Modèles Linéaires Appliqués

Arthur Charpentier

Automne 2Q20

OLS #11 (régression sur des variables qualitatives)



Loi Multinomiale

$$\mathbf{Y}=(Y_1,\cdots,Y_d)\sim\mathcal{M}(\mathbf{p})$$
 où $\mathbf{p}=(p_1,\cdots,p_d)$ si
$$Y_1+\cdots+Y_d=1 \text{ et } Y_j\sim\mathcal{B}(p_j), \ \forall j\in\{1,\cdots,d\}$$

i.e.
$$\mathbf{Y} = (\mathbf{1}_{C_1}, \mathbf{1}_{C_2}, \cdots, \mathbf{1}_{C_d})$$

$$\mathbf{Y} = (Y_1, \cdots, Y_d) \sim \mathcal{M}(n, \mathbf{p})$$
 où $\mathbf{p} = (p_1, \cdots, p_d)$ si

$$Y_1 + \cdots + Y_d = n \text{ et } Y_j \sim \mathcal{B}(n, p_j), \ \forall j \in \{1, \cdots, d\}$$

cf loi multinomiale. Pour

$$(y_1, \dots, y_d) \in S_{d,n} = \{(y_1, \dots, y_d) \in \mathbb{N}^d : (y_1 + \dots + y_d = n)\}$$

$$\mathbb{P}[(Y_1, \dots, Y_d) = (y_1, \dots, y_d)] = \frac{n!}{v_1! \dots v_d!} p_1^{y_1} \dots p_d^{y_d}$$

Example: $\mathbf{Y} = (Y_0, Y_1) \sim \mathcal{M}(n, \mathbf{p})$ où $\mathbf{p} = (p_0, p_1)$.

On considère ici un variable qualitative, ou catégorielle, ou factorielle, x, prenant J modalités.

Pour $j = 1, \dots, J$, on note n_i le nombre d'observations de la jème modalité du facteur, $n = n_1 + \cdots + n_I$.

La variable d'intérêt est (toujours) Y. Notons $y_{i,j}$ la *i*ème observation du facteur j.

Example:

	y	182	161	161	177	157	170	167	186	178	171	175
	X	М	F	F	М	F	М	М	М	М	М	М
devient												
	У	182	177	170	167	186	178	171	175	161	161	157
	X	М	М	М	М	М	М	М	М	F	F	F
	1_{M}	1	1	1	1	1	1	1	1	0	0	0
	1 -	0	0	0	Ο	Ω	Ω	Ω	Ω	1	1	1

On suppose que
$$y_{i,j} = \beta_0 + \beta_j + \varepsilon_{i,j}$$
, où $\mathbb{E}[\varepsilon_{i,j}] = 0$, $cov(\varepsilon_{i,j}, \varepsilon_{i',j'}) = 0$ et $Var(\varepsilon_{i,j}) = \sigma^2$

On va disjoncter la variable de groupe $x \in \{1, 2, \dots, J\}$ en J indicatrices, $(\mathbf{1}_1, \mathbf{1}_2, \dots, \mathbf{1}_J)$

On peut noter $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ avec

$$\mathbf{y} = (y_{1,1}, \cdots, y_{n_1,1}, y_{1,2}, \cdots, y_{n_2,2}, \cdots, y_{1,J}, \cdots, y_{n_J,J})^{\mathsf{T}}$$

$$\triangleright \beta = (\beta_0, \beta_1, \cdots, \beta_J)^{\top}$$

X =
$$[\mathbf{1}_n, \mathbf{A}]$$
 où **A** est une matrice $n \times J$

$$\mathbf{A} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_J} \end{pmatrix}, \text{ et } \mathbf{X} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_2} & \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & & \vdots \\ \mathbf{1}_{n_J} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_J} \end{pmatrix}$$

Ce modèle n'est pas identifiable (X n'est pas de plein rang)

On va alors imposer une contrainte linéaire, (au choix)

- ightharpoonup imposer $\beta_0 = 0$
- imposer $\beta_{i_{\star}} = 0$ pour un certain j_{\star} (référence)
- $imposer \beta_1 + \beta_2 + \cdots + \beta_I = 0$
- imposer $n_1\beta_1 + n_2\beta_2 + \cdots + n_J\beta_J = 0$



▶ imposer $\beta_0 = 0$

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

$$\widehat{oldsymbol{eta}} = (\mathbf{A}^{ op}\mathbf{A})^{-1}\mathbf{A}^{ op}\mathbf{y} = (\overline{y}_{\cdot 1}, \cdots, \overline{y}_{\cdot J}) \in \mathbb{R}^J$$

avec
$$\overline{y}_{.j} = \frac{1}{n_j} \sum_{i=1}^{n_j} y_{i,j}$$

Estimateur sans biais de β , de variance σ^2 diag $(n_1^{-1}, \dots, n_J^{-1})$, minimale parmi les estimateurs linéaires sans biais de β

$$\widehat{\sigma}^2 = \frac{1}{n-J} \sum_{i=1}^{J} \sum_{j=1}^{n_j} (y_{i,j} - \overline{y}_{.j})^2$$

- \triangleright imposer $\beta_{i_{\star}} = 0$
- $\boldsymbol{\beta} = (\beta_0, \beta_1, \cdots, \beta_{i_*-1}, \beta_{i_*+1}, \cdots, \beta_J)^{\mathsf{T}} \in \mathbb{R}^J$
- $\mathbf{X} = [\mathbf{1}_n, \mathbf{A}_{-i}]$ où où \mathbf{A}_{-i} est une matrice $n \times (J-1)$

$$\mathbf{A}_{-j} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots & & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_{j_\star-1}} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{1}_{n_{j_\star+1}} & \cdots & \mathbf{0} \\ \vdots & \vdots & & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_J} \end{pmatrix}$$

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

$$\widehat{oldsymbol{eta}} = (\mathbf{X}^{ op}\mathbf{X})^{-1}\mathbf{X}^{ op}\mathbf{y} = (\widehat{eta}_0, \widehat{eta}_1, \cdots, \widehat{eta}_J)$$

 $\operatorname{avec} \widehat{\beta}_{j_{\star}} = \overline{y}_{\cdot j_{\star}} = \frac{1}{n_{j_{\star}}} \sum_{i=1}^{n_{j_{\star}}} y_{i,j_{\star}} \text{ et } \widehat{\beta}_{j} = \overline{y}_{\cdot j} - \overline{y}_{\cdot j_{\star}}$

Alors

$$\boldsymbol{\beta} = (\beta_0, \beta_1, \cdots, \beta_{J-1})^{\mathsf{T}} \in \mathbb{R}^J$$

▶
$$\mathbf{X} = [\mathbf{1}_n, \mathbf{A}_{\star}]$$
 où \mathbf{A}_{\star} est une matrice $n \times (J-1)$

$$\mathbf{A} = \begin{pmatrix} \mathbf{1}_{n_1} - \frac{n_1}{n_J} \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n_2} - \frac{n_2}{n_J} \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_{J-1}} - \frac{n_{J-1}}{n_J} \mathbf{1}_{n_{J-1}} \end{pmatrix}$$

Alors

$$\widehat{\boldsymbol{\beta}} = (\widehat{\beta}_0, \widehat{\beta}_1, \cdots, \widehat{\beta}_{J-1})$$

avec
$$\widehat{eta}_0 = \overline{y} = rac{1}{n} \sum_{j=1}^J n_j \overline{y}_{\cdot j}$$
 et $\widehat{eta}_j = \overline{y}_{\cdot j} - \overline{y}$

Modèle équivalent au précédent: même résidus, mêmes prévision, estimateur sans biais, etc

- $imposer \beta_1 + \beta_2 + \cdots + \beta_J = 0$
- $m{\beta} = (eta_0, eta_1, \cdots, eta_{J-1})^{\top} \in \mathbb{R}^J$

Alors

$$\widehat{\boldsymbol{\beta}} = (\widehat{\beta}_0, \widehat{\beta}_1, \cdots, \widehat{\beta}_{J-1})$$

$$\operatorname{avec} \widehat{\beta}_0 = \widetilde{y} = \frac{1}{J} \sum_{i=1}^J \overline{y}_{\cdot j} \text{ et } \widehat{\beta}_j = \overline{y}_{\cdot j} - \widetilde{y}$$

Modèle équivalent au précédent: même résidus, mêmes prévision, estimateur sans biais, etc.

Tous ces modèles sont équivalents !



Pour rappel, la formule de décomposition de la variance s'écrit

$$TSS = RSS + ESS$$

Soit

$$||\mathbf{y} - \overline{y}\mathbf{1}||^2 = ||\mathbf{y} - \widehat{\mathbf{y}}||^2 + ||\widehat{\mathbf{y}} - \overline{y}\mathbf{1}||^2$$

On suppose ε Gaussien. On veut tester

 $H_0: \beta_1=\beta_2=\cdots=\beta_J=0.$

Test de Fisher: sous H_0

$$F = \frac{ESS/(J-1)}{RSS/(n-J)} \sim \mathcal{F}(J-1, n-J)$$



ANOVA à deux facteurs

On suppose ici qu'il est possible d'appartenir aux groupes $j \in \{1, 2, \dots, J\}$ et $k \in \{1, 2, \dots, K\}$. Soit $n_{i,k}$ le nombre d'observations dans ce cas.

$$\overline{y}_{\cdot jk} = \frac{1}{n_{jk}} \sum_{i=1}^{n_{jk}} y_{ijk}, \ \overline{y}_{\cdot j} = \frac{1}{n_{j}} \sum_{k=1}^{K} \sum_{i=1}^{n_{jk}} y_{ijk}, \ \overline{y}_{\cdot \cdot k} = \frac{1}{n_{\cdot k}} \sum_{j=1}^{J} \sum_{i=1}^{n_{jk}} y_{ijk}$$

Le modèle s'écrit ici

$$y_{ijk} = \mu + \alpha_j + \beta_k + \gamma_{jk} + \varepsilon_{ijk}$$

avec
$$\mathbb{E}[\varepsilon_{ijk}] = 0$$
, $\text{cov}(\varepsilon_{ijk}, \varepsilon_{i'j'k'}) = 0$ et $\text{Var}(\varepsilon_{i,j,k}) = \sigma^2$



Pour simplifier les notations, supposons $n_{ik} = n/(JK)$ (constant), $\forall i, k$.

On va disjoncter le couple de lois multinomiales

$$(X_1,X_2) \in \{1,2,\cdots,J\} \times \{1,2,\cdots,K\}$$
 en $(J+K+JK)$ indicatrices

$$(\underbrace{\mathbf{1}_1,\mathbf{1}_2,\cdots,\mathbf{1}_J}_{\mathsf{A}},\underbrace{\mathbf{1}_1,\mathbf{1}_2,\cdots,\mathbf{1}_K}_{\mathsf{B}},\underbrace{\mathbf{1}_{11},\mathbf{1}_{12},\cdots,\mathbf{1}_{JK}}_{\mathsf{C}})$$

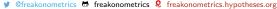
de telle sorte que

$$y = (1, \mathbf{A}, \mathbf{B}, \mathbf{C}) \begin{pmatrix} \mu \\ lpha \\ eta \\ \gamma \end{pmatrix} + arepsilon$$

Là encore, le modèle n'est pas identifiable car X est $n \times (1 + J + K + JK)$ mais de rang JK: il faut un ensemble de 1 + J + K contraintes (linéaires).

- ▶ imposer $\mu = 0$, $\alpha_i = 0 \ \forall j$ et $\beta_k = 0 \ \forall k$
- imposer $\alpha_{i_{+}}=0$, $\beta_{k_{+}}=0$ et $\gamma_{i_{+}k}=\gamma_{ik_{+}}=0$ pour j_{\star} et k_{\star} (modalités de références)
- imposer $\alpha_1 + \alpha_2 + \cdots + \alpha_I = 0$, $\beta_1 + \beta_2 + \cdots + \beta_K = 0$ et

$$\gamma_{j1} + \gamma_{j2} + \cdots + \gamma_{jK} = 0, \ \forall j, \gamma_{1k} + \gamma_{2k} + \cdots + \gamma_{Jk} = 0, \ \forall k$$







$$\mu = 0, \alpha_i = 0 \ \forall i \ \text{et } \beta_k = 0 \ \forall k$$

alors $y_{ijk} = \gamma_{jk} + \varepsilon_{ijk}$. L'estimateur par moindres carrés est

$$\widehat{\gamma}_{jk} = \overline{y}_{jk} = \frac{1}{n_{jk}} \sum_{i=1}^{n_{jk}} y_{ijk}$$

Et

$$\widehat{\sigma}^2 = \frac{1}{n - JK} \sum_{i=1}^{J} \sum_{k=1}^{K} \sum_{i=1}^{n_{jk}} (y_{ijk} - \overline{y}_{jk})^2$$

Alors

$$\widehat{\mu} = \overline{y}, \ \widehat{\alpha}_j = \overline{y}_{j.} - \overline{y}, \ \widehat{\beta}_k = \overline{y}_{.,k} - \overline{y}$$
$$\widehat{\gamma}_{jk} = \overline{y}_{jk} - \overline{y}_{j.} - \overline{y}_{.k} + \overline{y}$$



La décomposition de la variance s'écrit ici

$$TSS = RSS + ESS_A + ESS_B + ESS_C$$

οù

$$ESS_A = \frac{n}{J} \sum_{j=1}^J \widehat{\alpha}_j^2, \ ESS_B = \frac{n}{K} \sum_{k=K}^J \widehat{\beta}_k^2 \ \text{ and } ESS_C = \frac{n}{JK} \sum_{j=1}^J \sum_{k=K}^J \widehat{\gamma}_{jk}^2$$

Le test multiple $H_0: \gamma_{jk} = 0$, $\forall j, k$ est appelé test de l'interaction, et

$$F = \frac{ESS_C/((J-1)(K-1))}{RSS/(n-JK)} \sim \mathcal{F}((J-1)(K-1), n-JK)$$

sous H_0 .



Le test multiple $H_0: \alpha_i = 0$, $\forall j$ est appelé test de l'effet du facteur A. et

$$F = \frac{ESS_A/(J-1)}{RSS/(n-JK)} \sim \mathcal{F}(J-1, n-JK)$$

sous H_0 .

Le test multiple $H_0: \beta_i = 0$, $\forall j$ est appelé test de l'effet du facteur B, et

$$F = \frac{ESS_B/(K-1)}{RSS/(n-JK)} \sim \mathcal{F}(K-1, n-JK)$$

sous H_0 .