

Modèles à un facteur aléatoire

Cours: Joseph Salmon

Scribes: JOLY Julien - Mohamed Sahardid - Anas El Benna

1 Modèles à deux facteurs avec interactions

Modèle :

$$y_{i,j,k} = \mu^* + \alpha_i^* + \beta_j^* + \gamma_{i,j}^* + \varepsilon_{i,j,k}, \text{ pour } i = 1, \dots, I \text{ et } j = 1, \dots, J . \quad (1)$$

Remarque 1.1. Ce modèle est compliqué si les classes sont déséquilibrées, on supposera donc dans la suite que

$$\forall i \in \llbracket 1, I \rrbracket, j \in \llbracket 1, J \rrbracket, \quad n_{i,j} = K . \quad (2)$$

Notation :

$$\hat{y}_{i,j} = \frac{1}{K} \sum_{k=1}^K y_{i,j,k} = \bar{y}_{i,j} \quad (3)$$

Contraintes sur l'estimateur : (XXX définition de l'estimateur manquante :)

$$\sum_{i=1}^I \alpha_i = 0 \quad (4)$$

$$\sum_{j=1}^J \beta_j = 0 \quad (5)$$

$$\forall i \in \llbracket 1, I \rrbracket, \quad \sum_{j=1}^J \gamma_{i,j} = 0 \quad (6)$$

$$\forall j \in \llbracket 1, J \rrbracket, \quad \sum_{i=1}^I \gamma_{i,j} = 0 . \quad (7)$$

Mise en forme matricielle : On pose : $n = I \times J \times K$. On peut donc écrire une matrice de design :

$$X = [\mathbf{1}_n, \underbrace{\mathbf{1}_{C_1}, \dots, \mathbf{1}_{C_I}}_{\text{Facteur 1}}, \underbrace{\mathbf{1}_{D_1}, \dots, \mathbf{1}_{D_J}}_{\text{Facteur 2}}, \underbrace{\mathbf{1}_{C_1 \cap D_1}, \dots, \mathbf{1}_{C_I \cap D_J}}_{I.J \text{ termes}}] \quad (8)$$

Nous avons donc $I \times J$ classes pour leurs interactions. En effet :

$$\text{rang}(X) = 1 + (I - 1) + (J - 1) + (I - 1)(J - 1) = I \times J . \quad (9)$$

L'estimateur associé : Sous les contraintes (4), (5), (6) et (7) on obtient donc comme estimateur

$$\begin{aligned}
 & \arg \min_{(\mu, \alpha, \beta) \in \mathbb{R} \times \mathbb{R}^I \times \mathbb{R}^J \times \mathbb{R}^{IJ}} \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (y_{i,j,k} - \mu - \alpha_i - \beta_j - \gamma_{i,j})^2 \\
 & \text{s.c.} \quad \sum_{i=1}^I \alpha_i = 0, \\
 & \quad \sum_{j=1}^J \beta_j = 0, \\
 & \quad \forall i \in \llbracket 1, I \rrbracket, \quad \sum_{j=1}^J \gamma_{i,j} = 0 \\
 & \quad \forall j \in \llbracket 1, J \rrbracket, \quad \sum_{i=1}^I \gamma_{i,j} = 0.
 \end{aligned} \tag{10}$$

On peut alors former le Lagrangien du problème

$$\begin{aligned}
 \mathcal{L}(\mu, \alpha, \beta, \gamma, \lambda) &= \frac{1}{2} \sum_i \sum_j \sum_k (y_{i,j,k} - \alpha_i - \beta_j - \gamma_{i,j})^2 \\
 &+ \lambda_\alpha \sum_i \alpha_i + \lambda_\beta \sum_j \beta_j + \sum_i \nu_i \left(\sum_j \gamma_{i,j} \right) + \sum_j \eta_j \left(\sum_i \gamma_{i,j} \right).
 \end{aligned} \tag{11}$$

Les conditions du premier ordre impliquent :

$$\begin{cases} \hat{\mu} = \bar{y}_n = \frac{1}{IJK} \sum_i \sum_j \sum_k y_{i,j,k} \\ \hat{\alpha}_i = \bar{y}_{i,.,.} - \bar{y}_n \\ \hat{\beta}_j = \bar{y}_{:,j,} - \bar{y}_n \\ \hat{\gamma}_{i,j} = \bar{y}_{i,j,} - \bar{y}_{i,.,.} - \bar{y}_{:,j,} + \bar{y}_n. \end{cases} \tag{12}$$

(XXX détails des calculs manquants.)

Test : Nous voulons tester l'hypothèse nulle :

$$(H_0) : \beta_1^* = \dots = \beta_J^* = 0'' \tag{13}$$

Ce test permet de tester si le deuxième facteur a un effet ou non.

$$F_{obs} = \frac{\frac{1}{J-1} \sum_{j=1}^J (\bar{y}_{:,j,} - \bar{y}_n)^2}{\frac{1}{n-IJ} \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (\bar{y}_n - y_{i,j,k})^2} \tag{14}$$

Sous $H_0 : F_{obs} \sim F_{n-IJ}^{J-1}$. (XXX test de Fisher) On peut fixer le niveau du test à α et considérer les quantiles $F_{n-IJ}^{J-1}(1-\alpha)$.

Remarque 1.2. Par symétrie, en testant " $\alpha_1^* = \dots = \alpha_I^*$ ", le premier facteur n'a pas d'effet sous cette hypothèse.

$$\bar{y}_k = \frac{1}{IJK} \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K y_{i,j,k} \tag{15}$$

- Test :

$(H_0) : \gamma_{i,j} = 0 ; \quad \forall i, j$

$$F_{obs} = \frac{\frac{1}{(I-1)(J-1)} \sum_i \sum_j \hat{\gamma}_{i,j}^2}{\frac{1}{n-IJ} \sum_i \sum_j \sum_k (\bar{y}_n - y_{i,j,k})^2}$$

Sous $H_0 : F_{obs} \sim F_{n-IJ}^{(J-1)(I-1)}$

2 Modèle à un facteur aléatoire

Nous considérons ici les situations où par exemple les traitements sont des échantillons aléatoires provenant d'une large population de traitements. Cela peut sembler assez spécial à première vue, mais c'est en fait très naturel dans de nombreuses situations. Pensez par exemple à un échantillon aléatoire de classes scolaires qui a été tiré de toutes les classes scolaires d'un pays. Un autre exemple pourrait être des machines qui ont été choisies au hasard parmi une grande population de machines. Généralement, nous souhaitons faire une déclaration sur certaines propriétés de l'ensemble de la population et non sur celle des individus observés (ici : classes scolaires ou machines).

2.1 Modèle

$$y_{i,j} = \mu + A_j + \varepsilon_{i,j} . \quad (16)$$

- $\mu \in \mathbb{R}$: **effet fixe**
- $A_j \sim \mathcal{N}(0, \sigma_A^2)$ et *iid* , avec $j = 1, \dots, J$ (les J niveaux); **effet aléatoire**
- $\varepsilon_{i,j} \sim \mathcal{N}(0, \sigma_\varepsilon^2)$ *iid*
- $\forall i, j$, $A_j \perp\!\!\!\perp \varepsilon_{i,j}$
- $n = \sum_{j=1}^J n_j$, où n_j est le nombre d'observations de modalité J .

Regardons maintenant l'espérance et la variance associées au modèle. On a :

$$\mathbb{E}[y_{i,j}] = \mathbb{E}[\mu + A_j + \varepsilon_{i,j}] = \mathbb{E}[\mu] = \mu . \quad (17)$$

Concernant la variance :

$$\begin{aligned} \text{Var}[y_{i,j}] &= \text{Var}[\mu + A_j + \varepsilon_{i,j}] \\ &= \text{Var}[A_j + \varepsilon_{i,j}] \\ &= \text{Var}[A_j] + \text{Var}[\varepsilon_{i,j}] \\ &= \sigma_A^2 + \sigma_\varepsilon^2 . \end{aligned} \quad (18)$$

Nous considérons maintenant la covariance du modèle. Pour j fixé et $i \neq i'$:

$$\text{Cov}(y_{i,j}; y_{i'j}) = \text{Cov}(\mu + A_j + \varepsilon_{i,j}, \mu + A_j + \varepsilon_{i'j}) \quad (19)$$

$$= \text{Cov}(A_j + \varepsilon_{i,j}, A_j + \varepsilon_{i'j}) \quad (20)$$

$$= \text{Cov}(A_j, A_j) + \text{Cov}(\varepsilon_{i,j}, \varepsilon_{i'j}) + \text{Cov}(A_j, \varepsilon_{i,j}) + \text{Cov}(A_j, \varepsilon_{i'j}) \quad (21)$$

$$= \text{Cov}(y_{i,j}; y_{i'j}) \quad (22)$$

$$= \sigma_A^2 . \quad (23)$$

Avec $j \neq j'$:

$$\text{Cov}(y_{i,j}; y_{i'j'}) = \text{Cov}(A_j + \varepsilon_{i,j}, A_{j'} + \varepsilon_{i'j'}) = 0 . \quad (24)$$

Remarque 2.1. J niveaux modélisent (souvent) un effet aléatoire de niveaux qui font partie d'une plus grande population (dont on observe une sous-partie).

Exemples :

- cadre médical : cohorte (sous-partie de la population), essais cliniques, etc.
- étude animale, végétale, etc.

Écriture matricielle :

On a le modèle suivant :

$$y = \nu \mathbf{1}_n + \varepsilon \quad (25)$$

$$\varepsilon \in \mathbb{R}^n \sim \mathcal{N}(0, \sigma^2 \text{Id}_n) . \quad (26)$$

Nous posons alors : $Z = [\mathbf{1}_{C_1}, \mathbf{1}_{C_2}, \dots, \mathbf{1}_{C_J}] \in \mathbb{R}^{n \times J}$ avec $C_1 \cup C_2 \cup \dots \cup C_n = \llbracket 1, n \rrbracket$.

Remarquons également que Z est déterministe.

- $A = (A_1, \dots, A_J)^\top \in \mathbb{R}^J$ et tel que $A \sim \mathcal{N}(0, \sigma_A^2 \text{Id}_J)$.

On a alors :

$$\text{Var}(ZA) = Z \text{Var}(A) Z^\top = \sigma_A^2 (ZZ^\top) . \quad (27)$$

$$ZZ^\top = \sum_{j=1}^J \mathbf{1}_{C_j} \mathbf{1}_{C_j}^\top \in \mathbb{R}^{n \times n} . \quad (28)$$

Donc :

$$\text{Var}(y) = \text{Var}(ZA) + \text{Var}(\varepsilon) = \sigma_A^2 ZZ^\top + \sigma^2 \text{Id} . \quad (29)$$

Exemple : Posons $J = 2$, $n_1 = 3$ et $n_2 = 2$, nous avons alors :

$$Z = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}, A = \begin{pmatrix} A_1 \\ A_2 \end{pmatrix} \text{ et } ZZ^\top = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 \end{pmatrix}. \text{ Donc,}$$

$$\text{Var}(y) = \begin{pmatrix} \sigma_A^2 & \sigma_A^2 & \sigma_A^2 & 0 & 0 \\ \sigma_A^2 & \sigma_A^2 & \sigma_A^2 & 0 & 0 \\ \sigma_A^2 & \sigma_A^2 & \sigma_A^2 & 0 & 0 \\ 0 & 0 & 0 & \sigma_A^2 & \sigma_A^2 \\ 0 & 0 & 0 & \sigma_A^2 & \sigma_A^2 \end{pmatrix} + \sigma_\epsilon^2 \text{Id} = \begin{pmatrix} \sigma_A^2 + \sigma_\epsilon^2 & \sigma_A^2 + \sigma_\epsilon^2 & \sigma_A^2 + \sigma_\epsilon^2 & 0 & 0 \\ \sigma_A^2 + \sigma_\epsilon^2 & \sigma_A^2 + \sigma_\epsilon^2 & \sigma_A^2 + \sigma_\epsilon^2 & 0 & 0 \\ \sigma_A^2 + \sigma_\epsilon^2 & \sigma_A^2 + \sigma_\epsilon^2 & \sigma_A^2 + \sigma_\epsilon^2 & 0 & 0 \\ 0 & 0 & 0 & \sigma_A^2 + \sigma_\epsilon^2 & \sigma_A^2 + \sigma_\epsilon^2 \\ 0 & 0 & 0 & \sigma_A^2 + \sigma_\epsilon^2 & \sigma_A^2 + \sigma_\epsilon^2 \end{pmatrix} \quad (30)$$

Bibliographie

- “ANOVA : A Short Intro Using R”, Lukas Meier, <https://stat.ethz.ch/~meier/teaching/anova/>
- “Modèles à effets aléatoires et modèles mixtes”, Philippe Besse, <https://www.math.univ-toulouse.fr/~besse/Wikistat/pdf/st-m-modmixt6-modmixt.pdf>