbf{V} la variance de \mathbf{Y} qui se calcule par :

$$\begin{aligned}
 \mathbf{V} &= \text{Var}(\mathbf{Y}) \\
 &= \text{Var}(\mathbf{Z}\mathbf{A}) + \text{Var}(\mathbf{U}) \\
 &= \sum_{k=1}^K (\sigma_k^2 \mathbf{Z}_k \mathbf{Z}_k') + \sigma^2 \mathbf{I}_n \\
 &= \mathbf{Z}_k \mathbf{G} \mathbf{Z}_k' + \sigma^2 \mathbf{I}_n
 \end{aligned}$$

où $\mathbf{G} = \text{diag}(\sigma_1^2 \mathbf{I}_{q_1}, \dots, \sigma_K^2 \mathbf{I}_{q_K})$. On obtient alors $\mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\beta, \mathbf{V})$.

Les composantes de Y ne sont ainsi pas indépendantes au sein de chaque niveau l d'un facteur aléatoire k donné. Ceci est évident si on observe un exemple de ce à quoi peut ressembler \mathbf{Z} :

$$\mathbf{Z} = \begin{pmatrix} 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 \end{pmatrix} \quad \begin{array}{l} \uparrow n=7 \\ \downarrow \end{array}$$

\mathbf{Z}_1 \mathbf{Z}_2
avec $q_1=2$ *avec $q_2=3$*

Dans cet exemple, le vecteur $\mathbf{Z}\mathbf{A}$ aura la forme $(\xi_{11} + \xi_{21}, \xi_{11} + \xi_{21}, \xi_{11} + \xi_{22}, \xi_{12} + \xi_{22}, \xi_{12} + \xi_{22}, \xi_{12} + \xi_{23}, \xi_{12} + \xi_{23})'$ où les $\xi_{kl} \sim \mathcal{N}(0, \sigma_k^2)$ ce qui induit les dépendances intra-niveau des Y .

5.2 Estimation des β

L'expression que l'on obtient dans le cas général pour $\hat{\beta}$ fait intervenir l'estimation de la matrice des variances-covariances \mathbf{V} de \mathbf{Y} . Cette expression obtenue est fournie par la méthode des moindres carrés généralisés notée $GLSE(\beta)$ (pour Generalized Least Squares Estimator) :

$$\hat{\beta} = \arg \min_{\beta} (\mathbf{Y} - \mathbf{X}\hat{\beta})' \hat{\mathbf{V}}^{-1} (\mathbf{Y} - \mathbf{X}\hat{\beta})$$

où $\hat{\mathbf{V}} = \sum_{k=1}^K \hat{\sigma}_k^2 \mathbf{Z}_k \mathbf{Z}_k' + \hat{\sigma}^2 \mathbf{I}_n$ et les $\hat{\sigma}_k^2$ et $\hat{\sigma}^2$ sont les composantes de variances. On a alors :

$$\hat{\beta} = GLSE(\beta) = (\mathbf{X}' \hat{\mathbf{V}}^{-1} \mathbf{X})^{-1} \mathbf{X}' \hat{\mathbf{V}}^{-1} \mathbf{Y}$$

et il est nécessaire d'estimer les composantes de covariance, ce qui se fait typiquement par maximum de vraisemblance.

On remarquera que dans le cas équilibré où tous les q_k ont la même valeur, on a plus simplement :

$$\hat{\beta} = OLSE(\beta) = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}$$

où $OLSE$ signifie *Ordinary Least Squares Estimator*.

5.3 Estimation de \mathbf{V}

Pour estimer \mathbf{V} par maximum de vraisemblance, on note d'abord $\Psi = (\hat{\sigma}_1^2, \dots, \hat{\sigma}_K^2, \hat{\sigma}^2)$ les paramètres à estimer dont dépend \mathbf{V} . La log-vraisemblance du modèle mixte gaussien s'écrit :

$$l(y, \beta, \mathbf{V}(\Psi)) = -\frac{1}{2} \log(\det(\mathbf{V}(\Psi))) - \frac{1}{2} (\mathbf{Y} - \mathbf{X}\beta)' (\mathbf{V}(\Psi))^{-1} (\mathbf{Y} - \mathbf{X}\beta)$$

On en déduit le système de p équations :

$$\frac{\partial l}{\partial \beta} = \mathbf{X}' \mathbf{V}^{-1} - \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \beta$$

dont découlent les équations normales pour $\hat{\beta}$.

On remarque ensuite que :

$$\frac{\partial \mathbf{V}}{\partial \sigma_k^2} = \mathbf{Z}_k \mathbf{Z}_k'$$

On déduit alors que pour chaque σ_k^2 :

$$\frac{\partial l}{\partial \sigma_k^2} = -\frac{1}{2} \text{tr}(\mathbf{V}(\Psi) \mathbf{Z}_k \mathbf{Z}_k') + \frac{1}{2} (\mathbf{Y} - \mathbf{X}\beta)' (\mathbf{V}(\Psi))^{-1} \mathbf{Z}_k \mathbf{Z}_k' (\mathbf{V}(\Psi))^{-1} (\mathbf{Y} - \mathbf{X}\beta)$$

On obtient ainsi un système de $K + 1 + p$ équations non linéaires à $K + 1 + p$ inconnues que l'on résoud par une méthode numérique itérative. Ces procédures numériques fournissent en plus, à la convergence, la matrice des variances-covariances asymptotiques des estimateurs.

5.4 Tests de significativité des facteurs

Ces tests sont standards dans le cas équilibré (tous les q_k sont égaux), mais deviennent assez problématiques dans le cas déséquilibré. Dans le cas équilibré le test de Fisher sur les variances est en effet valable (comme dans ANOVA sous-section [4.3](#) ou ANCOVA sous-section [4.6](#)). Il n'y a cependant pas de test exact, ni même de test asymptotique, qui permette de tester les effets, que ce soient les effets fixes ou les effets aléatoires, dans un modèle mixte avec un plan déséquilibré. Il existe seulement des tests approchés (dont on ne contrôle pas réellement le niveau, et encore moins la puissance)

Chapitre 6

Ouvertures

6.1 Régression logistique

Modèle

On se pose maintenant dans le cas où une variable qualitative Y a 2 modalités : 1 ou bien 0. Les modèles de régression précédents adaptés à l'explication d'une variable quantitative ne s'appliquent plus directement car le régresseur linéaire usuel $\mathbf{X}\beta$ ne prend pas des valeurs simplement binaires. Si l'on ne connaît que Y , on pourra estimer le paramètre Π de la loi de Bernoulli, $\mathbb{P}(Y = 1) = \Pi$ et $\mathbb{P}(Y = 0) = 1 - \Pi$, en calculant la moyenne empirique de $\mathbb{P}(Y = 1)$. On va cependant s'intéresser ici au cas où Y est lié à n observations en dimension p .

On note $\mathbb{P}(Y = 1|X)$ la loi conditionnelle que Y soit égal à 1 sachant X . On suppose alors que :

$$\ln \frac{\mathbb{P}(Y = 1|X)}{1 - \mathbb{P}(Y = 1|X)} = \beta_0 + \sum_{j=1}^p \beta_j X_j \quad (6.1)$$

où les X_j sont les p composantes de X . Il est intéressant de remarquer que $\ln(\mathbb{P}/(1 - \mathbb{P}(.)))$ est une fonction de $\mathbb{P}(.)$ strictement croissante qui :

- tend vers $-\infty$ quand $\mathbb{P}(.)$ se rapproche de 0,
- vaut 0 pour $\mathbb{P}(.) = 0.5$,
- tend vers $+\infty$ quand $\mathbb{P}(.)$ se rapproche de 1.

On en conclue que Y à plutôt des chances de valoir 0 si $\beta_0 + \sum_{j=1}^p \beta_j X_j$ est négatif, et que Y à plutôt des chances de valoir 1 si $\beta_0 + \sum_{j=1}^p \beta_j X_j$ est positif. Notons aussi que le modèle de régression est dit logistique car la loi de probabilité est modélisée à partir d'une loi logistique. Ce modèle est extrêmement populaire en apprentissage automatique car il se montre performant quand on n'a que deux classes à distinguer, et il passe facilement à l'échelle.

Apprentissage des β_j

Après transformation de l'équation, on obtient :

$$\mathbb{P}(Y = 1|X) = \frac{e^{\beta_0 + \sum_{j=1}^p \beta_j X_j}}{1 + e^{\beta_0 + \sum_{j=1}^p \beta_j X_j}}$$

que l'on notera pour une observation i , $i = 1, \dots, n$:

$$p(y_i = 1 | x_i^1, \dots, x_i^p) = \frac{e^{\beta_0 + \sum_{j=1}^p \beta_j x_i^j}}{1 + e^{\beta_0 + \sum_{j=1}^p \beta_j x_i^j}}.$$

Pour résoudre estimer les β_j à l'aide d'un jeu de n observations y_i, x_i^1, \dots, x_i^p , $i = 1, \dots, n$ et de la méthode du maximum de vraisemblance, on note la contribution à la vraisemblance de l'observation i :

$$(p(y_i = 1 | x_i^1, \dots, x_i^p))^{y_i} \cdot (1 - p(y_i = 1 | x_i^1, \dots, x_i^p))^{1-y_i}$$

qui vaut $p(y_i = 1 | x_i^1, \dots, x_i^p)$ si $y_i = 1$ et qui vaut $p(y_i = 0 | x_i^1, \dots, x_i^p)$ si $y_i = 0$. La vraisemblance des observations s'écrit alors :

$$L(\beta) = \prod_{i=1}^n \left[(p(y_i = 1 | x_i^1, \dots, x_i^p))^{y_i} \cdot (1 - p(y_i = 1 | x_i^1, \dots, x_i^p))^{1-y_i} \right]$$

Les paramètres β_j qui maximisent cette quantité sont les estimateurs du maximum de vraisemblance de la régression logistique. Ils seront estimés typiquement en utilisant une méthode itérative. Pour des raisons numériques la log-vraisemblance, *i.e.* $n^{-1} \log(L)$, sera aussi maximisé plutôt que L .

Prédiction

Une fois les β_j appris, on se réfère à Eq. (6.1) et son interprétation pour prédire le label d'un y_0 en fonction de x_0^j observés. On calculera simplement $\beta_0 + \sum_{j=1}^p \beta_j x_0^j$. Si le signe est positif alors $\hat{y}_0 = 1$ et si le signe est négatif alors $\hat{y}_0 = 0$.

Sélection de modèle

Notons enfin qu'il est possible et même classique de sélectionner un modèle en régression logistique en pénalisant les β_j lors de la maximisation de la Log-vraisemblance, typiquement avec une méthode de type Lasso. On résoudra alors le problème suivant :

$$\hat{\beta} = \arg \max_{\beta} (n^{-1} \log(L(\beta)) - \lambda |\beta|_1)$$

où comme pour la régression linéaire multiple, on sera amené à trouver un λ qui offrira un bon compromis entre pouvoir prédictif et explicabilité du modèle.

6.2 Méthode Partial Least Squares

Intuition

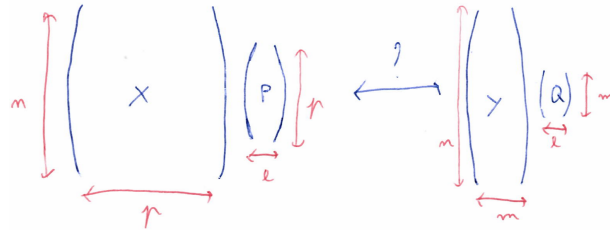
On a vu dans le cours de statistique que l'Analyse en Composantes Principales (ACP) était un outil essentiel pour explorer un ensemble d'observations $X_i = (x_i^1, \dots, x_i^p)$, $i = 1, \dots, n$ regroupées en ligne dans une matrice \mathbf{X} . L'ACP consiste en effet à maximiser la variance des projections des observations X_i , ce qui permet entre autres d'expliquer comment les variables interagissent

entre elles. Plus spécifiquement, le 1er vecteur propre v_1 est celui qui maximise la variance des projections des X_i . En supposant que les X_i sont centrés (et idéalement réduits), cela signifie que

$$v_1 = \arg \max_{v \text{ t.q. } |v|_2=1} \sum_{i=1}^n (X_i v)^2$$

Le 2ème vecteur propre v_2 est choisi suivant le même principe, une fois enlevée l'influence de v_1 dans \mathbf{X} ; et ainsi de suite.

L'idée de la méthode *Partial Least Squares* (PLS) est relativement similaire, mais maintenant on s'intéresse au lien entre \mathbf{X} et une matrice $n \times m$ de réponses \mathbf{Y} . Pour chaque observation X_i de \mathbf{X} , la matrice \mathbf{Y} contient une réponse Y_i en dimension m . Si $m = 1$, on a les mêmes données d'entrée que dans le cadre de la régression linéaire multiple (Section 2.2). L'approche d'analyse est cependant totalement différente : On cherche les transformations linéaires \mathbf{P} et \mathbf{Q} de \mathbf{X} et de \mathbf{Y} (si $m > 1$), respectivement, telles que : La 1ère colonne de \mathbf{P} est celle qui projete les X_i de manière à séparer au mieux les y_i projetés par la première colonne de \mathbf{Q} ; et ainsi de suite. Cette idée est schématisée ci-dessous :



Modèle

La méthode PLS (Partial Least Squares) repose toujours sur une hypothèse de modèle linéaire $\mathbf{Y} = \mathbf{X}\beta + \mathbf{U}$, où \mathbf{U} modélise le bruit. L'approche utilisée pour estimer le lien entre \mathbf{Y} et les variables explicatives de \mathbf{X} est cependant différente de celle du modèle linéaire classique. En particulier, le modèle sur le bruit est totalement différent et va dépendre de la covariance entre des combinaisons linéaires de \mathbf{X} et \mathbf{Y} .

Plus spécifiquement, on suppose :

$$\mathbf{X} = \mathbf{TP}' + \mathbf{E}$$

$$\mathbf{Y} = \mathbf{UQ}' + \mathbf{F}$$

où

- \mathbf{X} est la matrice $n \times p$ de prédicteurs. Elle est supposée centrée/réduite,
- \mathbf{Y} est la matrice $n \times m$ de réponses. Elle est supposée centrée/réduite,
- \mathbf{P} et \mathbf{Q} sont respectivement des matrices $p \times l$ et $m \times l$ de projection. Leurs colonnes sont orthonormés.
- \mathbf{T} et \mathbf{U} sont les projections de \mathbf{X} et de \mathbf{Y} respectivement par \mathbf{P} et \mathbf{Q} . Elles sont de taille $n \times l$.
- \mathbf{E} et \mathbf{F} sont des termes d'erreur de même taille que \mathbf{X} et \mathbf{Y} . Ils sont supposés *i.i.d.* et distribués suivant une loi normale.

Les projections de \mathbf{X} et de \mathbf{Y} dans \mathbf{T} et \mathbf{U} sont aussi toutes deux de même taille $n \times l$ avec $l \leq p$. La PLS consiste alors à calculer les projecteurs \mathbf{P} et \mathbf{Q} qui maximisent la covariance entre \mathbf{T} et \mathbf{U} . On dénote $\bar{\mathbf{T}}_j$ et $\bar{\mathbf{U}}_j$ la moyenne des valeurs des colonnes j de \mathbf{T} et \mathbf{U} . On maximise alors $\sum_{j=1}^l \sum_{i=1}^n (\mathbf{T}_{ij} - \bar{\mathbf{T}}_j)(\mathbf{U}_{ij} - \bar{\mathbf{U}}_j)$.

Estimation

Géométriquement, la régression PLS consiste à calculer une projection des \mathbf{X} sur un hyperplan qui est à la fois une bonne estimation de \mathbf{X} et dont les projections sont de bons prédicteurs des \mathbf{Y} . En vue de la définition d'une stratégie d'optimisation, le problème peut être vu sous la forme plus classique $\mathbf{Y} = \mathbf{X}\hat{\mathbf{B}} + \mathbf{B}_0$. Nous donnons Alg. [1](#) l'algorithme PLS1 qui permet de résoudre le problème pour $m = 1$, c'est à dire \mathbf{Y} est un vecteur colonne. Dans ce cas là, $\hat{\mathbf{B}}$ est un vecteur de taille p dont l'interprétation est similaire aux vecteurs $\hat{\beta}$ du modèle linéaire multiple mais avec un modèle sous-jacent différent.

Alg. 1 Fonction $PLS1(\mathbf{X}, \mathbf{y}, l)$

```

1:  $\mathbf{X}^{(0)} \leftarrow \mathbf{X}$ 
2:  $\mathbf{w}^{(0)} \leftarrow \mathbf{X}'\mathbf{y}/|\mathbf{X}'\mathbf{y}|_2$ 
3: for  $k = 0, \dots, l-1$  do
4:    $\mathbf{t}^{(k)} \leftarrow \mathbf{X}^{(k)}\mathbf{w}^{(k)}$ 
5:    $t_k \leftarrow \mathbf{t}^{(k)'}\mathbf{t}^{(k)}$ 
6:    $\mathbf{t}^{(k)} \leftarrow \mathbf{t}^{(k)}/t_k$ 
7:    $\mathbf{p}^{(k)} \leftarrow \mathbf{X}^{(k)'}\mathbf{t}^{(k)}$ 
8:    $q_k \leftarrow \mathbf{y}'\mathbf{t}^{(k)}$ 
9:   if  $q_k = 0$  then
10:     $l \leftarrow k$  et sort de la boucle for (toute la variabilité est capturée).
11:   end if
12:   if  $k < (l-1)$  then
13:      $\mathbf{X}^{(k+1)} \leftarrow \mathbf{X}^{(k)} - t_k\mathbf{t}^{(k)}\mathbf{p}^{(k)'}$ 
14:      $\mathbf{w}^{(k+1)} \leftarrow \mathbf{X}^{(k+1)'}\mathbf{y}/|\mathbf{X}^{(k+1)'}\mathbf{y}|_2$ 
15:   end if
16: end for
17:  $\mathbf{W}$  est la matrice composée des colonnes  $\mathbf{w}^{(0)}, \mathbf{w}^{(1)}, \dots, \mathbf{w}^{(l-1)}$ .
18:  $\mathbf{P}$  est la matrice composée des colonnes  $\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(l-1)}$ .
19:  $\mathbf{q}$  est le vecteur composé des scalaires  $q_0, q_1, \dots, q_{l-1}$ .
20:  $\mathbf{B} \leftarrow \mathbf{W}(\mathbf{P}'\mathbf{W})^{-1}\mathbf{q}$ 
21:  $\mathbf{B}_0 \leftarrow q_0 - \mathbf{P}^{(0)'}\mathbf{B}$ 
22: return  $\mathbf{B}, \mathbf{B}_0$ 
```

Sparse PLS

On a vu Chapitre [3](#) l'intérêt pratique des méthodes de régularisation, telles que LASSO, qui sélectionnent des modèles parcimonieux (sparse). En plus de bien contraindre le problème de régression, ces modèles sont en effet simples à interpréter, même pour des non-spécialistes de l'analyse de données.

Ce principe peut aussi s'appliquer dans le cas de la PLS, afin de trouver à la fois des changements de bases qui mettent en lien les \mathbf{X} et \mathbf{Y} de manière optimale, et qui permettent de regrouper des blocs de variables ayant une influence similaire lorsque \mathbf{X} et \mathbf{Y} sont mis en lien. La méthode de la *sparse PLS* est alors extrêmement puissante d'un point de vue pratique.

En posant par exemple une pénalisation L_1 sur les éléments de la base \mathbf{W} avec une pondération λ_W , Alg. 1 sera légèrement modifié en remplaçant les lignes 2 et 14 par : $\mathbf{w}^{(new)} \leftarrow \mathbf{X} \mathbf{y}$, puis $\mathbf{w}^{(new)} = \mathbf{w}^{(new)} - \lambda_W \text{sign}(\mathbf{w}^{(new)})$ tout en mettant les différents éléments à 0 si leur signe change (comme pour la régularisation LASSO), puis enfin en normalisant $\mathbf{w}^{(new)} = \mathbf{w}^{(new)} / \|\mathbf{w}^{(new)}\|_2$. Seuls les éléments de $\mathbf{w}^{(new)}$ ayant une réelle influence, en fonction de la *pression* de λ_P , auront alors une valeur non nulle et seront sélectionnés.