

Econométrie des Variables Qualitatives

Emmanuel Duguet

Version 5
2008

TABLE DES MATIÈRES

1	Les variables qualitatives explicatives	6
1.1	Modèle sans terme constant	6
1.2	Modèle avec un terme constant	10
1.3	Modèle avec variables explicatives	11
1.4	Modèle avec produits croisés	12
1.4.1	Cas dichotomique	12
1.4.2	Cas polytomique	13
1.4.3	Cas dichotomique	14
2	Les variables qualitatives expliquées	16
2.1	Variables dichotomiques	16
2.2	Variables polytomiques ordonnées	18
2.3	Variables de comptage	19
2.4	Variables censurées ou tronquées	21
3	Le maximum de vraisemblance	22
3.1	Définitions et propriétés	22
3.2	Estimation	30
3.3	Les moindres carrés ordinaires	34
4	Les algorithmes d'optimisation	38
4.1	Présentation des algorithmes	38
4.2	Les méthodes de gradient	39
4.2.1	Algorithme de Newton-Raphson	40
4.2.2	Algorithme de Berndt-Hall-Hall-Hausman	41
4.2.3	Algorithme du score	42
4.2.4	Algorithme de Levenberg-Marquardt	42
4.3	Méthodologie de programmation	43
5	Les variables dichotomiques	45
5.1	Cas général	45
5.2	Le modèle Logit	48

5.3	Le modèle Probit (ou Normit)	50
5.4	Interprétation et comparaison des coefficients	52
5.4.1	Le modèle Probit	52
5.4.2	Le modèle Logit	53
5.4.3	Comparaison des coefficients des modèles Logit et Probit	54
5.5	Les aides à l'interprétation	54
5.5.1	Variables explicatives binaires	55
5.5.2	Variables explicatives quantitatives	57
5.6	Application : la participation des femmes au marché du travail	58
6	Les variables polytomiques	64
6.1	Cas général	64
6.2	Les variables ordonnées	66
6.2.1	Cas général	66
6.2.2	Le modèle Probit ordonné	67
6.3	Les variables non ordonnées	69
6.3.1	Cas général	69
6.3.2	Le modèle logistique multinomial	69
7	Le pseudo maximum de vraisemblance	73
7.1	Le pseudo maximum de vraisemblance à l'ordre 1	73
7.1.1	La famille exponentielle linéaire à l'ordre 1	73
7.1.2	Estimation	77
7.1.3	Matrice de covariance robuste à l'hétéroscédasticité de forme inconnue	80
7.2	Le pseudo maximum de vraisemblance quasi généralisé	82
7.2.1	La famille exponentielle quasi-généralisée	82
7.2.2	Estimation	83
7.2.3	Les moindres carrés pondérés	83
8	Les variables entières	85
8.1	Le modèle de Poisson	85
8.1.1	Introduction	85
8.1.2	Estimation	87
8.2	Le modèle binomial négatif	89
8.2.1	Estimation par le maximum de vraisemblance	90
8.2.2	Estimation par le pseudo maximum de vraisemblance quasi généralisé	92
8.3	Le modèle avec décision	95
8.4	Le modèle avec saut	96

9	Les variables de durée	98
9.1	Terminologie	99
9.2	Lois usuelles	101
9.2.1	La loi exponentielle	101
9.2.2	La loi de Weibull	102
9.2.3	La loi Gamma généralisée	104
9.2.4	La loi log-normale	105
9.3	Modélisation en logarithmes	107
9.3.1	Rappels	108
9.3.2	Modèle exponentiel et loi de Gumbel	108
9.3.3	Modèle exponentiel et loi exponentielle	110
9.3.4	Modèle de Weibull	111
9.3.5	Modèle Gamma	111
9.3.6	Modèle Gamma généralisé	112
9.3.7	Modèle log-normal	113
9.4	Calcul des moments	114
9.4.1	Fonction génératrice des moments	114
9.4.2	Moments des lois usuelles	115
9.4.3	Résumé	123
9.5	Introduction des variables explicatives	124
9.5.1	Modèles à hasards proportionnels	124
9.5.2	Le modèle exponentiel	125
9.6	Ecriture de la vraisemblance	126
9.6.1	Modèle exponentiel	126
9.6.2	Modèle de Weibull	128
9.6.3	Modèle log-normal	129
9.6.4	Généralisation	130
10	Les variables tronquées	132
10.1	Le modèle tronqué	132
10.2	Le modèle Tobit	135
10.2.1	Estimation	135
10.2.2	Valeur initiale	137
10.2.3	Retour aux paramètres structurels	138
10.3	Le modèle Tobit généralisé	138
10.3.1	Définition	138
10.3.2	Estimation	139
10.3.3	Valeur initiale	141
10.3.4	Amélioration de l'estimation	141
10.3.5	Programmation	142
11	Estimation de modèles à plusieurs équations	144
11.1	Estimation de la forme réduite	144
11.2	Estimation de la forme structurelle	146

A	Moments empiriques et moments théoriques	149
A.1	Moments empiriques des vecteurs	149
A.1.1	Moyenne arithmétique	149
A.1.2	Variance empirique	150
A.1.3	Ecart-type empirique	150
A.1.4	Covariance empirique	151
A.1.5	Corrélation empirique	152
A.2	Moments empiriques des matrices	152
A.2.1	Moyenne arithmétique	152
A.2.2	Matrice de covariance empirique	152
A.3	Convergence en probabilité	156
A.4	Inégalité de Bienaymé-Chebichev	157
A.5	La loi faible des grands nombres	159
A.6	Théorème de la limite centrale	161
B	Algèbre linéaire	162
B.1	Calcul matriciel	162
B.2	Matrices définies positives	163
B.3	Produits de Kronecker	164
C	La loi normale	166
C.1	Loi normale univariée tronquée	167
C.2	Loi normale bivariée	168
C.3	Loi normale conditionnelle	168
C.4	Loi normale bivariée tronquée	170
D	Simplification du calcul des dérivées	171

CHAPITRE 1

Les variables qualitatives explicatives

Les variables qualitatives explicatives sont très nombreuses lorsque l'on étudie les thèmes de l'économie du travail ou de l'innovation. Le but de cette section est d'exposer l'interprétation des coefficients de ces variables dans le cas du modèle linéaire. Ce thème s'étend aux cas où la variable expliquée est qualitative.

Une première utilisation, très répandue, des variables qualitatives consiste à les utiliser sous forme d'indicatrices dans une régression linéaire. Elles servent à indiquer des effets fixes pour indiquer une appartenance à un groupe en général (e.g., région, industrie, catégorie socio professionnelle, niveau de diplôme). Les coefficients de ces variables qualitatives ne s'interprètent plus comme des dérivées par rapport aux variables explicatives, car les dérivées n'existent plus, mais comme un écart moyen par rapport à une modalité de référence. Une seconde utilisation de ces variables qualitatives consiste à découper une variable continue en intervalles puis à examiner la forme de la relation qu'elle entretient avec la variable expliquée. Il s'agit ici d'une approximation par intervalle d'une fonction inconnue.

1.1 Modèle sans terme constant

Nous allons prendre comme exemple introductif une variable qualitative polytomique possédant p modalités. On considère un échantillon de N individus; sans perte de généralité, on suppose que chaque individu appartient à un seul groupe et il y a p groupes différents.¹ Pour sim-

¹Dans le cas où des individus appartiennent à plusieurs groupes dans les données de départ, il est possible de redéfinir la variable qualitative de sorte que tous les individus

plifier l'analyse, on a défini ces groupes de manière à ce qu'ils soient disjoints. On note G_j l'ensemble des indices des individus du groupe j , avec $j = 1, \dots, p$. On remarque que $\bigcup_{j=1}^p G_j = \{1, \dots, N\}$. On considère l'estimation d'un modèle linéaire de la forme suivante :

$$y_i = \sum_{j=1}^p b_j D_{ji} + u_i,$$

$$E(u_i) = 0, E(u_i^2) = \sigma_u^2, E(u_i u_j) = 0 \quad \forall i \neq j, i = 1, \dots, N$$

où y_i est la variable expliquée, u_i la perturbation du modèle et les variables D_{ji} sont des variables qualitatives dichotomiques définies par :

$$D_{ji} = \begin{cases} 1 & \text{si } i \in G_j \\ 0 & \text{si } i \notin G_j \end{cases} \quad i = 1, \dots, N$$

La modélisation de base consiste donc à remplacer la variable qualitative d'appartenance à un groupe par p variables dichotomiques (D_{1i}, \dots, D_{pi}) définies par chacune de ses modalités $j \in \{1, \dots, p\}$. On remarque les propriétés suivantes des variables dichotomiques, qui montrent que le codage binaire $\{0, 1\}$ est le plus pertinent :

1. $D_{ji}^2 = D_{ji}$ puisque $0^2 = 0$ et $1^2 = 1$;
2. $D_{ji} D_{ki} = 0 \quad \forall j \neq k$, car un individu i ne peut pas appartenir à deux groupes à la fois;
3. $\sum_{i=1}^N D_{ji} = \sum_{i \notin G_j} 0 + \sum_{i \in G_j} 1 = N_j$, le nombre d'individus présents dans le groupe j ;
4. $1/N \sum_{i=1}^N D_{ji} = N_j/N$, la fraction des individus du groupe j dans la population totale. Dans le cas des variables dichotomiques, la moyenne arithmétique sert donc à calculer des pourcentages.

En utilisant les propriétés de la perturbation, on voit que :

$$E(y_i | D) = b_j \quad \text{si } i \in G_j,$$

ainsi les coefficients de régression s'interprètent comme les espérances conditionnelles de la variable expliquée dans le groupe j . Ce n'est pas le cas des variables explicatives quantitatives. On peut également interpréter la différence de deux coefficients comme la différence des espérances conditionnelles entre deux groupes :

$$b_j - b_k = E(y_i | i \in G_j) - E(y_i | i \in G_k).$$

appartiennent à un seul groupe.

L'estimation est facilitée en écrivant le modèle individu par individu.

On pose :

$$D_i = (D_{1i}, \dots, D_{ji}, \dots, D_{pi}), \quad i = 1, \dots, N$$

(1,p)

et l'on écrit le vecteur des paramètres en colonne :

$$b = \begin{pmatrix} b_1 \\ \vdots \\ b_p \end{pmatrix}.$$

On obtient donc le modèle linéaire suivant :

$$y_i = D_i b + u_i, \quad i = 1, \dots, N.$$

L'estimateur des moindres carrés ordinaires de b est donc défini par :

$$\hat{b} = \left(\sum_{i=1}^N D_i' D_i \right)^{-1} \sum_{i=1}^N D_i' y_i.$$

La matrice $\sum_{i=1}^N D_i' D_i$ est diagonale et donne les nombres d'observations dans chaque groupe. En effet, en utilisant les propriétés 1 et 2 :

$$\begin{aligned} D_i' D_i &= \begin{pmatrix} D_{1i} \\ D_{ji} \\ D_{pi} \end{pmatrix} (D_{i1}, \dots, D_{ij}, \dots, D_{ip}) \\ &= \begin{pmatrix} D_{1i}^2 & \cdots & D_{1i} D_{ji} & \cdots & D_{1i} D_{pi} \\ \vdots & \ddots & \vdots & & \vdots \\ D_{1i} D_{ji} & \cdots & D_{ji}^2 & \cdots & D_{ji} D_{pi} \\ \vdots & & \vdots & \ddots & \vdots \\ D_{pi} D_{1i} & \cdots & D_{pi} D_{ji} & \cdots & D_{pi}^2 \end{pmatrix} \\ &= \begin{pmatrix} D_{1i} & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & D_{ji} & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & D_{pi} \end{pmatrix} \end{aligned}$$

en conséquence, en utilisant la propriété 3 :

$$\begin{aligned} \sum_{i=1}^N D'_i D_i &= \begin{pmatrix} \sum_{i=1}^N D_{1i} & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & \sum_{i=1}^N D_{ji} & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & \sum_{i=1}^N D_{pi} \end{pmatrix} \\ &= \begin{pmatrix} N_1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & N_j & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & N_p \end{pmatrix}, \end{aligned}$$

ce qui implique :

$$\left(\sum_{i=1}^N D'_i D_i \right)^{-1} = \begin{pmatrix} 1/N_1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & 1/N_j & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 1/N_p \end{pmatrix}$$

La seconde partie de l'estimateur des moindres carrés ordinaires est égale à :

$$\sum_{i=1}^N D'_i y_i = \begin{pmatrix} \sum_{i=1}^N D_{1i} y_i \\ \vdots \\ \sum_{i=1}^N D_{ji} y_i \\ \vdots \\ \sum_{i=1}^N D_{pi} y_i \end{pmatrix} = \begin{pmatrix} \sum_{i \notin G_1} 0 \times y_i + \sum_{i \in G_1} 1 \times y_i \\ \vdots \\ \sum_{i \notin G_j} 0 \times y_i + \sum_{i \in G_j} 1 \times y_i \\ \vdots \\ \sum_{i \notin G_p} 0 \times y_i + \sum_{i \in G_p} 1 \times y_i \end{pmatrix}$$

Dans l'ensemble on obtient donc les moyennes arithmétiques des p groupes :

$$\hat{b} = \begin{pmatrix} 1/N_1 \sum_{i \in G_1} y_i \\ \vdots \\ 1/N_j \sum_{i \in G_j} y_i \\ \vdots \\ 1/N_p \sum_{i \in G_p} y_i \end{pmatrix} \triangleq \begin{pmatrix} \bar{y}_1 \\ \vdots \\ \bar{y}_j \\ \vdots \\ \bar{y}_p \end{pmatrix}$$

1.2 Modèle avec un terme constant

Ici il est inutile de refaire les calculs. En effet, les moindres carrés ordinaires reviennent à faire une projection orthogonale du vecteur des observations de la variable expliquée y sur le sous-espace vectoriel engendré par les vecteurs correspondants des variables explicatives, noté $\text{Im}(D_1, \dots, D_p)$. Ces vecteurs sont linéairement indépendants et forment donc une base de cet espace vectoriel. Pour trouver les coefficients du modèle avec terme constant, il faut avoir en tête les deux éléments suivants :

1. Le terme constant, noté $e_{(N,1)}$ est égal à la somme des vecteurs D_j : $e = \sum_{j=1}^p D_j$.
2. La décomposition d'un vecteur y en une base est unique, et les coefficients des moindres carrés ordinaires sont les coordonnées du vecteur y dans la base (D_1, \dots, D_p) .

La première propriété implique que, dans un modèle avec terme constant, il faut retirer un des vecteur D_j de la liste des variables explicatives pour éviter une multicollinéarité parfaite. La seconde propriété permet de calculer les nouveaux estimateurs des MCO en fonction de \hat{b} . Si on retire la modalité k de la liste des groupes, on estime le modèle :

$$y = c_0 e + c_1 D_1 + \dots + c_{k-1} D_{k-1} + c_{k+1} D_{k+1} + \dots + c_p D_p + u,$$

après estimation de ce modèle par les moindres carrés ordinaires, on obtient une prévision :

$$\hat{y} = \hat{c}_0 e + \hat{c}_1 D_1 + \dots + \hat{c}_{k-1} D_{k-1} + \hat{c}_{k+1} D_{k+1} + \dots + \hat{c}_p D_p,$$

en remplaçant la constante par sa valeur, $e = \sum_{j=1}^p D_j$, on obtient la formulation équivalente :

$$\begin{aligned} \hat{y} &= \hat{c}_0 (D_1 + \dots + D_p) + \hat{c}_1 D_1 + \dots + \hat{c}_{k-1} D_{k-1} + \hat{c}_{k+1} D_{k+1} + \dots + \hat{c}_p D_p \\ &= (\hat{c}_0 + \hat{c}_1) D_1 + \dots + (\hat{c}_0 + \hat{c}_{k-1}) D_{k-1} + \hat{c}_0 D_k + (\hat{c}_0 + \hat{c}_{k+1}) D_{k+1} + \\ &\quad \dots + (\hat{c}_0 + \hat{c}_p) D_p. \end{aligned}$$

La prévision du modèle de départ est égale à :

$$\hat{y} = \hat{b}_1 D_1 + \dots + \hat{b}_{k-1} D_{k-1} + \hat{b}_k D_k + \hat{b}_{k+1} D_{k+1} + \dots + \hat{b}_p D_p,$$

en utilisant l'unicité de la décomposition en une base, on obtient :

$$\left. \begin{array}{l} \widehat{c}_0 + \widehat{c}_1 = \widehat{b}_1 \\ \vdots \\ \widehat{c}_0 + \widehat{c}_{k-1} = \widehat{b}_{k-1} \\ \widehat{c}_0 = \widehat{b}_k \\ \widehat{c}_0 + \widehat{c}_{k+1} = \widehat{b}_{k+1} \\ \vdots \\ \widehat{c}_0 + \widehat{c}_p = \widehat{b}_p \end{array} \right\} \Leftrightarrow \left\{ \begin{array}{l} \widehat{c}_0 = \widehat{b}_k \\ \widehat{c}_1 = \widehat{b}_1 - \widehat{b}_k \\ \vdots \\ \widehat{c}_{k-1} = \widehat{b}_{k-1} - \widehat{b}_k \\ \widehat{c}_{k+1} = \widehat{b}_{k+1} - \widehat{b}_k \\ \vdots \\ \widehat{c}_p = \widehat{b}_p - \widehat{b}_k \end{array} \right.$$

La constante du nouveau modèle représente l'effet de l'indicatrice qui a été enlevée de la régression, et les autres coefficients l'écart entre le coefficient de l'indicatrice courante et de l'indicatrice enlevée. Ainsi l'indicatrice qui a été enlevée correspond à la modalité de référence. C'est la raison pour laquelle il faut indiquer explicitement les modalités des indicatrices enlevées dans les tableaux de régression, elle sont indispensables à l'interprétation.

Remarque 1.1 Le test de Fisher sur le modèle avec terme constant revient à tester ici l'égalité jointe des moyennes entre tous les groupes. En effet, le test correspond à l'hypothèse nulle $H_0 : c_1 = \dots = c_p = 0 \Leftrightarrow H_0 : E(y_j) - E(y_k) = 0 \forall j \neq k$. On notera qu'on ne teste pas la nullité du terme constant du modèle c_0 .

Remarque 1.2 On peut utiliser un simple test de Student pour tester l'égalité des moyennes entre un groupe donné, k , et un autre groupe. Il suffit de mettre un terme constant dans le modèle et d'enlever l'indicatrice du groupe dont on teste l'égalité de la moyenne avec les autres groupes.

1.3 Modèle avec variables explicatives

On introduit maintenant un autre jeu de variables, dont la matrice est notée X , dans le modèle de départ :

$$y_i = \underset{(1,m)(m,1)}{X_i} a + \underset{(1,p)(p,1)}{D_i} b + u_i,$$

on a clairement :

$$E(y_i | X_i, D_{ji} = 1) = X_i a + b_j,$$

de sorte que les coefficients b_j représentent les écarts de moyenne conditionnelle entre deux groupes :

$$\begin{aligned} E(y_i | X_i, D_{ji} = 1) - E(y_i | X_i, D_{ki} = 1) &= (X_i a + b_j) - (X_i a + b_k) \\ &= b_j - b_k. \end{aligned}$$

Les résultats de la section précédente sont donc toujours valables. Le terme constant représente le coefficient de l'indicatrice qui a été retirée et les coefficients des autres indicatrices doivent s'interpréter en écart au coefficient de l'indicatrice retirée.

1.4 Modèle avec produits croisés

1.4.1 Cas dichotomique

On peut introduire les produits croisés de manière naturelle à partir du modèle suivant. Considérons que des individus bénéficient d'une mesure d'aide que nous supposons affectée au hasard (i.e., sans biais de sélection). On note :

$$D_i = \begin{cases} 1 & \text{si l'individu } i \text{ est aidé} \\ 0 & \text{sinon} \end{cases}$$

Une fois cette mesure attribuée, on examine une variable de performance, liée aux objectifs de l'aide, que l'on note y_i . En théorie, pour chaque individu, cette mesure peut prendre deux valeurs :

- y_{0i} : la valeur de y_i si l'individu i n'est pas aidé;
- y_{1i} : la valeur de y_i si l'individu i est aidé.

Ce que l'on cherche à évaluer est l'effet de la mesure, noté :

$$\alpha = E(y_{1i} - y_{0i}),$$

qui représente la moyenne des variations de performance associée à la mesure, prise sur l'ensemble des individus. On peut aller plus loin, en introduisant un modèle explicatif des performances potentielles des individus :

$$\begin{aligned} y_{0i} &= a_0 + X_i c_0 + u_{0i} \\ y_{1i} &= a_1 + X_i c_1 + u_{1i}, \end{aligned}$$

où X_i représente les déterminants de la performance. Les coefficients a_1 et a_0 représentent les niveaux moyens de performances en $X_i = 0$, selon que l'on est aidé ou non. Pour obtenir un modèle empiriquement estimable, il faut l'écrire en fonction de quantités observables. Or, on n'observe que y_{0i} lorsque $D_i = 0$ et seulement y_{1i} si $D_i = 1$. La seule variable observable est :

$$y_i = T_i y_{1i} + (1 - T_i) y_{0i} = \begin{cases} y_{1i} & \text{si } T_i = 1 \\ y_{0i} & \text{si } T_i = 0 \end{cases}$$

En conséquence, le modèle économétrique s'écrit :

$$\begin{aligned} y_i &= T_i (a_1 + X_i c_1 + u_{1i}) + (1 - T_i) (a_0 + X_i c_0 + u_{0i}) \\ &= a_0 + X_i c_0 + T_i \underbrace{(a_1 - a_0)}_a + T_i X_i \underbrace{(c_1 - c_0)}_c + u_i, \end{aligned}$$

où a est le coefficient de l'aide et c le vecteur des coefficients des variables explicatives. La perturbation est égale à :

$$u_i = T_i u_{1i} + (1 - T_i) u_{0i}.$$

Ce modèle fait apparaître un produit croisé entre la variable d'aide T_i et les variables explicatives de la performance X_i . L'estimation de ce modèle permet d'évaluer l'effet de la mesure car :

$$\begin{aligned} \delta &= \text{E}(y_{1i} - y_{0i}) \\ &= \text{E}(a_1 - a_0 + X(c_1 - c_0)) \\ &= a_1 - a_0 + \text{E}(X)(c_1 - c_0) \\ &= a + \text{E}(X)c \end{aligned}$$

que l'on peut estimer sans biais et de manière convergente par :

$$\hat{\delta} = \hat{a} + \bar{X} \hat{c},$$

on remarque que lorsque les variables X sont centrées avant de prendre les produits croisés ($\bar{X} = 0$), l'estimateur $\hat{\delta}$ est obtenu directement par le coefficient de la variable indicatrice d'aide \hat{a} dans la régression avec produits croisés. On remarque également que ce modèle suppose qu'il existe une hétéroscédasticité par bloc car :

$$V(u_i) = \begin{cases} V(u_{0i}) & \text{si } T_i = 0 \\ V(u_{1i}) & \text{si } T_i = 1 \end{cases}$$

et dans le cas où la mesure d'aide affecte également la variance, $V(u_{0i}) \neq V(u_{1i})$, il faut estimer le modèle par les moindres carrés pondérés.

1.4.2 Cas polytomique

On introduit maintenant, en plus des variables explicatives et des indicatrices, les produits croisés des indicatrices et des variables explicatives. On a donc :

$$y_i = D_i b + (X_i \otimes D_i) c + u_i,$$

avec :

$$c_{(mp,1)} = \begin{pmatrix} c_1 \\ \vdots \\ c_j \\ \vdots \\ c_p \end{pmatrix},$$

le terme en X_i a été retiré puisque $\sum_{j=1}^p X_i D_{ji} = X_i$. L'espérance conditionnelle dans le groupe j devient maintenant :

$$E(y_i | X_i, D_{ji} = 1) = X_i c_j + b_j,$$

d'où la différence entre les groupes j et k :

$$\begin{aligned} \gamma_i &\triangleq E(y_i | X_i, D_{ji} = 1) - E(y_i | X_i, D_{ki} = 1) \\ &= (X_i c_j + b_j) - (X_i c_k + b_k) \\ &= X_i (c_j - c_k) + b_j - b_k, \end{aligned}$$

l'effet varie en fonction des caractéristiques de l'individu i au sein du groupe j . Ce modèle autorise donc une hétérogénéité individuelle au sein de chaque groupe. L'effet moyen est égal à :

$$\begin{aligned} \bar{\gamma} &= \frac{1}{N} \sum_{i=1}^N \gamma_i \\ &= \bar{X} (c_j - c_k) + b_j - b_k. \end{aligned}$$

On peut toutefois estimer directement la partie de l'écart entre les groupes qui ne dépend pas des variables explicatives du modèle, $b_j - b_k$, en utilisant la méthode suivante. On centre la variable X avant de faire les produits croisés; avec cette convention $\bar{X} = 0$ et on obtient directement la différence entre les groupes par $b_j - b_k$. Cette dernière mesure l'écart de moyenne entre les groupes une fois que l'on a éliminé l'effet des variables de X sur ces moyennes.

1.4.3 Cas dichotomique

On considère une variable dichotomique $T_i \in \{0, 1\}$ dont on veut connaître l'effet sur y_i . La variable T_i peut être une caractéristique individuelle ou une mesure de politique économique individuelle. La modalité $T_i = 1$ correspond aux individus qui ont bénéficié de la mesure. La variable y_i est alors une mesure de performance choisie en fonction de l'objectif de politique économique. Les variables explicatives de la performance X_i sont centrées avant de prendre les produits croisés, de sorte que l'on a $\bar{X} = 0$. Le modèle, facilement généralisable, est donné par :

$$E(y_i | X_i, T_i) = d + X_i a + T_i b + (X_i \otimes T_i) c_1, \quad (1.1)$$

où d est le terme constant du modèle. On obtient les espérances suivantes :

$$\begin{aligned} E(y_i | X_i, T_i = 0) &= d + X_i a, \\ E(y_i | X_i, T_i = 1) &= d + b + X_i (a + c_1), \end{aligned}$$

d'où l'effet de T_i sur y_i :

$$\gamma_i = E(y_i|X_i, T_i = 1) - E(y_i|X_i, T_i = 0) = b + X_i c_1,$$

après estimation par les moindres carrés ordinaires on obtient :

$$\hat{\gamma}_i = \hat{b} + X_i \hat{c}_1,$$

d'où l'effet moyen de T_i sur l'échantillon :

$$\hat{\gamma} = \frac{1}{N} \sum_{i=1}^N \hat{\gamma}_i = \hat{b} + \left(\frac{1}{N} \sum_{i=1}^N X_i \right) \hat{c}_1 = \hat{b}.$$

Quand on centre les variables explicatives X_i , on peut donc obtenir directement l'effet moyen de T_i par son coefficient dans la régression (1.1).

CHAPITRE 2

Les variables qualitatives expliquées

Les bases de données microéconomiques comprennent invariablement des données tronquées, connues seulement par intervalle ou encore de type purement qualitatif. Par exemple, dans l'enquête Innovation du SESSI, on connaît le simple fait d'avoir réalisé une innovation ou encore une appréciation de l'entreprise sur l'importance d'un déterminant de l'innovation. Ce n'est pas toujours un inconvénient, car certains phénomènes ne sont pas quantifiables de façon objective. Dans l'enquête Emploi de l'INSEE on peut connaître le fait qu'un individu travaille et le nombre d'heures travaillées par les individus qui ont travaillé sur la période d'enquête. Mais on peut préférer créer une variable qualitative de type "pas d'emploi, temps partiel, temps plein" si l'objet de l'étude le justifie. Malgré la perte d'information inévitable quand on passe par exemple, de la valeur exacte d'une variable à sa connaissance par intervalle, il est toujours possible d'étudier ses effets, même si cela comporte certaines limites. Dans cette section, nous donnons quelques exemples de variables qualitatives et leur représentation en économétrie.

2.1 Variables dichotomiques

Une variable dichotomique est une variable qui ne peut prendre que deux modalités exclusives l'une de l'autre, comme "Oui/Non" ou "Inférieur ou égal à/Strictement supérieur à". Par convention, on code une modalité à 0 et l'autre à 1. Ce n'est pas une obligation, mais cette présentation permet de simplifier la présentation du problème. La variable associée est appelée une indicatrice, dans l'exemple suivant la variable y_i est une

indicatrice d'innovation :

$$y_i = \begin{cases} 1 & \text{si l'entreprise } i \text{ a innové} \\ 0 & \text{sinon} \end{cases}$$

Ce que l'on peut étudier à partir de ce type de variable, ce sont les déterminants de la décision d'innover. Cette décision se traduit en termes statistiques par une probabilité d'innover. Ainsi, on va rechercher quelles sont les variables qui réduisent ou au contraire augmentent la probabilité d'innover. Il faut donc construire un modèle qui nous permette d'estimer l'effet d'un ensemble de variables sur la probabilité qu'une entreprise innove. Pour cela on construit ce que l'on appelle un modèle latent, c'est à dire inobservable. On peut faire une analogie entre ce modèle latent et le modèle théorique qui sert de base à l'étude.

Si l'on pouvait mesurer l'innovation sous forme quantitative, on utiliserait le modèle linéaire standard. Mais, soit pour des raisons de collecte des données soit parce que le phénomène que l'on étudie n'est pas quantifiable, on ne dispose que d'une information qualitative sur celui-ci. Dans notre exemple, on sait juste si l'entreprise a innové ou non. Le modèle latent est le modèle linéaire standard :

$$\pi_i^* = X_i b + u_i, \quad i = 1, \dots, N$$

Que représente ce modèle ? La variable endogène π^* est inobservable. On peut l'interpréter ici comme l'espérance de profit associé à l'introduction d'une innovation, compte-tenu d'un effet de remplacement des anciens produits. Les variables explicatives X sont les déterminants de ce gain et le paramètre b mesure l'importance de ces déterminants. Il nous faut maintenant expliquer comment on passe de π^* à y , c'est à dire de ce qui n'est pas observable (π_i^*, u_i) à ce qui l'est (y_i, X_i), car seules ces dernières informations peuvent être utilisées en pratique.

Il est raisonnable de penser que toutes les entreprises cherchent à améliorer leurs produits et leurs procédés de production, même marginalement. Le résultat anticipé de cette activité est justement représentée dans notre modèle par π^* . Mais on n'observera la mise en oeuvre d'une innovation que si cette activité procure des gains significatifs, s'ils dépassent un certain seuil. Notons que ce gain ne sera significatif que si l'innovation l'est également et que cette notion de seuil correspond bien à la définition retenue dans les enquêtes sur l'innovation.¹ Seules sont considérées dans les enquêtes les améliorations significatives de produit

¹ Ainsi dans l'enquête du SESSI sur "l'innovation technologique dans l'industrie", annexée à l'Enquête Annuelle d'Entreprise de 1990, l'innovation de produit est définie sur le questionnaire par : "Un produit est considéré comme technologiquement innovant s'il donne lieu à la création d'un nouveau marché ou s'il peut se distinguer substantiellement de produits précédemment fabriqués, d'un point de vue technologique ou par les prestations rendues à l'utilisateur. Ne sont pas concernées

et de procéder. Soit le seuil π_0 , qui peut dépendre de chaque industrie, on a :

$$y_i = \begin{cases} 1 & \text{si } \pi_i^* > \pi_0 \\ 0 & \text{si } \pi_i^* \leq \pi_0 \end{cases}$$

Ceci implique que l'on peut maintenant calculer la probabilité d'innover. Elle est égale à :

$$\Pr [y_i = 1] = \Pr [\pi_i^* > \pi_0].$$

Il reste alors à faire une hypothèse sur la distribution conditionnelle de π^* sachant X pour obtenir une forme fonctionnelle précise. Selon l'hypothèse que l'on fait, on obtient un modèle Logit (loi logistique) ou un modèle Probit (loi normale). Cette liste n'est bien sûr pas limitative et chaque hypothèse de distribution mène à un modèle différent. Des tests sont alors nécessaires pour trancher.

2.2 Variables polytomiques ordonnées

Cette fois-ci, la variable qualitative que l'on observe peut prendre plus de deux modalités qui sont ordonnées entre elles². Elles peuvent être définies aussi bien par rapport à une quantité que traduire une appréciation. Par exemple, dans l'enquête Innovation du Ministère de l'Industrie (SESSI) le pourcentage de produits de moins de cinq ans d'âge dans le chiffre d'affaires est donné sous la forme suivante : entre 0 et 10% , de 10% à 30%, de 30% à 70% et plus de 70% . Dans cette même enquête, l'importance de la recherche développement du groupe auquel appartient l'entreprise comme déterminant de l'innovation est donnée sous la forme : "pas du tout", "un peu", "moyennement" et "beaucoup". Dans les deux cas, les modalités traduisent un ordre, qui indique l'intensité de la variable. Le modèle latent représente alors la vraie valeur de la variable, qui n'est pas observable. Cette variable, que l'on cherche à expliquer, est représentée par le modèle latent linéaire :

$$y_i^* = X_i b + u_i, \quad i = 1, \dots, N.$$

La variable observable, qualitative, prend maintenant une forme plus

les innovations de nature purement esthétique ou de style (design); en revanche sont concernées, mais isolées, les innovations de conditionnement ou d'emballage." Pour une présentation de l'enquête et un exemplaire du questionnaire, voir François (1991).

²Il existe également des variables qualitatives non ordonnées qui représentent des choix. Les plus connues représentent le choix de mode du transport comme : véhicule individuel , bus , métro.

générale :

$$y_i = \begin{cases} 1 & \text{si } \alpha_0 < y_i^* \leq \alpha_1 \\ 2 & \text{si } \alpha_1 < y_i^* \leq \alpha_2 \\ \vdots & \\ r & \text{si } \alpha_{r-1} < y_i^* \leq \alpha_r \end{cases}$$

Les bornes délimitent les valeurs α_0 et α_r que peut prendre la variable y_i^* . Pour une variable réelle, on adopte la convention $\alpha_0 = \{-\infty\}$ et $\alpha_r = \{+\infty\}$. Plus généralement les bornes peuvent être connues ou inconnues. Pour le pourcentage d'innovation décrit plus haut, elles sont égales à 0, 0.1, 0.3, 0.7 et 1. Dans d'autre cas, les variables sont toujours ordonnées mais on ne connaît pas les seuils. C'est le cas quand les personnes interrogées répondent à une question par "Pas du tout, un peu, moyennement, beaucoup". Pourtant les seuils théoriques existent bien puisque l'on peut ordonner les modalités, on suppose simplement qu'ils sont *constants* au sein d'une population donnée. Dans les deux cas, seuils connus ou inconnus, on peut estimer un modèle pour trouver les déterminants de y_i^* . Cette fois-ci, la probabilité d'observer la modalité j est donnée par :

$$\Pr[y_i = j] = \Pr[\alpha_{j-1} < y_i^* \leq \alpha_j] = \Pr[y_i^* \leq \alpha_j] - \Pr[y_i^* \leq \alpha_{j-1}], \\ j = 1, \dots, r.$$

Une fois que l'on a spécifié la loi conditionnelle de y^* sachant X , on peut procéder aux estimations à partir des variables observables (y_i, X_i) . Les modèles polytomiques ordonnés peuvent être utilisés pour ce genre de variable endogène. Si la loi des perturbations u_i est normale, on obtient un modèle Probit polytomique ordonné. D'autres hypothèses sur la loi des perturbations u_i donnent d'autres modèles.

2.3 Variables de comptage

Certaines données d'innovation sont discrètes. Ainsi le nombre de brevets n'est pas une donnée quantitative de même nature que les dépenses de recherche et développement. Il s'agit d'une variable qui ne prend que des valeurs entières. Qui plus est, il s'agit du comptage d'événements relativement rares. Sur une année, en France, on compte beaucoup d'entreprises qui ne déposent pas de brevet. Il peut s'agir du résultat d'une décision mais également du simple fait que l'entreprise n'a pas trouvé d'innovation brevetable durant l'année écoulée. La variable expliquée prend ses valeurs dans l'ensemble des entiers naturels $y_i \in \{0, 1, 2, \dots\}$.

Ce processus est par nature aléatoire et, comme pour les autres variables, on modélise son espérance mathématique. Ici toutefois, cette espérance est toujours strictement positive et l'on prend donc une forme

exponentielle :

$$E(y_i | X_i, b) = \exp(X_i b + u_i) > 0.$$

Cette espérance mathématique est alors supposée être celle d'une loi de Poisson, utilisée pour représenter les variables endogènes discrètes positives ou nulles. Notons bien qu'il y a *deux sources* d'aléas dans cette dernière modélisation. La première vient de l'erreur que l'on fait sur la moyenne de la variable expliquée, représentée par $\exp(u_i)$, la seconde vient du tirage dans une loi de Poisson dont la moyenne est aléatoire. Dans les modèles usuels, tout l'aléa provient de l'erreur que l'on fait sur la moyenne.

Lorsqu'il n'y a pas de perturbation dans la moyenne ($V(\exp u_i) = 0, \forall i$), on parle du modèle de Poisson *homogène*, dans le cas inverse il s'agit du modèle de Poisson *hétérogène*. Notons que l'on peut faire un parallèle entre les données de comptage et les données de durée, car une donnée de comptage donne le nombre d'événements qui se sont produits pendant une durée donnée. On montre que la loi de Poisson homogène correspond à une loi de durée exponentielle.

2.4 Variables censurées ou tronquées

Une variable censurée ou tronquée est une variable dont on observe la réalisation pour certains individus seulement. La troncature peut provenir soit du processus de collecte des données soit d'une décision prise par ces mêmes individus. C'est ce dernier cas qui nous intéresse.³ Prenons le cas de l'activité de recherche et développement : une entreprise doit à la fois décider si elle investit ou non dans un programme de recherche et combien elle y investit. Ces deux décisions sont étroitement reliées. Le processus de décision est représenté par une variable latente, qui peut être le critère de décision π^* . Les déterminants de cette décision sont notés X_1 . On pose :

$$\pi_i^* = X_{1i}b_1 + u_{1i}, \quad i = 1, \dots, N.$$

Cette première variable latente génère une variable qualitative dichotomique :⁴

$$y_i = \begin{cases} 1 & \text{si } \pi_i^* \geq 0 \\ 0 & \text{si } \pi_i^* < 0 \end{cases}$$

Cette indicatrice nous dit si l'entreprise a investi en recherche ou non. Mais elle détermine également s'il est possible d'observer le montant investi en recherche, représenté par une seconde variable r^* . L'investissement en recherche r^* est expliqué par le modèle :

$$r_i^* = X_{2i}b_2 + u_{2i}, \quad i = 1, \dots, N,$$

où X_2 contient les déterminants du montant investi en recherche. On admet de plus que les deux variables latentes, π_i^* et r_i^* sont corrélées entre elles. Cette corrélation provient du fait que l'on obtient généralement r^* en maximisant le profit π^* , ce qui implique que les deux variables sont déterminées simultanément. La variable de recherche observable, notée r est donc donnée par :

$$r_i = \begin{cases} r_i^* & \text{si } \pi_i^* \geq 0 \\ \text{manquant} & \text{si } \pi_i^* < 0 \end{cases}$$

où "manquant" indique une valeur manquante dans la base de données. Lorsque les perturbations u_{1i} et u_{2i} suivent une loi normale bivariée on obtient le modèle *tobit généralisé* de Heckman.

³Le cas des censures exogènes correspond au modèle tobit simple, le lecteur pourra trouver une présentation de ce modèle dans Gouriéroux (1989) et Maddala (1983). Le cas que nous présentons ici est celui d'une censure endogène aboutissant au modèle tobit généralisé, étudié à l'origine par Heckman (1976, 1979).

⁴Le seuil peut être mis à 0 sans perte de généralité tant que les variables explicatives contiennent un terme constant.

CHAPITRE 3

Le maximum de vraisemblance

Le maximum de vraisemblance est une méthode d'estimation qui repose sur la distribution conditionnelle des variables que l'on étudie. Intuitivement, elle consiste à estimer un paramètre inconnu en choisissant la valeur de ce paramètre qui maximise la "probabilité" d'observer l'échantillon que l'on observe effectivement. La vraisemblance de l'échantillon est soit la probabilité d'observer l'échantillon (cas discret) soit la densité correspondante (cas continu).

3.1 Définitions et propriétés

Plus généralement, on suppose que la variable expliquée y admet une distribution conditionnelle par rapport aux variables explicatives X dont la densité conditionnelle ou probabilité conditionnelle est notée $f(y|X; \theta)$ où θ est le paramètre que l'on cherche à estimer. On suppose ici que les N observations présentes dans l'échantillon $y = (y_1, \dots, y_N)$ sont indépendantes. La vraisemblance de l'échantillon, notée L , est définie par :

$$L(y|X; \theta) = \prod_{i=1}^N f(y_i|X_i; \theta).$$

Voici quelques exemples, pour des modèles sans variable explicative.

Exemple 3.1 *Loi normale.* Soit un échantillon de variables réelles (y_1, \dots, y_N) iid selon une loi normale $N(\theta, \omega)$ où ω est un nombre positif connu. Sa densité est donnée par :

$$f(y; \theta) = \frac{1}{\sqrt{2\pi\omega}} \exp \left\{ -\frac{1}{2\omega} (y - \theta)^2 \right\}.$$

La log-vraisemblance de cet échantillon est égale à :

$$\ell(y|X; \theta) = \sum_{i=1}^N \ln f(y_i|\theta) = -\frac{N}{2} \ln(2\pi\omega) - \frac{1}{2\omega} \sum_{i=1}^N (y_i - \theta)^2.$$

Exemple 3.2 *Loi de Poisson.* Soit un échantillon de variables entières positives (y_1, \dots, y_N) iid selon une loi de Poisson de paramètre θ :

$$f(y; \theta) = \frac{\exp(-\theta) \theta^y}{y!}.$$

La log-vraisemblance de cet échantillon est égale à :

$$\ell(y|X; \theta) = \sum_{i=1}^N \ln f(y_i|\theta) = -N\theta + \ln(\theta) \sum_{i=1}^N y_i - \sum_{i=1}^N \ln(y_i!).$$

Exemple 3.3 *Loi de Bernoulli.* Soit un échantillon de variables dichotomiques (y_1, \dots, y_N) iid selon une loi de Bernoulli de paramètre θ . Les probabilités sont égales à :

$$f(y; \theta) = \theta^y (1 - \theta)^{1-y} = \begin{cases} \theta & \text{si } y = 1 \\ 1 - \theta & \text{si } y = 0 \end{cases}.$$

La log-vraisemblance de cet échantillon est égale à :

$$\ell(y|X; \theta) = \sum_{i=1}^N \ln f(y_i|\theta) = \ln\left(\frac{\theta}{1-\theta}\right) \sum_{i=1}^N y_i + N \ln(1-\theta).$$

La méthode du maximum de vraisemblance consiste à estimer θ par $\hat{\theta}_n$ tel que :

$$L(y|X; \hat{\theta}_n) \geq L(y|X; \theta) \quad \forall \theta \in \Theta,$$

où Θ est l'ensemble des valeurs admissibles du paramètre θ . Cet estimateur est appelé estimateur du maximum de vraisemblance de θ ou en abrégé EMV de θ . Notons ici que cette inégalité est équivalente à :

$$\ln L(y|X; \hat{\theta}_n) \geq \ln L(y|X; \theta) \quad \forall \theta \in \Theta,$$

de sorte que l'on peut maximiser la log-vraisemblance $\ln L$ au lieu de la vraisemblance L . Cette méthode permet de simplifier l'écriture des dérivées de la fonction objectif, car la dérivée d'une somme est plus simple que la dérivée d'un produit.

PROPRIÉTÉ 3.1 *Sous les hypothèses de régularité habituelles (Gouriéroux et Monfort, 1989, ch. VII, p. 192), que nous supposons vérifiées par la*

suite, les estimateurs du maximum de vraisemblance sont convergents, asymptotiquement normaux et asymptotiquement efficaces (i.e., à variance minimale parmi les estimateurs convergents) :

$$\sqrt{N} \left(\widehat{\theta}_N - \theta \right) \xrightarrow[N \rightarrow +\infty]{L} N[0, I_1^{-1}(\theta)],$$

où $I_1(\theta)$ est la matrice d'information de Fisher définie par :

$$I_1(\theta) = \mathbb{E}_X \mathbb{V}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \right] = \mathbb{E}_X \mathbb{E}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \frac{\partial \ln f(y|X, \theta)}{\partial \theta'} \right].$$

De plus, en vertu de l'égalité de la matrice d'information, on peut aussi utiliser la matrice $J_1(\theta)$:

$$J_1(\theta) = \mathbb{E}_X \mathbb{E}_y \left[- \frac{\partial^2 \ln f(y|X, \theta)}{\partial \theta \partial \theta'} \right],$$

car on a $I_1(\theta) = J_1(\theta)$. La distribution de $\widehat{\theta}_N$ peut donc être approximée par :

$$\widehat{\theta}_N \overset{A}{\rightsquigarrow} N \left(\theta, \frac{1}{N} \times I_1(\widehat{\theta}_N)^{-1} \right) = N \left(\theta, I_N(\widehat{\theta}_N)^{-1} \right).$$

où $\overset{A}{\rightsquigarrow}$ désigne une distribution asymptotique (i.e., utilisable pour de grands échantillons).

Remarque 3.1 La matrice d'information de l'ensemble de l'échantillon est définie par $I_N(\theta) = N \times I_1(\theta)$, on a donc :

$$\frac{1}{N} I_1^{-1}(\theta) = (N \times I_1(\theta))^{-1} = I_N(\theta)^{-1}.$$

Pour voir d'où vient le résultat de normalité asymptotique, il suffit de partir de la définition de l'estimateur du maximum de vraisemblance. Cette définition est implicite et donnée par la condition du premier ordre :

$$\frac{\partial \ln L(y|X, \widehat{\theta})}{\partial \theta} = \sum_{i=1}^N \frac{\partial \ln f(y_i|X_i, \widehat{\theta})}{\partial \theta} = 0. \quad (3.1)$$

En effectuant un développement limité de $\partial \ln L / \partial \theta$ au voisinage de $\widehat{\theta}$ on obtient :

$$\frac{\partial \ln L(y|X, \widehat{\theta})}{\partial \theta} \underset{A}{=} \frac{\partial \ln L(y|X, \theta)}{\partial \theta} + \frac{\partial^2 \ln L(y|X, \theta)}{\partial \theta \partial \theta'} (\widehat{\theta} - \theta).$$

On remarque ici que ce développement limité devient exact quand $N \rightarrow +\infty$. La condition (3.1) implique que :

$$0 \underset{A}{=} \frac{\partial \ln L(y|X, \theta)}{\partial \theta} + \frac{\partial^2 \ln L(y|X, \theta)}{\partial \theta \partial \theta'} (\widehat{\theta} - \theta),$$

de sorte que l'on peut écrire :

$$\begin{aligned}\hat{\theta} - \theta &\stackrel{A}{=} \left[-\frac{\partial^2 \ln L(y|X, \theta)}{\partial \theta \partial \theta'} \right]^{-1} \frac{\partial \ln L(y|X, \theta)}{\partial \theta} \\ \Leftrightarrow \sqrt{N} (\hat{\theta} - \theta) &\stackrel{A}{=} \left[-\frac{1}{N} \frac{\partial^2 \ln L(y|X, \theta)}{\partial \theta \partial \theta'} \right]^{-1} \frac{1}{\sqrt{N}} \frac{\partial \ln L(y|X, \theta)}{\partial \theta}\end{aligned}$$

La première quantité du membre de droite de l'équation est une moyenne qui converge en probabilité vers l'espérance mathématique correspondante. En appliquant la loi des grands nombres :

$$\begin{aligned}-\frac{1}{N} \frac{\partial^2 \ln L(y|X, \theta)}{\partial \theta \partial \theta'} &= \frac{1}{N} \sum_{i=1}^N \left[-\frac{\partial^2 \ln f(y_i|X_i, \theta)}{\partial \theta \partial \theta'} \right] \\ &\xrightarrow{p} \underset{X}{\mathbb{E}} \underset{y}{\mathbb{E}} \left[-\frac{\partial^2 \ln f(y|X, \theta)}{\partial \theta \partial \theta'} \right] \stackrel{A}{=} J_1(\theta).\end{aligned}$$

Le second terme du membre de droite de l'équation suit, asymptotiquement, une loi normale. On peut écrire :

$$\begin{aligned}\frac{1}{\sqrt{N}} \frac{\partial \ln L(y|X, \theta)}{\partial \theta} &= \frac{1}{\sqrt{N}} \sum_{i=1}^N \frac{\partial \ln f(y_i|X_i, \theta)}{\partial \theta} \\ &= \frac{1}{N} \sum_{i=1}^N \sqrt{N} \frac{\partial \ln f(y_i|X_i, \theta)}{\partial \theta} \\ &= \frac{1}{N} \sum_{i=1}^N z_i,\end{aligned}$$

où z_i est la variable dont on cherche la distribution. On a alors, sous les hypothèses usuelles :

$$\sqrt{N} (\bar{z} - \mathbb{E}(\bar{z})) \xrightarrow[N \rightarrow +\infty]{L} N(0, V(\bar{z})), \quad (3.2)$$

où p est le nombre d'éléments de θ . Pour appliquer le théorème de la limite centrale, on a besoin de l'espérance et de la variance de z_i . Pour trouver l'espérance de z_i on utilise la propriété suivante :

PROPRIÉTÉ 3.2 *Soit $f(y|X, \theta)$ la densité conditionnelle de la variable expliquée. Elle vérifie la propriété suivante :*

$$\underset{y}{\mathbb{E}} \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \right] = 0.$$

Preuve :

$$\begin{aligned}
\mathbb{E}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \right] &= \mathbb{E}_y \left[\frac{1}{f(y|X, \theta)} \frac{\partial f(y|X, \theta)}{\partial \theta} \right] \\
&= \int \frac{1}{f(y|X, \theta)} \frac{\partial f(y|X, \theta)}{\partial \theta} f(y|X, \theta) dy \\
&= \int \frac{\partial f(y|X, \theta)}{\partial \theta} dy \\
&= \frac{\partial}{\partial \theta} \underbrace{\int f(y|X, \theta) dy}_1 \\
&= 0,
\end{aligned}$$

□

On voit que :

$$\mathbb{E}(z_i) = \sqrt{N} \mathbb{E}_X \mathbb{E}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \right] = 0,$$

Passons maintenant à la variance; nous avons besoin de la quantité suivante :

$$\begin{aligned}
\mathbb{V}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \right] &= \mathbb{E}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \frac{\partial \ln f(y|X, \theta)}{\partial \theta'} \right] \\
&\quad - \underbrace{\mathbb{E}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \right] \mathbb{E}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta'} \right]}_0 \\
&= \mathbb{E}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \frac{\partial \ln f(y|X, \theta)}{\partial \theta'} \right]
\end{aligned}$$

Pour calculer la variance de z_i , on utilise la propriété suivante :

PROPRIÉTÉ 3.3 Soit $f(y|X, \theta)$ la densité conditionnelle de la variable expliquée. Elle vérifie la propriété suivante :

$$\mathbb{E}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \frac{\partial \ln f(y|X, \theta)}{\partial \theta'} \right] = \mathbb{E}_y \left[-\frac{\partial^2 \ln f(y|X, \theta)}{\partial \theta \partial \theta'} \right].$$

Preuve :

On dérive la relation suivante par rapport à θ' :

$$\int \frac{\partial \ln f(y|X, \theta)}{\partial \theta} f(y|X, \theta) dy = 0$$

$$\Rightarrow \int \left\{ \frac{\partial^2 \ln f(y|X, \theta)}{\partial \theta \partial \theta'} f(y|X, \theta) + \frac{\partial \ln f(y|X, \theta)}{\partial \theta} \frac{\partial f(y|X, \theta)}{\partial \theta'} \right\} dy = 0$$

or

$$\frac{\partial \ln f(y|X, \theta)}{\partial \theta} = \frac{1}{f(y|X, \theta)} \frac{\partial f(y|X, \theta)}{\partial \theta'}$$

$$\Leftrightarrow \frac{\partial f(y|X, \theta)}{\partial \theta'} = \frac{\partial \ln f(y|X, \theta)}{\partial \theta} f(y|X, \theta),$$

en remplaçant dans la relation (3.3), on obtient :

$$\int \left\{ \frac{\partial^2 \ln f(y|X, \theta)}{\partial \theta \partial \theta'} + \frac{\partial \ln f(y|X, \theta)}{\partial \theta} \frac{\partial \ln f(y|X, \theta)}{\partial \theta'} \right\} f(y|X, \theta) dy = 0 \quad (3.3)$$

$$\Leftrightarrow \int \frac{\partial \ln f(y|X, \theta)}{\partial \theta} \frac{\partial \ln f(y|X, \theta)}{\partial \theta'} f(y|X, \theta) dy =$$

$$- \int \frac{\partial^2 \ln f(y|X, \theta)}{\partial \theta \partial \theta'} f(y|X, \theta) dy \quad (3.4)$$

$$\Leftrightarrow \mathbb{E}_y \left[\frac{\partial \ln f(y|X, \theta)}{\partial \theta} \frac{\partial \ln f(y|X, \theta)}{\partial \theta'} \right] = \mathbb{E}_y \left[- \frac{\partial^2 \ln f(y|X, \theta)}{\partial \theta \partial \theta'} \right].$$

□

La variance de z_i est donnée par la formule de la variance totale :

$$\begin{aligned}
V(z_i) &= \mathbb{E}_X \mathbb{V}_y [z_i] + \underbrace{\mathbb{V}_X \mathbb{E}_y [z_i]}_0 \\
&= N \mathbb{E}_X \left[\mathbb{V}_y \left(\frac{\partial \ln f(y_i | X_i, \theta)}{\partial \theta} \right) \right] \\
&= N \mathbb{E}_X \mathbb{E}_y \left[\frac{\partial \ln f(y_i | X_i, \theta)}{\partial \theta} \frac{\partial \ln f(y_i | X_i, \theta)}{\partial \theta'} \right] \\
&= N \mathbf{I}_1(\theta),
\end{aligned}$$

donc

$$\begin{aligned}
V \left(\frac{1}{N} \sum_{i=1}^N z_i \right) &= \frac{1}{N^2} \sum_{i=1}^N V(z_i) \\
&= \frac{N V(z_i)}{N^2} \\
&= \frac{N^2 \mathbf{I}_1(\theta)}{N^2} \\
&= \mathbf{I}_1(\theta),
\end{aligned}$$

qui est une quantité finie. Globalement, on trouve que :

$$\frac{1}{\sqrt{N}} \frac{\partial \ln L(y|X, \theta)}{\partial \theta} = \bar{z} \xrightarrow[N \rightarrow +\infty]{L} N(0, \mathbf{I}_1(\theta)),$$

et l'on déduit de (3.2) que :

$$\sqrt{N} (\hat{\theta} - \theta) \stackrel{A}{\underset{\sim}{\rightrightarrows}} \mathbf{J}_1(\theta)^{-1} \bar{z},$$

converge en loi vers une distribution normale d'espérance nulle et de variance :

$$\mathbf{J}_1(\theta)^{-1} V(\bar{z}) \mathbf{J}_1(\theta)^{-1} = \mathbf{J}_1(\theta)^{-1} \mathbf{I}_1(\theta) \mathbf{J}_1(\theta)^{-1} = \mathbf{J}_1(\theta)^{-1} = \mathbf{I}_1(\theta)^{-1}.$$

□

Une dernière propriété est utile, celle de l'invariance fonctionnelle. Elle permet de retrouver l'estimateur du maximum de vraisemblance après un changement de paramètres.

Si l'on effectue un changement de paramètre du type $\tau = h(\theta)$, où h est une fonction inversible choisie par l'économètre, on a :

$$L(y|X, h^{-1}(\hat{\tau}_n)) = L(y|X, \hat{\theta}_n) \geq L(y|X, \theta) = L(y|X, h^{-1}(\tau)), \forall \theta \in \Theta$$

donc $\hat{\tau}_n$ est l'estimateur du maximum de vraisemblance de τ . Il n'est donc pas nécessaire de réestimer le modèle quand on effectue un changement de paramètre.

PROPRIÉTÉ 3.4 (Invariance fonctionnelle)

Soit $\hat{\theta}_N$ un estimateur du maximum de vraisemblance de θ et $\tau = h(\theta)$ un changement de paramètre. L'estimateur du maximum de vraisemblance de τ est donné par $\hat{\tau}_N = h(\hat{\theta}_N)$.

Cette propriété est très pratique car certaines log-vraisemblances ne sont concaves que par rapport à un changement de paramètres bien précis (e.g., dans le modèle tobit généralisé). On est alors certain que l'optimum ne dépend pas de ce changement de paramètre, et que les algorithmes courants convergent vers ce maximum. On effectue donc toujours les changements de paramètres de ce type quand ils existent. Mais encore faut-il remonter des nouveaux paramètres τ aux paramètres structurels du modèle θ . Pour cela on utilise le théorème de Slutsky.

THÉORÈME 3.1 (Slutsky)

Soit h une fonction de classe C^1 (dérivable et de dérivée première continue), une relation entre deux paramètres $\tau = h(\theta)$, et $\hat{\theta}_N$ estimateur convergent de θ vérifiant :

$$\sqrt{N}(\hat{\theta}_N - \theta) \xrightarrow[N \rightarrow +\infty]{L} N(0, \Omega_{\hat{\theta}})$$

alors,

$$\sqrt{N}(h(\hat{\theta}_N) - h(\theta)) \xrightarrow[N \rightarrow +\infty]{L} N\left[0, \frac{\partial h}{\partial \theta}(\theta) \Omega_{\hat{\theta}} \frac{\partial h}{\partial \theta'}(\theta)\right]$$

Notons bien que ce théorème s'applique même si $\hat{\theta}_N$ n'est pas un estimateur du maximum de vraisemblance et même si la fonction h n'est pas inversible. Dans la pratique, on estimera la variance asymptotique de $h(\hat{\theta}_n)$ par :

$$\widehat{\text{Vas}}[h(\hat{\theta}_N)] = \frac{1}{N} \frac{\partial h}{\partial \theta}(\hat{\theta}_N) \hat{\Omega}_{\hat{\theta}} \frac{\partial h}{\partial \theta'}(\hat{\theta}_N),$$

où $\hat{\Omega}_{\hat{\theta}}$ est un estimateur convergent de $\Omega_{\hat{\theta}}$.

Preuve :

Pour comprendre ce résultat, il suffit de faire un développement limité de $h(\hat{\theta}_N)$ au voisinage de θ :

$$\begin{aligned} h(\hat{\theta}_N) &\stackrel{A}{=} h(\theta) + \frac{\partial h}{\partial \theta}(\theta) (\hat{\theta}_N - \theta) \\ &\Leftrightarrow \sqrt{N} \left(h(\hat{\theta}_N) - h(\theta) \right) \stackrel{A}{=} \frac{\partial h}{\partial \theta}(\theta) \sqrt{N} (\hat{\theta}_N - \theta), \end{aligned}$$

cette expression est une transformation linéaire de $\sqrt{N}(\hat{\theta}_N - \theta)$, en conséquence elle suit asymptotiquement une loi normale. Son espérance mathématique est égale à :

$$\begin{aligned} \mathbb{E} \left(\sqrt{N} \left(h(\hat{\theta}_N) - h(\theta) \right) \right) &= \mathbb{E} \left(\frac{\partial h}{\partial \theta}(\theta) \sqrt{N} (\hat{\theta}_N - \theta) \right) \\ &= \frac{\partial h}{\partial \theta}(\theta) \underbrace{\mathbb{E} \left(\sqrt{N} (\hat{\theta}_N - \theta) \right)}_0 = 0, \end{aligned}$$

et sa variance est égale à :

$$\begin{aligned} \mathbb{V} \left(\sqrt{N} \left(h(\hat{\theta}_N) - h(\theta) \right) \right) &= \mathbb{V} \left(\frac{\partial h}{\partial \theta}(\theta) \sqrt{N} (\hat{\theta}_N - \theta) \right) \\ &= \frac{\partial h}{\partial \theta}(\theta) \underbrace{\mathbb{V} \left(\sqrt{N} (\hat{\theta}_N - \theta) \right)}_{\Omega_{\hat{\theta}}} \frac{\partial h}{\partial \theta'}(\theta). \end{aligned}$$

□

3.2 Estimation

Les points candidats à un maximum sont obtenus par la résolution des conditions du premier ordre. En effet, dans les cas usuels les conditions du premier ordre fournissent un maximum local sous réserve de vérification de la condition du second ordre. Il faut alors rechercher numériquement les maxima locaux et prendre celui qui fournit la valeur la plus élevée de la vraisemblance. Toutefois, la plupart des modèles que nous verrons dans ce cours possèdent une log-vraisemblance concave. Dans ce cas particulier, le maximum est unique et donné par les conditions du premier ordre. La log-vraisemblance est égale à :

$$\ell(y|X, \theta) = \ln L(y|X, \theta).$$

Quand la solution est unique, on cherche la solution du problème d'optimisation :¹

$$\hat{\theta}_N = \arg \max_{\theta} \ell(y|X, \theta).$$

La condition du premier ordre pour un maximum local est donnée par la nullité du score :

$$\frac{\partial \ell}{\partial \theta} (y|X, \hat{\theta}_N) = 0,$$

et la condition du second ordre par un hessien défini négatif :

$$\frac{\partial^2 \ell}{\partial \theta \partial \theta'} (y|X, \hat{\theta}_N) \ll 0,$$

où \ll désigne l'infériorité au sens des matrices.

En général l'expression de $\hat{\theta}_N$ ne peut pas être obtenue directement en fonction des observations, c'est-à-dire sous forme explicite², il faut donc recourir à des algorithmes d'optimisation numérique pour effectuer une estimation par le maximum de vraisemblance. Une fois cette valeur obtenue, il nous faut estimer la matrice de covariance de $\hat{\theta}_N$.

En utilisant la loi des grands nombres, on peut estimer les moments théoriques par les moments empiriques correspondants, soit :

$$\begin{aligned} \mathbf{I}_1(\theta) &= \mathbb{E}_X \mathbb{E}_y \left[\frac{\partial \ln f}{\partial \theta} (y|X, \theta) \frac{\partial \ln f}{\partial \theta'} (y|X, \theta) \right] \\ \text{par } \hat{\mathbf{I}}_1(\hat{\theta}_N) &= \frac{1}{N} \sum_{i=1}^N \frac{\partial \ln f}{\partial \theta} (y_i|X_i, \hat{\theta}_N) \frac{\partial \ln f}{\partial \theta'} (y_i|X_i, \hat{\theta}_N), \\ \text{et } \mathbf{J}_1(\theta) &= \mathbb{E}_X \mathbb{E}_y \left[-\frac{\partial^2 \ln f}{\partial \theta \partial \theta'} (y|X, \theta) \right] \\ \text{par } \hat{\mathbf{J}}_1(\hat{\theta}_N) &= -\frac{1}{N} \sum_{i=1}^N \frac{\partial^2 \ln f}{\partial \theta \partial \theta'} (y_i|X_i, \hat{\theta}_N). \end{aligned}$$

Il suffit de prendre l'inverse de l'une de ces deux matrices pour obtenir un estimateur convergent de la matrice de covariance asymptotique de $\sqrt{N}(\hat{\theta}_N - \theta)$. Voici quelques exemples, pour des modèles sans variable explicative.

¹Quand la solution n'est pas unique, on a $\hat{\theta}_n \in \arg \max_{\theta} \ell(y|x, \theta)$ car il y a un ensemble de solutions.

²Le modèle linéaire standard est une exception à cette règle.

Exemple 3.4 Loi normale $N(\theta, \omega)$. On maximise la log-vraisemblance :

$$\ell(y|X, \theta) = -\frac{N}{2} \ln(2\pi\omega) - \frac{1}{2\omega} \sum_{i=1}^N (y_i - \theta)^2,$$

ce qui donne la condition du premier ordre :

$$\frac{\partial \ell}{\partial \theta} (y|\hat{\theta}_N) = \frac{1}{\omega} \sum_{i=1}^N (y_i - \hat{\theta}_N),$$

qui permet d'obtenir l'estimateur du maximum de vraisemblance :

$$\hat{\theta}_N = \frac{1}{N} \sum_{i=1}^N y_i$$

avec

$$\frac{\partial^2 \ell}{\partial \theta^2} (y|\theta) = -\frac{N}{\omega} < 0 \quad \forall \theta.$$

Pour trouver la distribution asymptotique de $\hat{\theta}_N$, on peut utiliser soit $I_N(\theta)$ soit $J_N(\theta)$. On a :

$$\begin{aligned} I_N(\theta) &= \sum_{i=1}^N \mathbb{E}_y \left[\frac{\partial \ln f}{\partial \theta} (y_i|\theta)^2 \right] \\ &= \frac{1}{\omega^2} \sum_{i=1}^N \mathbb{E}_y \left[(y_i - \theta)^2 \right] \\ &= \frac{N}{\omega}, \end{aligned}$$

de même :

$$J_N(\theta) = \sum_{i=1}^N \mathbb{E}_y \left[-\frac{\partial^2 \ln f}{\partial \theta^2} (y_i|\theta) \right] = \frac{N}{\omega}.$$

et l'on obtient :

$$\hat{\theta}_N \overset{A}{\rightsquigarrow} N(\theta, \omega/N).$$

Exemple 3.5 Loi de Poisson $P(\theta)$. On maximise la log-vraisemblance :

$$\ell(y|X, \theta) = -N\theta + \ln(\theta) \sum_{i=1}^N y_i - \sum_{i=1}^N \ln(y_i!),$$

ce qui donne la condition du premier ordre :

$$\frac{\partial \ell}{\partial \theta} (y|\hat{\theta}_N) = -N + \frac{1}{\hat{\theta}_N} \sum_{i=1}^N y_i = 0,$$

qui permet d'obtenir l'estimateur du maximum de vraisemblance :

$$\hat{\theta}_N = \frac{1}{N} \sum_{i=1}^N y_i$$

avec

$$\frac{\partial^2 \ell}{\partial \theta^2}(y|\theta) = -\frac{1}{\theta^2} \sum_{i=1}^N y_i < 0 \quad \forall \theta.$$

Pour trouver la distribution asymptotique de $\hat{\theta}_N$, on peut utiliser soit $I_N(\theta)$ soit $J_N(\theta)$. On a :

$$\begin{aligned} I_N(\theta) &= \sum_{i=1}^N V_y \left[\frac{\partial \ln f}{\partial \theta}(y_i|\theta) \right] \\ &= \frac{1}{\theta^2} \sum_{i=1}^N \underbrace{V_y[y_i]}_{\theta} \\ &= \frac{N}{\theta}, \end{aligned}$$

de même :

$$\begin{aligned} J_N(\theta) &= \sum_{i=1}^N E_y \left[-\frac{\partial^2 \ln f}{\partial \theta^2}(y_i|\theta) \right] \\ &= \frac{1}{\theta^2} \sum_{i=1}^N \underbrace{E_y(y_i)}_{\theta} \\ &= \frac{N}{\theta}. \end{aligned}$$

et l'on obtient :

$$\hat{\theta}_N \overset{A}{\rightsquigarrow} N(\theta, \theta/N).$$

Exemple 3.6 Loi de Bernoulli $B(\theta)$. On maximise la log-vraisemblance :

$$\ell(y|X; \theta) = \ln \left(\frac{\theta}{1-\theta} \right) \sum_{i=1}^N y_i + N \ln(1-\theta).$$

ce qui donne la condition du premier ordre :

$$\frac{\partial \ell}{\partial \theta}(y|\hat{\theta}_N) = \frac{1}{\hat{\theta}_N(1-\hat{\theta}_N)} \sum_{i=1}^N y_i - \frac{N}{1-\hat{\theta}_N} = 0,$$

qui permet d'obtenir l'estimateur du maximum de vraisemblance :

$$\hat{\theta}_N = \frac{1}{N} \sum_{i=1}^N y_i$$

avec

$$\begin{aligned} \frac{\partial^2 \ell}{\partial \theta^2} (y|\theta) &= -\frac{1-2\theta}{\theta^2(1-\theta)^2} \sum_{i=1}^N y_i - \frac{N}{(1-\theta)^2} \\ &= -\frac{(1-\theta)^2 \sum_{i=1}^N y_i + \theta^2 \sum_{i=1}^N (1-y_i)}{\theta^2(1-\theta)^2} < 0 \quad \forall \theta, \end{aligned}$$

en utilisant $N = \sum_{i=1}^N y_i + \sum_{i=1}^N (1-y_i)$. Pour trouver la distribution asymptotique de $\hat{\theta}_N$, on peut utiliser soit $I_N(\theta)$ soit $J_N(\theta)$. On a :

$$\begin{aligned} I_N(\theta) &= \sum_{i=1}^N \underset{y}{V} \left[\frac{\partial \ln f}{\partial \theta} (y_i|\theta) \right] \\ &= \frac{1}{\theta^2(1-\theta)^2} \sum_{i=1}^N \underbrace{\underset{\theta(1-\theta)}{V} [y_i]} \\ &= \frac{N}{\theta(1-\theta)}, \end{aligned}$$

de même :

$$\begin{aligned} J_N(\theta) &= \sum_{i=1}^N \underset{y}{E} \left[-\frac{\partial^2 \ln f}{\partial \theta^2} (y_i|\theta) \right] \\ &= \frac{1}{\theta^2(1-\theta)^2} \sum_{i=1}^N \left[(1-2\theta) \underset{y}{E} (y_i) + \theta^2 \right] \\ &= \frac{N}{\theta(1-\theta)}. \end{aligned}$$

et l'on obtient :

$$\hat{\theta}_N \overset{A}{\rightsquigarrow} N(\theta, \theta(1-\theta)/N).$$

3.3 Les moindres carrés ordinaires

Soit le modèle linéaire standard :

$$y_i = X_i b + u_i, \quad i = 1, \dots, N$$

où l'indice i désigne une observation, N la taille de l'échantillon, y_i la variable expliquée, X_i le vecteur des variables explicatives, b le vecteur de leurs coefficients et u_i la perturbation du modèle. On suppose que les perturbations sont indépendantes et identiquement distribuées (*iid*) selon une loi normale $N(0, \omega)$. Ces hypothèses impliquent que les variables expliquées y_i , transformations linéaires des u_i , sont indépendantes et suivent des lois normales $N(X_i b, \omega)$. En conséquence la densité conditionnelle de la i -ème observation est égale à :

$$f(y_i | X_i; b, \omega) = \frac{1}{\sqrt{2\pi\omega}} \exp \left\{ -\frac{(y_i - X_i b)^2}{2\omega} \right\}.$$

Le paramètre à estimer est donc :

$$\theta = \begin{pmatrix} b \\ \omega \end{pmatrix}.$$

Comme les observations sont indépendantes, la vraisemblance de l'échantillon, notée L , est égale au produit des densités individuelles :

$$\begin{aligned} L(y|X; \theta) &= \prod_{i=1}^N f(y_i | X_i; \theta) \\ &= (2\pi\omega)^{-\frac{N}{2}} \exp \left\{ -\frac{1}{2\omega} \sum_{i=1}^N (y_i - X_i b)^2 \right\}. \end{aligned}$$

La méthode du maximum de vraisemblance consiste à choisir des valeurs de b et ω qui rendent cette densité jointe la plus grande possible. La log-vraisemblance de l'échantillon est égale à :

$$\ell(y|X; \theta) = -\frac{N}{2} \ln(2\pi\omega) - \frac{1}{2\omega} \sum_{i=1}^N (y_i - X_i b)^2.$$

On pose :

$$\hat{\theta}_N = \begin{pmatrix} \hat{b}_N \\ \hat{\omega}_N \end{pmatrix}.$$

Les conditions du premier ordre sont données par :

$$\begin{aligned} \frac{\partial \ell}{\partial \theta} (y|X, \hat{\theta}_N) &= \begin{bmatrix} \frac{\partial \ell}{\partial b} (y|X, \hat{\theta}_N) \\ \frac{\partial \ell}{\partial \omega} (y|X, \hat{\theta}_N) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\ \Leftrightarrow \begin{bmatrix} \frac{1}{\hat{\omega}_N} \sum_{i=1}^N X_i' (y_i - X_i \hat{b}_N) \\ -\frac{N}{2\hat{\omega}_N} + \frac{1}{2\hat{\omega}_N^2} \sum_{i=1}^N (y_i - X_i \hat{b}_N)^2 \end{bmatrix} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \end{aligned}$$

ce qui implique :

$$\hat{b}_N = \left(\sum_{i=1}^N X_i' X_i \right)^{-1} \sum_{i=1}^N X_i' y_i \quad \text{et} \quad \hat{\omega}_N = \frac{1}{N} \sum_{i=1}^N (y_i - X_i \hat{b}_N)^2.$$

On retrouve l'estimateur des moindres carrés ordinaires de b . Par contre l'estimateur de la variance ne comporte pas de correction pour les degrés de liberté.³ On estime la matrice de covariance asymptotique à partir de :

$$\begin{aligned} J_N(\theta) &= \sum_{i=1}^N E_{\theta} \begin{bmatrix} \frac{1}{\omega} X_i' X_i & \frac{1}{\omega^2} X_i' (y_i - X_i b) \\ \frac{1}{\omega^2} X_i (y_i - X_i b) & -\frac{1}{2\omega^2} + \frac{1}{\omega^3} (y_i - X_i b)^2 \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{\omega} \sum_{i=1}^N X_i' X_i & 0 \\ 0 & \frac{N}{2\omega^2} \end{bmatrix} \end{aligned}$$

La matrice de covariance asymptotique de $\hat{\theta}_N$ est donc estimée par :

$$\widehat{\text{Vas}}(\hat{\theta}) = \hat{J}_N(\hat{\theta}_N)^{-1} = \begin{bmatrix} \hat{\omega}_N \left[\sum_{i=1}^N X_i' X_i \right]^{-1} & 0 \\ 0 & \frac{2\hat{\omega}_N^2}{N} \end{bmatrix},$$

elle permet de construire un intervalle de confiance sur la variance. Pour raisonner sur l'écart-type, on utilise l'invariance fonctionnelle et le théorème de Slutsky. L'invariance fonctionnelle garantit que :

$$\hat{\sigma}_N = \sqrt{\widehat{\omega}_N},$$

³Habituellement, à distance finie, on estime ω par $1/(N-p) \sum (y_i - x_i \hat{b}_N)^2$, où p est le nombre de variables explicatives (y compris le terme constant). La distribution asymptotique de cet estimateur est la même que celle de $\hat{\omega}_N$.

est bien l'estimateur du maximum de vraisemblance de $\sigma = \sqrt{V(u_i)}$. Et le théorème de Slutsky permet de calculer la variance de l'estimateur de l'écart-type. On a une transformation :

$$h(\omega) = \sqrt{\omega} \Rightarrow h'(\omega) = \frac{1}{2\sqrt{\omega}},$$

on peut donc écrire que :

$$\begin{aligned} \widehat{\text{Vas}}(\hat{\sigma}) &= \frac{1}{2\sqrt{\hat{\omega}_N}} \frac{2\hat{\omega}_N^2}{N} \frac{1}{2\sqrt{\hat{\omega}_N}} \\ &= \frac{\hat{\omega}_N}{2N}. \end{aligned}$$

CHAPITRE 4

Les algorithmes d'optimisation

4.1 Présentation des algorithmes

Le problème que nous devons résoudre est de trouver une valeur numérique $\hat{\theta}_N$ qui résout un système d'équations de la forme :¹

$$s(\hat{\theta}_N) = 0$$

où $s(\cdot)$ est une fonction connue des observations et du paramètre θ . Elle dépend aussi bien des données que du modèle postulé (distribution, paramètres). Dans le cas du maximum de vraisemblance $s(\theta)$ est appelé le *score*, dans le cas du pseudo maximum de vraisemblance, on l'appelle le *pseudo score*. Quand la fonction à maximiser est concave, cette condition du premier ordre est suffisante pour un maximum global. Ce sera le cas du modèle Logit, avec lequel nous ferons une des premières applications. La technique que l'on utilise pour parvenir à la valeur $\hat{\theta}_N$ s'appelle un algorithme. On peut décomposer un algorithme en quatre grandes étapes.

1. *Une valeur initiale. Le choix de la valeur initiale n'est pas problématique quand l'objectif est concave. Dans ce cas, avec un algorithme croissant (voir plus loin), tout point de départ doit mener au maximum. Par contre, quand l'objectif n'est pas concave, ou pour accélérer la procédure quand l'objectif est concave, on prendra un estimateur convergent comme point de départ. Il est en effet possible dans certains cas de trouver un estimateur en deux étapes*

¹Un système d'équations car $s(\hat{\theta}_N)$ est un vecteur.

relativement facile à calculer. Il n'est généralement pas efficace, et c'est la raison pour laquelle on réalise une estimation supplémentaire. On note cette valeur initiale $\theta_{(0)}$.

2. Une règle d'itération. Une fois la valeur initiale fixée, il faut utiliser une règle qui permette de trouver une nouvelle valeur plus proche du maximum. Le pas de l'itération, défini comme la différence entre deux valeurs successives du paramètre, est déterminé selon différentes méthodes et constitue le coeur de l'algorithme. Nous utiliserons des méthodes dites de gradient et plus particulièrement les algorithmes de Newton-Raphson, du Score, de Berndt-Hall-Hall-Hausman et de Levenberg-Marquardt. La pratique montre qu'ils permettent de traiter la plupart des cas, même difficiles. On résume cette étape par la relation $\theta_{(p+1)} = \theta_{(p)} + M(\theta_{(p)})$ où p est l'itération, $p \in N$. La valeur du pas $M(\theta_{(p)})$ doit dépendre uniquement de la valeur du paramètre à l'étape précédente.
3. Une règle d'arrêt de l'algorithme. Au voisinage du maximum la fonction objectif ne doit plus varier, on peut donc baser l'arrêt de la procédure sur la différence entre deux valeurs successives de la fonction objectif. Une seconde condition porte sur le gradient, qui doit être nul (condition du premier ordre). On peut aussi utiliser des variantes, comme l'élasticité de l'objectif aux paramètres du modèle, qui est également nulle à l'optimum et possède l'avantage d'être insensible à un changement d'unité des variables (contrairement au gradient). Dans l'ensemble, tous ces critères sont équivalents à l'optimum.
4. Vérifier que l'on a bien atteint un maximum local. Ce problème ne se pose que lorsque l'objectif n'est pas globalement concave. La condition du second ordre pour un optimum local précise que le hessien doit être défini négatif au point en question. Il faudra donc le vérifier systématiquement. Ceci est d'autant plus important que dans le cas du maximum de vraisemblance l'inverse du hessien n'est autre qu'un estimateur convergent de la matrice de covariance du paramètre. En conséquence, si cette propriété n'était pas vérifiée, on obtiendrait un estimateur de la matrice de covariance qui n'est pas défini positif et, pour cette raison, inutilisable.

4.2 Les méthodes de gradient

Nous voulons maximiser une fonction $\ell(\theta)$, de gradient $s(\theta) = \partial\ell(\theta)/\partial\theta$ et de hessien $H(\theta) = \partial^2\ell(\theta)/\partial\theta\partial\theta'$. Nous disposons également d'une valeur initiale notée $\theta_{(0)}$. De même, on note $\theta_{(p)}$ la valeur du paramètre

à la p -ième itération. Un algorithme du gradient est une règle d'itération de la forme suivante :

$$\theta_{(p+1)} = \theta_{(p)} + W_{(p)} s(\theta_{(p)}),$$

où $s(\theta_{(p)})$ est le gradient de la fonction que l'on cherche à maximiser et $W_{(p)}$ une matrice qui dépend de l'algorithme particulier que l'on emploie. On vérifie que lorsque l'on a atteint le maximum, $s(\hat{\theta}_n) = 0$, et le pas de l'itération est nul. Toutefois, dans la pratique, il peut arriver que le pas d'une itération soit trop fort et dépasse le point qui donne le maximum, on modifie donc la règle précédente en l'écrivant :

$$\theta_{(p+1)} = \theta_{(p)} + \lambda W_{(p)} s(\theta_{(p)}),$$

où $\lambda \in [0, 1]$. Il est également possible de prendre $\lambda > 1$ au début de l'algorithme pour accélérer la convergence. Le lecteur intéressé par ce dernier cas peut consulter Gouriéroux et Monfort (1989, chap. XIII). La valeur de λ n'est réduite que lorsque $\ell(\theta_{(p+1)}) < \ell(\theta_{(p)})$.

Les algorithmes de gradient sont très employés car ils possèdent une propriété intéressante : ils sont *croissants*. Cette propriété signifie que, si $W_{(p)}$ est symétrique et définie positive, alors pour de petits accroissements du pas de l'itération (i.e., λ petit), l'algorithme mène toujours à une valeur supérieure ou égale de l'objectif soit $\ell(\theta_{(p+1)}) \geq \ell(\theta_{(p)})$. Quand la fonction est concave, ceci garantit que l'on parvienne au maximum. Les trois algorithmes de gradient les plus utilisés sont ceux de Newton-Raphson, de Berndt-Hall-Hall-Hausman et du score.

4.2.1 Algorithme de Newton-Raphson

Il consiste à effectuer une approximation quadratique de la fonction à maximiser, en chacun des points de l'itération. Dans ce cas, si le hessien est défini négatif, on obtient le maximum de l'approximation par la condition du premier ordre sur une forme quadratique, dont on peut calculer facilement l'expression analytique parce qu'elle est linéaire. La succession de maxima ainsi obtenue donne la solution du problème. Le développement limité au second ordre de $\ell(\theta)$ au voisinage de $\theta_{(p)}$ est égal à :

$$\ell(\theta) \simeq \ell(\theta_{(p)}) + s'(\theta_{(p)}) (\theta - \theta_{(p)}) + \frac{1}{2} (\theta - \theta_{(p)})' H(\theta_{(p)}) (\theta - \theta_{(p)}).$$

La maximisation de cette forme quadratique par rapport à θ donne la condition du premier ordre :

$$s(\theta_{(p)}) + H(\theta_{(p)}) (\theta - \theta_{(p)}) = 0 \Leftrightarrow \theta - \theta_{(p)} = -H(\theta_{(p)})^{-1} s(\theta_{(p)}),$$

de plus la dérivée seconde est égale à $H(\theta_{(p)})$, qui est définie négative lorsque l'objectif est concave au point $\theta_{(p)}$. Dans ce cas, on a bien un maximum local donné par les conditions du premier ordre. Dans l'ensemble l'itération est donnée par :

$$\theta_{(p+1)} = \theta_{(p)} - \lambda H(\theta_{(p)})^{-1} s(\theta_{(p)}).$$

Cet algorithme présente généralement un pas assez fort dans les premières itérations et, dans l'ensemble, s'avère assez rapide. Tout dépend toutefois si l'on travaille sur des dérivées secondes analytiques ou numériques, car ces dernières augmentent fortement le temps de calcul à chaque itération.

4.2.2 Algorithme de Berndt-Hall-Hausman

Cet algorithme, justifié dans le cas du maximum de vraisemblance, se base sur l'égalité de la matrice d'information :

$$I_1(\theta) = \mathbb{E}_x \mathbb{E}_y \left[\frac{\partial \ln f(\theta)}{\partial \theta} \frac{\partial \ln f(\theta)}{\partial \theta'} \right] = \mathbb{E}_x \mathbb{E}_y \left[-\frac{\partial^2 \ln f(\theta)}{\partial \theta \partial \theta'} \right] = J_1(\theta).$$

Pour l'estimation, on remplace les moments théoriques par les moments empiriques correspondants, ce qui suggère d'approximer les dérivées secondes par l'opposé des produits croisés des dérivées premières. Pour un échantillon de taille N on peut écrire la fonction objectif et ses dérivées sous la forme :

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^N \ln f(y_i|x_i, \theta), \quad s(\theta) = \sum_{i=1}^N \frac{\partial \ln f(y_i|x_i, \theta)}{\partial \theta} \\ \text{et } H(\theta) &= \sum_{i=1}^N \frac{\partial^2 \ln f(y_i|x_i, \theta)}{\partial \theta \partial \theta'} \end{aligned}$$

La méthode de Berndt-Hall-Hausman consiste à approximer

$$H(\theta) \quad \text{par} \quad \sum_{i=1}^N \frac{\partial \ln f(y_i|x_i, \theta)}{\partial \theta} \frac{\partial \ln f(y_i|x_i, \theta)}{\partial \theta'},$$

ce qui donne la règle d'itération suivante :

$$\theta_{(p+1)} = \theta_{(p)} - \lambda \left[\sum_{i=1}^N \frac{\partial \ln f(y_i|x_i, \theta_{(p)})}{\partial \theta} \frac{\partial \ln f(y_i|x_i, \theta_{(p)})}{\partial \theta'} \right]^{-1} s(\theta_{(p)}).$$

Cet algorithme ne nécessite que les dérivées au premier ordre et est donc facile à mettre en oeuvre. Toutefois, il implique généralement plus

d'itérations que l'algorithme de Newton-Raphson. Son principal défaut est qu'il ne permet pas de vérifier la négativité du hessien à chaque itération, ce qui peut s'avérer problématique en certains points $\theta_{(p)}$ lorsque l'objectif n'est pas globalement concave.

4.2.3 Algorithme du score

Il s'agit d'un raffinement de l'algorithme de Berndt-Hall-Hall-Hausman où l'on prend l'espérance mathématique des produits croisés du score, égale à l'information de Fisher dans le cas du maximum de vraisemblance, à la place de leurs produits croisés. On approxime :

$$H(\theta) \quad \text{par} \quad \sum_{i=1}^N \mathbb{E}_y \left[\frac{\partial \ln f(y_i|x_i, \theta)}{\partial \theta} \frac{\partial \ln f(y_i|x_i, \theta)}{\partial \theta'} \right],$$

ce qui donne la règle d'itération :

$$\theta_{(p+1)} = \theta_{(p)} - \lambda \left[\sum_{i=1}^N \mathbb{E}_y \left[\frac{\partial \ln f(y_i|x_i, \theta_{(p)})}{\partial \theta} \frac{\partial \ln f(y_i|x_i, \theta_{(p)})}{\partial \theta'} \right] \right]^{-1} s(\theta_{(p)}).$$

Dans le cas où les dérivées secondes ne dépendent pas de la variable endogène y , cet algorithme est identique à celui de Newton-Raphson (e.g., cas des modèles Logit et de Poisson).

4.2.4 Algorithme de Levenberg-Marquardt

Il s'agit d'une extension de l'algorithme de Newton-Raphson que l'on applique quand l'objectif à maximiser n'est pas globalement concave. Supposons que le hessien au point $\theta_{(p)}$ n'est pas défini négatif, l'algorithme n'est plus nécessairement croissant et l'on n'est plus sûr de parvenir à un maximum local. On pourrait alors penser à utiliser l'algorithme du score ou celui de Berndt-Hall-Hall-Hausman. Mais la pratique montre que l'on aboutit souvent en fait à une valeur propre nulle du paramètre qui pose problème, car elle rend impossible l'inversion nécessaire à l'itération. On utilise donc une modification du hessien qui est définie négative, ce qui garantit que la matrice qui détermine le pas soit définie positive. Plus précisément, on utilise

$$W_p = - [H(\theta_{(p)}) - (1 + \alpha) \mu_H \mathbf{I}_k]^{-1}$$

où \mathbf{I}_k est la matrice identité de dimension k (la dimension de θ), μ_H la plus grande valeur propre du hessien (fonction disponible sous SAS-IML), supposée positive en cas de problème (i.e., de non concavité locale de la fonction objectif) et $\alpha > 0$ un paramètre choisi par l'utilisateur.

L'intuition de la méthode est la suivante : puisque le hessien n'est pas défini négatif, on le remplace par une matrice définie négative qui est la plus proche possible du hessien original. Pour cela il suffit de retrancher au hessien la matrice identité multipliée par la plus grande valeur propre du hessien (positive par hypothèse). Toutefois, ceci nous donnerait une valeur propre nulle. Il faut donc retrancher un peu plus que cette quantité; c'est ce que détermine le paramètre α . En général une petite valeur de α suffit et l'algorithme est sensible à de petites variations de ce paramètre (i.e, ce qui fonctionne avec $\alpha = 0.1$ peut ne plus fonctionner du tout avec $\alpha = 0.2$).

Pour déterminer la valeur de α , on utilise le constat suivant : si α est trop élevé, la plus grande valeur propre du hessien tend à croître avec le nombre d'itérations, si la valeur de α est trop petite, elle tend vers 0. Il suffit de se placer entre les deux après quelques essais, et l'on retombe généralement dans une zone de concavité de la fonction objectif. En général, plusieurs estimations sur un même échantillon ne nécessitent qu'une valeur de α , que l'on garde constante pendant tout l'algorithme.

Cet algorithme, qui sert au modèle tobit généralisé et aux moindres carrés non linéaires, deux exemples d'objectifs non concaves, mène à la règle d'itération suivante :

$$\theta_{(p+1)} = \begin{cases} \theta_{(p)} - \lambda [H(\theta_{(p)}) - (1 + \alpha)\mu_H I_k]^{-1} s(\theta_{(p)}) & \text{si } \mu_H \geq 0 \\ \theta_{(p)} - \lambda H(\theta_{(p)})^{-1} s(\theta_{(p)}) & \text{si } \mu_H < 0 \end{cases} .$$

4.3 Méthodologie de programmation

Comment s'assurer que le programme que l'on a écrit ne comporte aucune erreur et qu'il soit pratique à utiliser? En prenant un certain nombre de précautions présentées dans cette section. Il y a quatre étapes qu'il faut prendre soin de bien effectuer.

1. *Vérifier la fonction objectif. A la fois par le calcul, mais également en consultant les ouvrages et les articles qui la donnent.*
2. *Vérifier le gradient. Par la même méthode que précédemment, mais également numériquement. Ainsi, on peut détecter des erreurs de recopie aussi bien sur la fonction objectif que sur son gradient. A cette étape, on utilise un algorithme basé sur les dérivées premières, de type Berndt-Hall-Hall-Hausman.*
3. *Vérifier le hessien. On utilise le gradient analytique pour calculer le hessien numérique, afin d'éviter le cumul des erreurs d'approximation. Ce problème est particulièrement sensible ici car les dérivées sont calculées à partir de quantités très petites par définition. On utilise*

un algorithme de Newton-Raphson ou de Levenberg-Marquardt. On peut conserver les dérivées secondes numériques si le calcul des dérivées analytiques est trop complexe, ou pour obtenir de premières estimations.

4. *Paramétrer le programme définitif. Ceci vise à éviter toute intervention sur le programme une fois qu'il a été vérifié, car c'est une source d'erreur potentielle. Par exemple, on peut paramétrer des programmes écrits en SAS-IML par des macro-variables, ce qui permet d'écrire des routines appelées macro-commandes. On peut alors utiliser ces routines sur toutes les bases de données et quel que soit le nombre de variables explicatives.*

CHAPITRE 5

Les variables dichotomiques

5.1 Cas général

On observe une variable dichotomique $y \in \{0, 1\}$. Par définition cette variable suit une loi de Bernoulli. À partir d'un modèle théorique on peut écrire que le paramètre de cette loi est égal à p . La probabilité d'une loi de Bernoulli de paramètre p peut s'écrire :

$$\Pr [y = k] = p^k (1 - p)^{1-k} = \begin{cases} p^0 (1 - p)^{1-0} = 1 - p & \text{si } k = 0 \\ p^1 (1 - p)^{1-1} = p & \text{si } k = 1 \end{cases}$$

et son logarithme est égal à :

$$\ln \Pr [y = k] = k \ln p + (1 - k) \ln (1 - p).$$

On remarque également que l'espérance de y est égale à :

$$E(y) = 0 \times (1 - p) + 1 \times p = p,$$

et que sa variance est égale à :

$$\begin{aligned} V(y) &= E\left((y - p)^2\right) \\ &= (0 - p)^2 \times (1 - p) + (1 - p)^2 \times p \\ &= p(1 - p)(p + 1 - p) \\ &= p(1 - p). \end{aligned}$$

On considère maintenant un échantillon de N variables aléatoires de Bernoulli (y_1, \dots, y_N) , indépendantes, et de paramètres (p_1, \dots, p_N) . Les paramètres des N lois sont différents parce que l'on considère un modèle conditionnel où chaque variable y_i possède une probabilité (i.e., une espérance conditionnelle) qui lui est propre. Ces probabilités dépendent donc de variables explicatives, regroupées dans le vecteur X , et dont les réalisations sont notées (X_1, \dots, X_N) . Pour bien montrer le caractère conditionnel du modèle, on pose :

$$p_i = p(X_i, \beta), \quad i = 1, \dots, N$$

où β est le paramètre que l'on cherche à estimer. La log-vraisemblance d'un échantillon (y_1, \dots, y_N) est donc égale à :

$$\ell(y|X, \beta) = \sum_{i=1}^N y_i \ln p(X_i, \beta) + (1 - y_i) \ln(1 - p(X_i, \beta)).$$

Le score est donc égal à :

$$\frac{\partial \ell}{\partial \beta}(y|X, \beta) = \sum_{i=1}^N \frac{\partial p_i}{\partial \beta} \frac{y_i - p_i}{p_i(1 - p_i)},$$

et le hessien à :

$$\begin{aligned} \frac{\partial^2 \ell}{\partial \beta \partial \beta'}(y|X, \beta) &= \sum_{i=1}^N \frac{\partial^2 p_i}{\partial \beta \partial \beta'} \frac{y_i - p_i}{p_i(1 - p_i)} \\ &\quad - \frac{\partial p_i}{\partial \beta} \frac{\partial p_i}{\partial \beta'} \left[\frac{(y_i - p_i)(1 - 2p_i)}{p_i^2(1 - p_i)^2} + \frac{1}{p_i(1 - p_i)} \right]. \end{aligned}$$

Pour appliquer l'algorithme du score, on remarque simplement que :

$$E(y_i|X) = p_i,$$

ce qui permet d'obtenir :

$$E_y \left[\frac{\partial^2 \ell}{\partial \beta \partial \beta'}(y|X, \beta) \right] = - \sum_{i=1}^N \frac{\partial p_i}{\partial \beta} \frac{\partial p_i}{\partial \beta'} \frac{1}{p_i(1 - p_i)}.$$

Cette représentation est valable pour tous les modèles dichotomiques; on a juste besoin de l'expression de la probabilité de réalisation d'un événement en fonction des variables explicatives.

Pour obtenir les modèles usuels, on suppose que la variable qualitative que l'on observe résulte d'un modèle latent qui porte sur une variable

continue, notée y_i^* . Cette variable inobservable est supposée décrite par un modèle linéaire standard donné par :

$$y_i^* = X_i\beta + u_i,$$

où u_i est une perturbation d'espérance nulle, sans perte de généralité tant que le modèle latent contient un terme constant. Pour pouvoir estimer ce modèle par le maximum de vraisemblance, il nous faut écrire la loi de la variable observable conditionnellement aux variables explicatives. Cette variable observable est définie par :

$$y_i = \begin{cases} 1 & \text{si } y_i^* > 0 \\ 0 & \text{si } y_i^* \leq 0 \end{cases}$$

Le fait de prendre la valeur 0 comme seuil de référence n'a aucune incidence sur les estimations tant que le modèle comporte un terme constant, car on peut alors utiliser la variable latente $y_i^* - c$, où c est le seuil de troncature constant. La loi suivie par y_i est une loi de Bernoulli de paramètre $p_i = \Pr[y_i^* > 0]$, mais contrairement au cas habituel rencontré en statistique (i.e., modèle marginal), la probabilité est différente pour chaque observation puisqu'elle dépend des variables explicatives (i.e., modèle conditionnel). Le paramètre de la loi de Bernoulli est défini par :

$$p_i = \Pr[y_i^* > 0] = \Pr[X_i\beta + u_i > 0] = \Pr[u_i > -X_i\beta] = 1 - F(-X_i\beta),$$

où F est la fonction de répartition des u_i , $i = 1, \dots, N$. La probabilité d'observer une réalisation y est donc donnée par :

$$\Pr[y_i = y] = p_i^y (1 - p_i)^{1-y}, \quad y \in \{0, 1\}.$$

De plus, on a $E(y_i|X_i, \beta) = p_i$ et $V(y_i|X_i, \beta) = p_i(1 - p_i)$. Comme les N observations sont supposées indépendantes, la vraisemblance de l'échantillon est donnée par le produit des probabilités individuelles :

$$L(y|X, \beta) = \prod_{i=1}^N p_i^{y_i} (1 - p_i)^{1-y_i},$$

d'où la log-vraisemblance :

$$\begin{aligned} \ell(y|X, \beta) &= \sum_{i=1}^N y_i \ln(p_i) + (1 - y_i) \ln(1 - p_i) \\ &= \sum_{i=1}^N y_i \ln[1 - F(-X_i\beta)] + (1 - y_i) \ln F(-X_i\beta). \end{aligned}$$

La forme spécifique prise par la probabilité p_i dépend directement de l'hypothèse faite sur la loi de la perturbation u . On peut essayer des

distributions différentes lors d'une étude, puis les comparer en effectuant des tests non emboîtés. Les deux lois les plus utilisées en pratique sont les lois normale et logistique. La première définit le modèle Probit, parfois appelé modèle Normit, tandis que la seconde définit le modèle Logit. Ces deux lois sont symétriques, on a donc $F(-z) = 1 - F(z)$, d'où :

$$p_i = F(X_i\beta).$$

Dans le cas du modèle Logit :

$$F(z) = \frac{1}{1 + \exp(-z)}.$$

Dans le cas du modèle Probit (ou Normit) :

$$F(z) = \Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{s^2}{2}\right) ds.$$

On remarque que dans tous les cas, la probabilité p_i est une fonction de la quantité réelle $m_i = X_i\beta$. Cette propriété résulte de la linéarité du modèle latent.

5.2 Le modèle Logit

Dans ce modèle on suppose que la perturbation u suit une loi logistique de fonction de répartition :

$$F(z) = \frac{1}{1 + \exp(-z)}.$$

On vérifie directement que $F(-z) = 1 - F(z)$, la loi est symétrique. Sa densité est donnée par :

$$\begin{aligned} f(z) &= \partial F(z) / \partial z \\ &= \frac{\exp(-z)}{[1 + \exp(-z)]^2} \\ &= \frac{1}{1 + \exp(-z)} \times \frac{\exp(-z)}{1 + \exp(-z)} \\ &= F(z) [1 - F(z)]. \end{aligned}$$

Cette dernière propriété permet de simplifier toutes les dérivées de la log-vraisemblance. Les moments de cette loi logistique sont donnés par :

$$E(u) = 0, \quad V(u) = \frac{\pi^2}{3},$$

elle est donc centrée, mais n'est pas réduite. Il faudra en tenir compte lors de la comparaison des coefficients des modèles Logit et Probit car, dans ce dernier cas, la loi est à la fois centrée et réduite.

Les observations sont indépendantes donc la log-vraisemblance de l'échantillon (y_1, \dots, y_n) est donnée par :

$$\ell(y|X, \beta) = \sum_{i=1}^n y_i \ln F(m_i) + (1 - y_i) \ln [1 - F(m_i)],$$

avec $m_i = X_i \beta$. Le vecteur du score est donné par :

$$\begin{aligned} s(y|X, \beta) &= \frac{\partial \ell}{\partial \beta}(y|X, \beta) \\ &= \sum_{i=1}^n X_i' \left[y_i \frac{f(m_i)}{F(m_i)} - (1 - y_i) \frac{f(m_i)}{1 - F(m_i)} \right] \\ &= \sum_{i=1}^n X_i' [y_i [1 - F(m_i)] - (1 - y_i) F(m_i)] \\ &= \sum_{i=1}^n X_i' [y_i - F(m_i)]. \end{aligned}$$

On vérifie que son espérance est nulle puisque $E(y_i|X_i, \beta) = F(m_i)$. Le hessien est donné par :

$$\begin{aligned} H(y|X, \beta) &= \frac{\partial^2 \ell}{\partial b \partial b'}(y|X, \beta) \\ &= \frac{\partial s}{\partial b'}(y|X, \beta) \\ &= - \sum_{i=1}^n X_i' X_i f(m_i). \end{aligned}$$

Comme f est une fonction strictement positive, le hessien est bien défini négatif. L'algorithme de Newton-Raphson sera donc croissant. On peut remarquer également que l'information de Fisher du modèle est donnée par :

$$I(y|X, \beta) = E[-H(y|X, \beta) | X, \beta] = \sum_{i=1}^n X_i' X_i f(m_i).$$

On peut également calculer cette quantité en utilisant la variance du

score :

$$\begin{aligned} V[s(y|X, \beta) | X, \beta] &= \sum_{i=1}^n X_i' X_i V[y_i - F(m_i) | X, \beta] \\ &= \sum_{i=1}^n X_i' X_i F(m_i) (1 - F(m_i)) \\ &= \sum_{i=1}^n X_i' X_i f(m_i). \end{aligned}$$

Donc l'algorithme du score est identique à celui de Newton Raphson. En conséquence, nous utiliserons seulement deux algorithmes pour estimer ce modèle : Berndt-Hall-Hausman et Newton-Raphson dont les matrices sont données respectivement par :

$$\begin{aligned} W_{B3H}^{-1} &= - \sum_{i=1}^n X_i' X_i [y_i - F(m_i)]^2, \\ W_{NR}^{-1} &= - \sum_{i=1}^n X_i' X_i f(m_i). \end{aligned}$$

5.3 Le modèle Probit (ou Normit)

Maintenant on suppose que la perturbation u suit une loi normale de fonction de répartition donnée par :

$$F(z) = \Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{s^2}{2}\right) ds.$$

Comme tous les programmes d'estimation possèdent la fonction de répartition de la loi normale centrée-réduite, nous garderons la notation

$\Phi(z)$. La densité est donnée par :¹

$$\begin{aligned}\varphi(z) &= \frac{d}{dz} \left[\int_{-\infty}^z \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{s^2}{2}\right) ds \right] \\ &= 1 \times \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) - 0 \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right).\end{aligned}$$

En dérivant cette densité, on voit directement que :

$$\varphi'(z) = -z \times \varphi(z).$$

Cette dernière relation permet de simplifier l'écriture du hessien et plus généralement, de trouver facilement les dérivées d'un ordre quelconque. Les moments de cette loi normale sont donnés par :

$$E(u) = 0, \quad V(u) = 1.$$

Les observations sont indépendantes donc la log-vraisemblance de l'échantillon $y = (y_1, \dots, y_n)$ est donnée par :

$$\ell(y|X, \beta) = \sum_{i=1}^N y_i \ln \Phi(m_i) + (1 - y_i) \ln [1 - \Phi(m_i)].$$

Le vecteur du score est donné par :

$$\begin{aligned}s(y|X, \beta) &= \frac{\partial \ell}{\partial \beta}(y|X, \beta) \\ &= \sum_{i=1}^N X_i' \left[y_i \frac{\varphi(m_i)}{\Phi(m_i)} - (1 - y_i) \frac{\varphi(m_i)}{1 - \Phi(m_i)} \right] \\ &= \sum_{i=1}^N X_i' \frac{\varphi(m_i)(y_i - \Phi(m_i))}{\Phi(m_i)(1 - \Phi(m_i))}.\end{aligned}$$

On vérifie que son espérance est nulle puisque $E(y_i|X_i, \beta) = \Phi(m_i)$. Pour simplifier les notations on pose $\varphi(m_i) = \varphi_i$ et $\Phi(m_i) = \Phi_i$. En

¹On applique le théorème de Leibniz :

$$\frac{d}{dz} \int_{a(z)}^{b(z)} f(t) dt = b'(z) f(b(z)) - a'(z) f(a(z)).$$

utilisant $\varphi'_i = -m_i \times \varphi_i$, le hessien est donné par :

$$\begin{aligned} H(y|X, \beta) &= \frac{\partial^2 \ell}{\partial \beta \partial \beta'} (y|X, \beta) \\ &= \frac{\partial s}{\partial \beta'} (y|X, \beta) \\ &= - \sum_{i=1}^N X'_i X_i \left\{ \frac{\varphi_i (y_i - \Phi_i)}{\Phi_i (1 - \Phi_i)} \left[m_i + \frac{\varphi_i (1 - 2\Phi_i)}{\Phi_i (1 - \Phi_i)} \right] - \frac{\varphi_i^2}{\Phi_i (1 - \Phi_i)} \right\}. \end{aligned}$$

On en déduit que l'information de Fisher du modèle est donnée par :

$$\begin{aligned} I_1(\beta) &= \mathbb{E}_X \mathbb{E}_y [-H(y|X, \beta) | X, \beta] \\ &= \mathbb{E}_X \mathbb{V}_y [s(y|X, \beta) | X, \beta] \\ &= \frac{1}{N} \sum_{i=1}^N X'_i X_i \frac{\varphi_i^2}{\Phi_i (1 - \Phi_i)}. \end{aligned}$$

On peut également calculer cette quantité en utilisant la variance du score, ce qui est plus pratique ici que de calculer la dérivée seconde. Clairement, l'algorithme du score est différent de l'algorithme de Newton-Raphson. Il nécessite moins de calcul et est également croissant.

On peut donc utiliser trois algorithmes pour estimer ce modèle : Berndt-Hall-Hausman, le Score et Newton-Raphson dont les matrices sont données respectivement par :

$$\begin{aligned} W_{BH}^{-1} &= - \sum_{i=1}^N X'_i X_i \frac{\varphi_i^2}{\Phi_i^2 (1 - \Phi_i)^2} [y_i - \Phi_i]^2, \\ W_{SC}^{-1} &= - \sum_{i=1}^N X'_i X_i \frac{\varphi_i^2}{\Phi_i (1 - \Phi_i)}, \\ W_{NR}^{-1} &= \sum_{i=1}^N X'_i X_i \left\{ \frac{\varphi_i (y_i - \Phi_i)}{\Phi_i (1 - \Phi_i)} \left[m_i + \frac{\varphi_i (1 - 2\Phi_i)}{\Phi_i (1 - \Phi_i)} \right] - \frac{\varphi_i^2}{\Phi_i (1 - \Phi_i)} \right\}. \end{aligned}$$

5.4 Interprétation et comparaison des coefficients

5.4.1 Le modèle Probit

En fait, les coefficients du modèle latent ne sont estimés qu'à une constante multiplicative et positive près : l'inverse de l'écart-type de la perturbation du modèle latent. Dans le cas du modèle Probit, le modèle

latent s'écrit :

$$z_i^* = X_i b + v_i, \quad v_i \overset{iid}{\rightsquigarrow} N(0, \sigma^2).$$

Le modèle que nous avons estimé s'écrit donc :

$$y_i^* = X_i \beta + u_i, \quad u_i \overset{iid}{\rightsquigarrow} N(0, 1)$$

avec $y_i^* = z_i^*/\sigma$, $\beta = b/\sigma$ et $u_i = v_i/\sigma$. Les paramètres b et σ ne sont pas identifiables, seule la fonction $\beta = b/\sigma$ de ces deux paramètres peut être estimée, ou toute fonction monotone de β . En conséquence, tous les coefficients d'un modèle Probit sont implicitement réduits par l'écart-type de la perturbation de la régression, ce qui a un certain nombre de conséquences sur leur interprétation :

1. On ne peut pas comparer les coefficients obtenus sur les régressions de deux variables dichotomiques endogènes différentes, car l'écart-type de la perturbation change avec le modèle latent.
2. Le signe du coefficient β est le même que celui de b car un écart-type est toujours positif.
3. Le ratio de deux coefficients extraits de β est identique au ratio des deux coefficients correspondants de b . On peut donc dire qu'un coefficient est deux fois plus grand qu'un autre.
4. La différence entre deux coefficients extraits de β n'est connue qu'à un facteur multiplicatif positif près, égal à σ^{-1} . On ne peut donc interpréter que le signe de la différence entre deux coefficients, pas la grandeur de l'écart. Par contre, on peut comparer deux écarts tirés de la même régression.

5.4.2 Le modèle Logit

La même interprétation reste valable, à ceci près qu'une loi logistique de paramètres $(0, \phi)$ admet pour espérance 0 et pour variance $\phi^2 \pi^2/3$. On peut donc écrire le modèle latent :²

$$z_i^* = X_i b + v_i, \quad v_i \overset{iid}{\rightsquigarrow} \Lambda(0, \phi)$$

²La fonction de répartition variable Z suivant une loi logistique de paramètres $\Lambda(\mu, \phi)$ est donnée par :

$$F(z) = \left[1 + \exp\left(-\frac{z - \mu}{\phi}\right) \right]^{-1}.$$

On a :

$$E(Z) = \mu \text{ et } V(Z) = \frac{\phi^2 \pi^2}{3}.$$

Le modèle que nous avons estimé s'écrit donc :

$$y_i^* = X_i\beta + u_i, \quad u_i \overset{iid}{\rightsquigarrow} \Lambda(0, 1)$$

avec $y_i^* = z_i^*/\phi$, $\beta = b/\phi$ et $u_i = v_i/\phi$. L'interprétation des coefficients du modèle Logit se fait donc de la même façon que celle du modèle Probit.

5.4.3 Comparaison des coefficients des modèles Logit et Probit

Dans le modèle Probit, on estime :

$$\beta_{\text{PROBIT}} = \frac{b}{\sqrt{V(u)}},$$

alors que dans le modèle Logit on estime :

$$\beta_{\text{LOGIT}} = \frac{b}{\phi} \text{ avec } \phi = \frac{\sqrt{3}}{\pi} \sqrt{V(u)},$$

on en déduit que :

$$\beta_{\text{PROBIT}} \simeq \frac{\sqrt{3}}{\pi} \beta_{\text{LOGIT}}.$$

Dans la pratique on utilisera donc les approximations suivantes :

$$\beta_{\text{PROBIT}} \simeq 0,55 \times \beta_{\text{LOGIT}} \text{ et } \beta_{\text{LOGIT}} \simeq 1,81 \times \beta_{\text{PROBIT}}.$$

On peut imprimer systématiquement les coefficients réduits en fin de programme pour faciliter les comparaisons. La même modification doit être effectuée sur les écarts-types asymptotiques des estimateurs. Il n'en demeure pas moins que les deux modèles sont différents et que la comparaison n'est qu'approximative car la loi logistique n'est qu'une approximation de la loi normale et admet notamment plus de valeurs extrêmes que cette dernière.³

5.5 Les aides à l'interprétation

Les coefficients des modèles Logit et Probit ne sont définis qu'à une constante multiplicative près, de sorte qu'ils ne sont pas directement interprétables. La méthode la plus simple pour obtenir des coefficients directement interprétables consiste à calculer l'impact d'une variables explicative directement sur la probabilité. Le cas le plus simple est celui

³Le coefficient d'applatissage de la loi logistique est de 1,2 au lieu de 1 pour la loi normale.

où la variable explicative est binaire; il suffit de comparer les deux états $\{0, 1\}$. Dans le cas d'une variables explicative quantitative il faut prendre deux points de référence; par exemple en comparant l'effet du passage du premier au troisième quartile. Deux types de mesures sont utilisées dans la littérature : d'une part, l'effet direct de la variable explicative sur la probabilité ou effet incrémental; d'autre part, l'effet d'une variable explicative sur le ratios des probabilités de réalisation de l'évènement (ou "odds ratio").

5.5.1 Variables explicatives binaires

Sans perte de généralité, considérons le cas d'une seule variable explicative binaire, $X \in \{0, 1\}$. La fonction suivante donne les chances que l'évènement $Y = 1$ se réalise par rapport à l'évènement $Y = 0$:

$$R(X) = \frac{\Pr[Y = 1|X]}{\Pr[Y = 0|X]} = \frac{F(\beta_0 + \beta_1 X)}{1 - F(\beta_0 + \beta_1 X)},$$

il s'agit du rapport des probabilités pour une même valeur de X ("the odds function"). Pour voir l'effet de X sur ce ratio, on utilise le "ratio des cotes" ("odds ratio") :

$$\psi_X = \frac{R(1)}{R(0)},$$

soit :

$$\psi_X = \frac{\frac{F(\beta_0 + \beta_1)}{1 - F(\beta_0 + \beta_1)}}{\frac{F(\beta_0)}{1 - F(\beta_0)}},$$

ce ratio indique la modification des chances d'obtenir l'évènement $Y = 1$ lorsque l'on passe du sous échantillon $X = 0$ au sous échantillon $X = 1$. Le coefficient $\beta = (\beta_0, \beta_1)'$ est généralement estimé avec de nombreuses autres variables de sorte qu'il s'agit d'un effet "toutes choses égales par ailleurs".

Si le modèle comporte plusieurs variables explicatives, on estime généralement un effet au point moyen. Dans ce cas le modèle avec les autres variables explicatives s'écrit :

$$E(Y) = F(\gamma_0 + \beta_1 X + Z\beta_2),$$

de sorte que le ratio des cotes peut s'écrire :

$$\psi_X = \frac{\frac{F(\gamma_0 + \beta_1 + \bar{Z}\beta_2)}{1 - F(\gamma_0 + \beta_1 + \bar{Z}\beta_2)}}{\frac{F(\gamma_0 + \bar{Z}\beta_2)}{1 - F(\gamma_0 + \bar{Z}\beta_2)}},$$

en posant :

$$\beta_0 = \gamma_0 + \bar{Z}\beta_2,$$

on retrouve la même formule que précédemment.

Le cas du modèle Probit s'obtient avec $F = \Phi$. Pour le modèle Logit l'expression se simplifie nettement :

$$\Pr [Y = 1|X] = \frac{\exp(\beta_0 + \beta_1 X)}{1 + \exp(\beta_0 + \beta_1 X)},$$

et

$$\Pr [Y = 0|X] = \frac{1}{1 + \exp(\beta_0 + \beta_1 X)},$$

de sorte que le rapport des probabilités est égal à :

$$R(X) = \exp(\beta_0 + \beta_1 X),$$

et que le rapport des cotes est égal à :

$$\psi = \frac{R(1)}{R(0)} = \frac{\exp(\beta_0 + \beta_1)}{\exp(\beta_0)} = \exp(\beta_1),$$

il suffit de prendre l'exponentielle de ce coefficient pour obtenir l'effet de la variable X sur le rapport des cotes. Ceci permet également de calculer un intervalle de confiance facilement en utilisant le théorème de Slutsky. Il ne s'agit toutefois pas de l'effet de la variable X sur $\Pr [Y = 1|X]$ mais sur le ratio $\Pr [Y = 1|X] / \Pr [Y = 0|X]$. Pour obtenir l'effet incrémental de X sur la probabilité que l'évènement $Y = 1$ se réalise, on doit calculer :

$$\Delta_X = \frac{\Pr [Y = 1|X = 1]}{\Pr [Y = 1|X = 0]} = \frac{F(\beta_0 + \beta_1)}{F(\beta_0)},$$

dans le cas du modèle Logit, on obtient donc :

$$\begin{aligned} \Delta_X &= \frac{\frac{\exp(\beta_0 + \beta_1)}{1 + \exp(\beta_0 + \beta_1)}}{\frac{\exp(\beta_0)}{1 + \exp(\beta_0)}} \\ &= \exp(\beta_1) \times \frac{1 + \exp(\beta_0)}{1 + \exp(\beta_0 + \beta_1)} \\ &= \frac{\exp(\beta_1) + \exp(\beta_0 + \beta_1)}{1 + \exp(\beta_0 + \beta_1)}, \end{aligned}$$

on voit que si $\beta_1 > 0$, alors $\exp(\beta_1) > 1$ et $\Delta_X > 1$.

5.5.2 Variables explicatives quantitatives

Avec des variables quantitatives, il est possible de calculer deux types de quantités : l'effet marginal qui donne la variation de la probabilité que $Y = 1$ and X augmente d'une unité; ou l'effet incrémental, qui mesure la variation de la probabilité que $Y = 1$ quand X passe d'une valeur de référence à une autre. Comme valeurs de référence, on peut prendre les quartiles ou des déciles extrêmes de X . La variation interdécile donne une bonne idée du potentiel d'influence de la variable X , puisqu'elle représente la variation de la probabilité quand on passe des 10% plus petites valeurs de X à ses 10% les plus élevées. Dans le cas d'une loi normale ceci serait équivalent à calculer la variation de la probabilité quand la variable X se situe entre $\bar{X} \pm 1.645 \times \sigma_X$.

L'effet marginal est donné par :

$$\delta_X = \frac{\partial}{\partial X} \Pr[Y = 1|X] = \beta_1 \times f(\beta_0 + \beta_1 X),$$

et il varie avec chaque valeur de X . On peut soit faire un graphique sur l'ensemble des valeurs de X , soit prendre un point de référence comme la moyenne ou la médiane. L'effet incrémental du passage de X de a à b est analogue au traitement sur variables qualitatives. On a :

$$\Delta_X = \frac{\Pr[Y = 1|X = b]}{\Pr[Y = 1|X = a]} = \frac{F(\beta_0 + \beta_1 b)}{F(\beta_0 + \beta_1 a)},$$

dans le cas du modèle Logit :

$$\Delta_X = \frac{\frac{\exp(\beta_0 + \beta_1 b)}{1 + \exp(\beta_0 + \beta_1 b)}}{\frac{\exp(\beta_0 + \beta_1 a)}{1 + \exp(\beta_0 + \beta_1 a)}} = \exp[\beta_1 (b - a)] \times \frac{1 + \exp(\beta_0 + \beta_1 a)}{1 + \exp(\beta_0 + \beta_1 b)},$$

dans ce cas, on remarque que le ratio des cotes se simplifie de manière importante avec le modèle Logit :

$$\psi_X = \frac{R(b)}{R(a)} = \frac{\exp(\beta_0 + \beta_1 b)}{\exp(\beta_0 + \beta_1 a)} = \exp[\beta_1 (b - a)] = \exp(\beta_1)^{b-a},$$

cette formule simplifie également le calcul des intervalles de confiance. On peut notamment utiliser les rapports des cotes pour calculer l'effet incrémental d'une variable quantitative et pas seulement d'une variable qualitative.

5.6 Application : la participation des femmes au marché du travail

Cette section présente une version simplifiée d'une équation de participation au marché du travail. La variable que l'on cherche à expliquer est dichotomique : une personne a un emploi ou non au moment de l'enquête. Les données sont issues de l'enquête "Jeunes et Carrières" réalisée par l'INSEE en 1997. L'échantillon comprend des données sur $N = 5425$ couples. On considère la participation des femmes au marché du travail que l'on explique par les déterminants suivants :

1. Age;
2. Nombre d'enfants;
3. Naissance l'année courante;
4. Activité des parents;
5. Nationalité de l'intéressée et de ses parents;
6. Niveau d'éducation (1 = sans diplôme ou première année de CAP; 2 = sans diplôme ou dernière année de CAP; 3 = CAP ou BEP; 4 = Baccalauréat professionnel; 5 = Baccalauréat général ou équivalent; 6 = BTS; 7 = Enseignement supérieur général).
7. Région d'habitation;
8. Les mêmes variables pour le conjoint;

L'estimation est réalisée sous SAS à partir de la procédure *logistic*. La syntaxe de base, si l'on veut expliquer une variable dichotomique y par les variables explicatives x_1, x_2 et x_3 , est la suivante :

SYNTAXE 5.1

```
proc logistic data=tab descending;
model y=x1 x2 x3;
run;
```

Sous cette forme la procédure va chercher les données dans le tableau *tab*. L'option *descending* est très importante car, par défaut, la procédure estime un modèle dichotomique avec une probabilité $p(X_i, \beta) = \Pr[y_i = 0]$ au lieu de $\Pr[y_i = 1]$. Dans le cas du modèle Logit, on a $\Pr[y_i = 0] = F(-X_i\beta)$ de sorte que sans l'option *descending* on estime $-\beta$ au lieu de β . L'option *descending* permet donc d'imposer $p(X_i, \beta) = \Pr[y_i = 1]$. L'instruction *model* sert à indiquer la variable (dichotomique) expliquée y et la liste des variables explicatives X_1, X_2 et X_3 . Par défaut, la procédure *logistic* permet d'estimer le modèle Logit, mais d'autres

modèles sont disponibles. Par exemple, pour estimer un modèle Probit, on utilise la syntaxe suivante :

SYNTAXE 5.2

```
proc logistic data=tab descending;
model y=x1 x2 x3/link=normit;
run;
```

L'option *link=* permet d'indiquer la distribution que l'on souhaite. Le modèle Normit correspond à la fonction de répartition de la loi normale. Pour obtenir le modèle de Weibull, on remplace l'option *link=normit* par *link=cloglog*.

Reprenons notre application. Afin d'estimer le modèle, on entre les commandes :

Programme 5.1

```
proc logistic descending data=tab;
model f_jc97=
f_age f_enf1 f_enf2 f_enf3 f_enf4 f_nai9697
f_mcspmiss f_magricul f_martisan f_mcadre
f_mprofint /*f_employe*/ f_mouvrier f_minactiv
f_pcspmiss f_pagricul f_partisan f_pcadre
f_pprofint f_pemploye /*f_pouvrier*/
/*f_francais*/ f_afrnord f_europ f_autrenat
/*f_pfrancai*/ f_pafrnord f_peurop f_pautrnat
/*f_ndip1*/ f_ndip2 f_ndip3 f_ndip4 f_ndip5 f_ndip6
f_ndip7
h_age
h_mcspmiss h_magricul h_martisan h_mcadre h_mprofint
/*h_employe*/ h_mouvrier h_minactiv
h_pcspmiss h_pagricul h_partisan h_pcadre h_pprofint
h_pemploye /*h_pouvrier*/
/*h_francais*/ h_afrnord h_europ h_autrenat
/*h_pfrancai*/ h_pafrnord h_peurop h_pautrnat
/*h_ndip1*/ h_ndip2 h_ndip3 h_ndip4 h_ndip5 h_ndip6
h_ndip7
/*h_ILEFR97*/ h_CHAMPA97 h_PICARD97 h_HAUTNO97
h_CENTRE97 h_BASSNO97 h_BOURGO97 h_NORDP97
h_LORRAI97 h_ALSA97 h_FRCOM97 h_PAYSLO97
h_BRETA97 h_POITOU97 h_AQUITA97 h_MIDIPY97
h_LIMOUS97 h_RHONEA97 h_AUVER97 h_LANGUE97
h_PROVEN97
```

```

/link=normit;
run;

```

les variables entre commentaires (début /* et fin */) indiquent la modalité de référence. On prend généralement la modalité la plus répandue. Le programme précédent produit la sortie :

Sortie 5.1

The LOGISTIC Procedure

Model Information

Data Set	WORK.TAB
Response Variable	f_jc97
Number of Response Levels	2
Number of Observations	5425
Model	binary probit
Optimization Technique	Fisher's scoring

Response Profile

Ordered Value	f_jc97	Total Frequency
1	1	3701
2	0	1724

Probability modeled is f_jc97=1.

Model Convergence Status

Convergence criterion (GCONV=1E-8) satisfied.

Model Fit Statistics

Criterion	Intercept Only	Intercept and Covariates
-----------	----------------	--------------------------

AIC	6785.321	5798.470
SC	6791.920	6313.174
-2 Log L	6783.321	5642.470

Testing Global Null Hypothesis: BETA=0

Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	1140.8513	77	<.0001
Score	1074.0382	77	<.0001
Wald	925.8740	77	<.0001

The LOGISTIC Procedure

Analysis of Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq
Intercept	1	-0.9183	0.1678	29.9387	<.0001
f_age	1	0.0522	0.00587	79.0265	<.0001
f_enf1	1	-0.2213	0.0665	11.0658	0.0009
f_enf2	1	-0.5780	0.0670	74.5248	<.0001
f_enf3	1	-0.9756	0.0777	157.8460	<.0001
f_enf4	1	-1.6081	0.1101	213.2488	<.0001
f_nai9697	1	-0.2890	0.0622	21.6082	<.0001
F_MCSPMISS	1	-0.2903	0.1115	6.7828	0.0092
F_MAGRICUL	1	0.1182	0.1154	1.0491	0.3057
F_MARTISAN	1	-0.1044	0.0941	1.2298	0.2674
F_MCADRE	1	-0.2240	0.1888	1.4089	0.2352
F_MPROFINT	1	-0.00686	0.0893	0.0059	0.9388
F_MOUVRIER	1	-0.0702	0.0596	1.3843	0.2394
F_MINACTIV	1	-0.1264	0.0521	5.8969	0.0152
F_PCSPMISS	1	-0.1194	0.1081	1.2211	0.2691
F_PAGRICUL	1	0.0425	0.0970	0.1916	0.6616
F_PARTISAN	1	0.00175	0.0727	0.0006	0.9808
F_PCADRE	1	-0.2211	0.0877	6.3645	0.0116
F_PPROFINT	1	-0.0366	0.0686	0.2848	0.5936
F_PEMPLOYE	1	-0.0746	0.0609	1.5025	0.2203
F_AFRNORD	1	-0.7277	0.2178	11.1613	0.0008
F_EUROP	1	-0.0673	0.1735	0.1506	0.6980
F_AUTRENAT	1	-0.7775	0.1747	19.7985	<.0001
F_PAFRNORD	1	-0.1896	0.1184	2.5623	0.1094

F_PEUROP	1	0.00814	0.0972	0.0070	0.9333
F_PAUTRNAT	1	-0.00406	0.0705	0.0033	0.9541
F_NDIP2	1	0.2343	0.0668	12.3118	0.0005
F_NDIP3	1	0.3749	0.0586	40.9611	<.0001
F_NDIP4	1	0.6551	0.0792	68.3899	<.0001
F_NDIP5	1	0.5770	0.0867	44.3076	<.0001
F_NDIP6	1	0.9994	0.0873	131.1030	<.0001
F_NDIP7	1	0.9218	0.0955	93.1190	<.0001
h_age	1	0.00383	0.00587	0.4265	0.5137
H_MCSPMISS	1	-0.1559	0.1092	2.0401	0.1532
H_MAGRICUL	1	-0.1517	0.1079	1.9778	0.1596
H_MARTISAN	1	-0.0929	0.0938	0.9817	0.3218
H_MCADRE	1	0.00537	0.1960	0.0008	0.9781
H_MPROFINT	1	0.0275	0.0887	0.0960	0.7566
H_MOUVRIER	1	0.00255	0.0623	0.0017	0.9674
H_MINACTIV	1	-0.0664	0.0526	1.5936	0.2068
H_PCSPMISS	1	0.0136	0.1063	0.0165	0.8979
H_PAGRICUL	1	0.2244	0.0931	5.8138	0.0159
H_PARTISAN	1	0.1023	0.0718	2.0322	0.1540
H_PCADRE	1	0.0246	0.0899	0.0750	0.7842
H_PPROFINT	1	0.0924	0.0704	1.7231	0.1893
H_PEMPLOYE	1	0.0339	0.0624	0.2957	0.5866
H_AFRNORD	1	-0.0567	0.2017	0.0791	0.7786
H_EUROP	1	0.1085	0.1661	0.4269	0.5135
H_AUTRENAT	1	0.0837	0.1951	0.1839	0.6681
H_PAFRNORD	1	-0.00706	0.1174	0.0036	0.9521
H_PEUROP	1	0.1157	0.1025	1.2733	0.2592
H_PAUTRNAT	1	-0.1482	0.0681	4.7431	0.0294
H_NDIP2	1	0.0360	0.0721	0.2496	0.6174
H_NDIP3	1	0.0432	0.0623	0.4807	0.4881
H_NDIP4	1	0.1063	0.0880	1.4589	0.2271
H_NDIP5	1	-0.1846	0.1099	2.8211	0.0930
H_NDIP6	1	-0.0549	0.0961	0.3268	0.5676
H_NDIP7	1	-0.1936	0.0983	3.8824	0.0488
H_CHAMPA97	1	-0.2672	0.1153	5.3710	0.0205
H_PICARD97	1	-0.1980	0.1167	2.8774	0.0898
H_HAUTN097	1	-0.1628	0.1068	2.3226	0.1275
H_CENTRE97	1	-0.0725	0.1096	0.4378	0.5082
H_BASSN097	1	-0.2480	0.1324	3.5110	0.0610
H_BOURG097	1	-0.2045	0.1145	3.1930	0.0740
H_NORDP97	1	-0.6270	0.0927	45.7330	<.0001
H_LORRAI97	1	-0.2985	0.1197	6.2144	0.0127
H_ALSA97	1	-0.3131	0.1078	8.4422	0.0037
H_FRCOM97	1	-0.0378	0.1199	0.0991	0.7529

H_PAYSLO97	1	-0.1121	0.0993	1.2733	0.2592
H_BRETA97	1	-0.3328	0.1011	10.8330	0.0010
H_POITOU97	1	-0.0622	0.1131	0.3026	0.5822
H_AQUITA97	1	-0.4379	0.1071	16.7092	<.0001
H_MIDIPY97	1	-0.2844	0.1331	4.5682	0.0326
H_LIMOUS97	1	-0.3096	0.1219	6.4570	0.0111
H_RHONEA97	1	-0.2625	0.0872	9.0686	0.0026
H_AUVER97	1	-0.1557	0.1170	1.7717	0.1832
H_LANGUE97	1	-0.6318	0.1229	26.4220	<.0001
H_PROVEN97	1	-0.5341	0.0963	30.7795	<.0001

Association of Predicted Probabilities and Observed Responses

Percent Concordant	76.2	Somers' D	0.527
Percent Discordant	23.5	Gamma	0.528
Percent Tied	0.2	Tau-a	0.228
Pairs	6380524	c	0.763

Les principaux déterminants de la participation des femmes au marché du travail sont le nombre d'enfants (effet négatif) et le niveau d'études (effet positif).

CHAPITRE 6

Les variables polytomiques

6.1 Cas général

Les variables polytomiques correspondent au cas où l'on observe plusieurs modalités en général, qu'elles soient ordonnées ou non. On suppose qu'une variable y_i peut prendre J modalités $y_i \in \{1, 2, \dots, J\}$. La probabilité que la variable y_i soit égale à la modalité j est notée :

$$p_{ji} = \Pr [y_i = j],$$

où les probabilités vérifient, pour chaque individu :

$$\sum_{j=1}^J p_{ji} = 1, \quad \forall i.$$

On définit également les J variables indicatrices suivantes pour chaque individu :

$$d_{ji} = \begin{cases} 1 & \text{si } y_i = j \\ 0 & \text{sinon} \end{cases} \quad j = 1, \dots, J.$$

Ces variables vérifient :

$$\sum_{j=1}^J d_{ji} = 1, \quad \forall i.$$

La log-vraisemblance d'un échantillon (y_1, \dots, y_N) s'écrit donc simplement :

$$\ell = \sum_{i=1}^N \sum_{j=1}^J d_{ji} \ln p_{ji}.$$

Le score est donc égal à :

$$\frac{\partial \ell}{\partial \theta} = \sum_{i=1}^N \sum_{j=1}^J d_{ji} \frac{\partial p_{ji}}{\partial \theta} \frac{1}{p_{ij}},$$

il est d'espérance nulle puisque $E(d_{ji}) = p_{ji}$, ce qui implique :

$$\forall i \quad E \left[\sum_{j=1}^J d_{ji} \frac{\partial p_{ji}}{\partial \theta} \frac{1}{p_{ij}} \right] = \sum_{j=1}^J \frac{\partial p_{ji}}{\partial \theta} = \frac{\partial}{\partial \theta} \underbrace{\sum_{j=1}^J p_{ji}}_1 = 0.$$

On peut appliquer l'algorithme de Berndt-Hall-Hall-Hausman à partir de la matrice :

$$W_{B3H}^{-1} = - \sum_{i=1}^N \sum_{j=1}^J \frac{d_{ji}}{p_{ji}^2} \frac{\partial p_{ji}}{\partial \theta} \frac{\partial p_{ji}}{\partial \theta'},$$

Le hessien est donné par :

$$\frac{\partial^2 \ell}{\partial \theta \partial \theta'} = \sum_{i=1}^N \sum_{j=1}^J \frac{d_{ji}}{p_{ji}} \left[\frac{\partial^2 p_{ji}}{\partial \theta \partial \theta'} - \frac{1}{p_{ji}} \frac{\partial p_{ji}}{\partial \theta} \frac{\partial p_{ji}}{\partial \theta'} \right],$$

ce qui permet d'employer l'algorithme de Newton-Raphson :

$$W_{NR}^{-1} = \sum_{i=1}^N \sum_{j=1}^J \frac{d_{ji}}{p_{ji}^2} \left[p_{ji} \frac{\partial^2 p_{ji}}{\partial \theta \partial \theta'} - \frac{\partial p_{ji}}{\partial \theta} \frac{\partial p_{ji}}{\partial \theta'} \right],$$

l'algorithme du score peut donc s'obtenir à partir de la matrice suivante :

$$\begin{aligned} E_d \left[\frac{\partial^2 \ell}{\partial \theta \partial \theta'} \right] &= \sum_{i=1}^N \sum_{j=1}^J \left\{ \frac{\partial^2 p_{ji}}{\partial \theta \partial \theta'} - \frac{1}{p_{ji}} \frac{\partial p_{ji}}{\partial \theta} \frac{\partial p_{ji}}{\partial \theta'} \right\} \\ &= - \sum_{i=1}^N \sum_{j=1}^J \frac{1}{p_{ji}} \frac{\partial p_{ji}}{\partial \theta} \frac{\partial p_{ji}}{\partial \theta'}, \end{aligned}$$

car

$$\forall i \quad \sum_{j=1}^J \frac{\partial p_{ji}}{\partial \theta} = 0 \Rightarrow \sum_{j=1}^J \frac{\partial^2 p_{ji}}{\partial \theta \partial \theta'} = 0.$$

En conséquence :

$$W_{SC}^{-1} = - \sum_{i=1}^N \sum_{j=1}^J \frac{1}{p_{ji}} \frac{\partial p_{ji}}{\partial \theta} \frac{\partial p_{ji}}{\partial \theta'}.$$

Les dérivées premières des probabilités suffisent donc pour obtenir un algorithme croissant.

6.2 Les variables ordonnées

6.2.1 Cas général

La variable expliquée que l'on observe est qualitative par classe. On suppose qu'il y a J classes numérotées de $j = 1$ à $j = J$. Les données observables sont définies par :

$$y_i = \begin{cases} 1 & \text{si } a_0 < y_i^* \leq a_1 \\ \vdots & \\ j & \text{si } a_{j-1} < y_i^* \leq a_j \\ \vdots & \\ J & \text{si } a_{J-1} < y_i^* \leq a_J \end{cases}$$

Remarque 6.1 *Dans le cas d'une variable latente réelle, on peut poser la convention $a_0 = \{-\infty\}$ et $a_J = \{+\infty\}$.*

Remarque 6.2 *Le modèle dichotomique comprend $J = 2$ classes et s'obtient comme le cas particulier où $a_0 = \{-\infty\}$, $a_1 = 0$ et $a_2 = \{+\infty\}$.*

Les variables polytomiques suivent, par définition, une distribution multinomiale dont les paramètres sont donnés par :

$$\{p_{ji}\}_{j=1,\dots,J} = \{\Pr[a_{j-1} < y_i^* \leq a_j]\}_{j=1,\dots,J}.$$

La probabilité associée à cette distribution est simplement :

$$\prod_{j=1}^J p_{ji}^{d_{ji}}, \quad \text{où } d_{ji} = \begin{cases} 1 & \text{si } a_{j-1} < y_i^* \leq a_j \\ 0 & \text{sinon} \end{cases}$$

Ces probabilités dépendent de la distribution suivie par la variable latente y_i^* . En supposant que cette perturbation suit une loi de fonction de répartition F qui dépend d'un paramètre θ , on obtient les probabilités suivantes :

$$\begin{aligned} p_{ji} &= \Pr[a_{j-1} < y_i^* \leq a_j] \\ &= \Pr[y_i^* \leq a_j] - \Pr[y_i^* \leq a_{j-1}] \\ &= F(a_j) - F(a_{j-1}). \end{aligned}$$

Remarque 6.3 *Dans le cas particulier où y^* est réelle, on a $F(a_0) = F(-\infty) = 0$ et $F(a_J) = F(+\infty) = 1$. Les probabilités s'écrivent donc $p_{1i} = F(a_1)$, $p_{ji} = F(a_j) - F(a_{j-1})$ si $2 \leq j \leq J-1$ et $p_{Ji} = 1 - F(a_{J-1})$.*

6.2.2 Le modèle Probit ordonné

Dans ce modèle on suppose que la variable latente y_i^* est générée par un modèle linéaire standard :

$$y_i^* = X_i b + u_i, \quad u_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2).$$

La variable y_i appartient à la classe j dès lors que :

$$\begin{aligned} a_{j-1} < y_i^* \leq a_j &\Leftrightarrow a_{j-1} < X_i b + u_i \leq a_j \\ &\Leftrightarrow \frac{a_{j-1} - X_i b}{\sigma} < \frac{u_i}{\sigma} \leq \frac{a_j - X_i b}{\sigma} \end{aligned}$$

où u_i/σ suit une loi normale centrée et réduite de fonction de répartition $\Phi(z)$. On en déduit que :

$$\begin{aligned} p_{ji} &= \Pr(y_{ji} = j) \\ &= \Phi\left(\frac{a_j - X_i b}{\sigma}\right) - \Phi\left(\frac{a_{j-1} - X_i b}{\sigma}\right) \\ &= \Phi\left(\frac{a_j - X_i b}{\sigma}\right) - \Phi\left(\frac{a_{j-1} - X_i b}{\sigma}\right). \end{aligned}$$

Pour simplifier la présentation, on effectue le changement de paramètre suivant :

$$\beta = \frac{b}{\sigma} \text{ et } h = \frac{1}{\sigma},$$

ce qui donne :

$$p_{ji} = \Phi(a_j h - X_i \beta) - \Phi(a_{j-1} h - X_i \beta).$$

Arrivé à ce stade, il faut distinguer le cas où les seuils a_j sont connus ou inconnus. Si les seuils sont connus, les deux paramètres h et β sont identifiables et l'on peut procéder à l'estimation par le maximum de vraisemblance.

Si les seuils sont inconnus, on ne peut pas estimer h et des termes en $a_j h$ apparaissent dans l'expression des probabilités. Ici, on peut remarquer que les quantités inconnues $a_j h$ dépendent uniquement de la classe j . Ceci revient à dire que chaque classe admet un terme constant différent et croissant avec l'ordre de la classe puisque $a_j > a_{j-1} \Rightarrow a_j h > a_{j-1} h$.

Remarque 6.4 Dans les modèles avec terme constant (le cas usuel) la constante de la première classe est prise comme référence, de sorte que a un terme constant global et les différentes estimations de $a_j h$ pour $j = 2, \dots, J - 1$ (car $a_J = \{+\infty\}$).

6.2.2.1 Estimation avec seuils connus

Pour estimer le modèle, on a juste besoin des dérivées des probabilités p_{ij} par rapport aux paramètres $\theta' = (\beta', h)$. Notons que l'on ne peut estimer β et h séparément que lorsque les seuils sont connus. On note $\phi(u)$ la densité de la loi Normale centrée et réduite, on utilise la propriété $\phi'(u) = -u\phi(u)$ et l'on pose la notation :

$$v_{ji} = a_j h - X_i \beta.$$

La forme générale de la probabilité est donc :

$$p_{ji} = \Phi(v_{ji}) - \Phi(v_{j-1i}),$$

et ses dérivées :

$$\begin{aligned} \frac{\partial p_{ji}}{\partial \beta} &= X'_i [\phi(v_{j-1i}) - \phi(v_{ji})], \\ \frac{\partial p_{ji}}{\partial h} &= a_j \phi(v_{ji}) - a_{j-1} \phi(v_{j-1i}), \end{aligned}$$

Pour calculer les probabilités correspondant aux premières et dernières modalités, il suffit de remarquer que $\phi(-\infty) = \phi(+\infty) = 0$. En conséquence, il suffit d'annuler les deux termes suivants dans les expressions ci-dessus :

$$\phi(v_{0i}) = 0 \quad \text{et} \quad \phi(v_{Ji}) = 0.$$

6.2.2.2 Estimation avec seuils inconnus

On ne peut plus estimer β et h séparément. En posant $\alpha_j = a_j h$, on doit estimer le modèle correspondant aux probabilités :

$$p_{ji} = \Phi(\alpha_j - X_i \beta) - \Phi(\alpha_{j-1} - X_i \beta).$$

Le paramètre à estimer est donc :

$$\theta' = (\beta', \alpha_2, \dots, \alpha_{J-1}).$$

et les dérivées de la probabilité sont données par :

$$\begin{aligned} \frac{\partial p_{ji}}{\partial \beta} &= X'_i [\phi(v_{j-1i}) - \phi(v_{ji})], \\ \frac{\partial p_{ji}}{\partial \alpha_j} &= \phi(v_{ji}), \end{aligned}$$

avec

$$\phi(v_{0i}) = 0 \quad \text{et} \quad \phi(v_{Ji}) = 0.$$

6.3 Les variables non ordonnées

6.3.1 Cas général

Dans ce cas, il n'est plus possible de mettre un ordre sur les modalités de la variable observable. Il peut s'agir, par exemple, d'un pays, d'un mode de transport, d'une couleur etc.. Dans ce cas on remplace l'ordre objectif (i.e. unanime) propre aux variables ordonnées par un ordre subjectif (i.e. propre à chacun) que l'on dérive d'une représentation en termes de comparaisons d'utilités. On considère qu'un individu se trouve confronté à un choix parmi J possibilités et que chaque possibilité $j \in \{1, \dots, J\}$ procure une utilité U_j définie par :

$$U_j = V_j(X_j) + \varepsilon_j, \quad j = 1, \dots, J$$

où $V_j(X_j)$ est la partie déterministe de l'utilité, c'est-à-dire la partie que l'on peut expliquer, X_j la liste des variables explicatives intervenant dans cette utilité et ε_j la partie aléatoire de l'utilité, indépendante des variables explicatives X_j . Un individu pris au hasard choisit la modalité k si :

$$U_k > U_j, \quad \forall j \neq k,$$

on observe donc le choix k par l'individu i avec la probabilité :

$$\begin{aligned} p_{ki} &= \Pr \left[\bigcap_{j \neq k} (U_{ki} > U_{ji}) \right] \\ &= \Pr \left[\bigcap_{j \neq k} (V_k(X_{ki}) + \varepsilon_{ki} > V_j(X_{ji}) + \varepsilon_{ji}) \right] \\ &= \Pr \left[\bigcap_{j \neq k} (\varepsilon_{ji} < V_k(X_{ki}) - V_j(X_{ji}) + \varepsilon_{ki}) \right]. \end{aligned}$$

La forme particulière de la probabilité p_{ji} dépend de la distribution jointe des $(\varepsilon_{1i}, \dots, \varepsilon_{Ji})$ et de la forme retenue pour la partie déterministe des fonctions d'utilité.

6.3.2 Le modèle logistique multinomial

Pour obtenir ce modèle, également appelé "multinomial Logit", on fait l'hypothèse que les $(\varepsilon_{1i}, \dots, \varepsilon_{Ji})$ sont indépendamment et identiquement distribués selon une loi de Gompertz de paramètres $(0, 1)$ de fonction de

répartition :¹

$$F(\varepsilon) = \exp(-\exp(-\varepsilon)).$$

La densité cette distribution est donc égale à :

$$f(\varepsilon) = \exp(-\varepsilon) \exp(-\exp(-\varepsilon)).$$

En notant $V_{ji} = V_j(X_{ji}), \forall j$, on a :

$$\begin{aligned} p_{ki} &= \mathbb{E}(d_{ki} = 1) \\ &= \mathbb{E} \left[\prod_{j \neq k} \Pr[\varepsilon_{ji} < V_{ki} - V_{ji} + \varepsilon_{ki}] \right], \\ &= \int_{-\infty}^{+\infty} \prod_{j \neq k} \Pr[\varepsilon_{ji} < V_{ki} - V_{ji} + \varepsilon_{ki}] f(\varepsilon_{ki}) d\varepsilon_{ki}. \end{aligned}$$

En développant, on obtient :

$$\begin{aligned} \prod_{j \neq k} \Pr[\varepsilon_{ji} < V_{ki} - V_{ji} + \varepsilon_{ki}] &= \prod_{j \neq k} \exp[-\exp(-(V_{ki} - V_{ji} + \varepsilon_{ki}))] \\ &= \prod_{j \neq k} \exp[-\exp(-\varepsilon_{ki}) \exp(V_{ji} - V_{ki})], \end{aligned}$$

ce qui implique :

$$p_{ki} = \int_{-\infty}^{+\infty} \prod_{j \neq k} \exp[-\exp(-\varepsilon_{ki}) \exp(V_{ji} - V_{ki})] \exp(-\varepsilon_{ki}) \exp(-\exp(-\varepsilon_{ki})) d\varepsilon_{ki},$$

on effectue donc le changement de variable :

$$z = \exp(-\varepsilon_{ki}),$$

ce qui implique :

$$d\varepsilon_{ki} = -\frac{dz}{z}, \quad \lim_{\varepsilon_{ki} \rightarrow -\infty} z = +\infty \quad \text{et} \quad \lim_{\varepsilon_{ki} \rightarrow +\infty} z = 0,$$

¹Une variable Z suit une loi de Gompertz de paramètres (μ, ϕ) si sa fonction de répartition s'écrit :

$$F(z) = 1 - \exp \left\{ -\exp \left(\frac{z - \mu}{\phi} \right) \right\}.$$

On obtient :

$$\mathbb{E}(Z) = \mu + 0.5772 \times \beta \quad \text{et} \quad \mathbb{V}(Z) = \frac{\phi^2 \pi^2}{6}.$$

d'où²

$$\begin{aligned}
 p_{ki} &= \int_0^{+\infty} \prod_{j \neq k} \exp(-z \exp(V_{ji} - V_{ki})) \exp(-z) dz \\
 &= \int_0^{+\infty} \exp\left(-z \sum_{j \neq k} \exp(V_{ji} - V_{ki})\right) \exp(-z) dz \\
 &= \int_0^{+\infty} \exp\left\{-z \left(1 + \sum_{j \neq k} \exp(V_{ji} - V_{ki})\right)\right\} dz.
 \end{aligned}$$

La fonction à intégrer est de la forme $\exp(-az)$ et admet pour primitive $-1/a \exp(-az)$, on obtient finalement :

$$\begin{aligned}
 p_{ki} &= \left[-\frac{\exp\left\{-z \left(1 + \sum_{j \neq k} \exp(V_{ji} - V_{ki})\right)\right\}}{1 + \sum_{j \neq k} \exp(V_{ji} - V_{ki})} \right]_0^{+\infty} \\
 &= \frac{1}{1 + \sum_{j \neq k} \exp(V_{ji} - V_{ki})}.
 \end{aligned}$$

Il est possible de réécrire cette probabilité en multipliant son numérateur et son dénominateur par $\exp(V_{ki})$, ce qui donne :

$$p_{ki} = \frac{\exp(V_{ki})}{\exp(V_{ki}) + \sum_{j \neq k} \exp(V_{ji})} = \frac{\exp(V_{ki})}{\sum_{j=1}^J \exp(V_{ji})}.$$

Pour obtenir le modèle logistique multinomial, on fait l'hypothèse supplémentaire que les utilités sont linéaires par rapport aux paramètres :

$$V_{ji}(X_j) = X_{ji}b_j, \quad j = 1, \dots, J,$$

ce qui donne finalement :

$$p_{ki} = \frac{\exp(X_{ki}b_k)}{\sum_{j=1}^J \exp(X_{ji}b_j)}, \quad k = 1, \dots, J.$$

Comme tous les paramètres du modèle ne sont pas identifiables, on doit imposer une contrainte sur les paramètres b_j . Dans le cas où les variables explicatives sont identiques pour toutes les modalités, on prend une modalité de référence, que l'on note $j = 1$; il s'agit généralement de la modalité la plus répandue. En effet dans ce cas particulier, on a :

$$p_{ki} = \frac{\exp(X_i b_k)}{\sum_{j=1}^J \exp(X_i b_j)} = \frac{\exp(X_i (b_k - b_1))}{1 + \sum_{j=2}^J \exp(X_i (b_j - b_1))},$$

²On utilise le fait qu'invertir les bornes de l'intégrale change son signe.

et tous les effets sont mesurés par rapport à l'utilité que procure le choix le plus courant. On voit clairement que seules les transformations de $b_j - b_1$ peuvent être estimées. On pose $\beta_j = b_j - b_1$, ce qui donne :

$$p_{ki} = \frac{\exp(X_i\beta_k)}{1 + \sum_{j=2}^J \exp(X_i\beta_j)}, \quad k = 1, \dots, J.$$

On remarque ici que pour $k = 1$:

$$p_{1i} = \frac{1}{1 + \sum_{j=2}^J \exp(X_i\beta_j)},$$

parce que $\beta_1 = b_1 - b_1 = 0$. La nullité du coefficient de la première modalité est donc la contrainte identifiante du modèle. On ne calcule donc pas de dérivée par rapport à β_1 . Les dérivées de la probabilité sont données simplement par :

$$k = 2, \dots, J, \quad \frac{\partial p_{ki}}{\partial \beta_k} = X'_i p_{ki} (1 - p_{ki}), \quad \frac{\partial p_{ki}}{\partial \beta_j} \Big|_{j \neq k} = -X'_i p_{ki} p_{ji}.$$

CHAPITRE 7

Le pseudo maximum de vraisemblance

Dans le cas général, l'estimateur du maximum de vraisemblance n'est convergent et asymptotiquement efficace que si l'hypothèse que l'on fait sur la loi conditionnelle de la variable expliquée y est juste. Sinon, il peut ne pas être convergent. Il existe une famille de distributions pour lesquelles une erreur de spécification de ce type ne remet pas en cause la convergence de l'estimateur du maximum de vraisemblance. Par contre il faut évaluer différemment sa matrice de covariance, c'est à dire changer les statistiques de test. Nous supposons donc dans cette section que nous ne connaissons *que* l'espérance conditionnelle de la variable expliquée $E(y|X, \theta)$.

7.1 Le pseudo maximum de vraisemblance à l'ordre 1

7.1.1 La famille exponentielle linéaire à l'ordre 1

DÉFINITION 7.1 *La famille exponentielle linéaire à l'ordre 1 désigne une famille de distributions dont la densité admet la forme suivante :*

$$f(y, m) = \exp \{A(m) + B(y) + C(m) \times y\} \quad \text{où} \quad E(y) = m.$$

Cette forme est vérifiée par de nombreuses lois usuelles, dont voici quelques exemples :

Exemple 7.1 *Loi normale. $y \rightsquigarrow N(m, \omega)$. Son espérance mathématique*

est égale à $m \in \mathbb{R}$ et la densité s'écrit :

$$\begin{aligned} \ln f(y, m) &= \ln \left[\frac{1}{\sqrt{2\pi\omega}} \exp \left(-\frac{1}{2\omega} (y - m)^2 \right) \right] \\ &= \underbrace{-\frac{m^2}{2\omega}}_{A(m)} - \underbrace{\frac{1}{2} \ln(2\pi\omega) - \frac{y^2}{2\omega}}_{B(y)} + y \times \underbrace{\frac{m}{\omega}}_{C(m)}. \end{aligned}$$

Exemple 7.2 Loi de Poisson. $y \rightsquigarrow P(m)$. Son espérance mathématique est égale à $m \in \mathbb{R}^+$ et les probabilités s'écrivent :

$$\begin{aligned} \ln f(y, m) &= \ln \left[\frac{\exp(-m) m^y}{y!} \right] \\ &= \underbrace{-m}_{A(m)} - \underbrace{\ln y!}_{B(y)} + y \times \underbrace{\ln m}_{C(m)}. \end{aligned}$$

Exemple 7.3 Loi de Bernoulli. $y \rightsquigarrow B(m)$. Son espérance mathématique est égale à $m \in]0, 1[$ et les probabilités s'écrivent :

$$\begin{aligned} \ln f(y, m) &= \ln \left[m^y (1 - m)^{1-y} \right] \\ &= \underbrace{\ln(1 - m)}_{A(m)} + \underbrace{0}_{B(y)} + y \times \underbrace{\ln \frac{m}{1 - m}}_{C(m)}. \end{aligned}$$

Exemple 7.4 Loi Gamma. $y \rightsquigarrow \Gamma(a, b)$. Ce cas est un peu plus compliqué que les précédents car l'espérance mathématique n'est pas un paramètre utilisé habituellement avec cette loi. Il faut donc la réécrire en fonction de $m = ab$. La forme usuelle de la probabilité est :

$$f(y) = \frac{y^{a-1} \exp(-y/b)}{b^a \Gamma(a)}, \quad a, b > 0, \quad y > 0$$

avec :¹

$$m = ab \Leftrightarrow b = m/a$$

en remplaçant dans la densité, on obtient :

$$\begin{aligned} \ln f(y, m) &= \ln \left[\frac{y^{a-1} \exp(-ya/m)}{(m/a)^a \Gamma(a)} \right] \\ &= \underbrace{-a \ln(m/a) - \ln \Gamma(a)}_{A(m)} + \underbrace{(a-1) \ln y}_{B(y)} + y \times \underbrace{\left(-\frac{a}{m} \right)}_{C(m)}. \end{aligned}$$

¹ Ici on peut prendre soit $b = m/a$ soit $a = m/b$, et l'on prend la première possibilité parce qu'elle mène aux calculs les plus simples.

Exemple 7.5 *Loi Binomiale négative.* $y \rightsquigarrow BN(r, p)$. Dans ce cas également, l'espérance mathématique n'est pas un paramètre utilisé habituellement avec cette loi. Il faut donc la réécrire en fonction de $m = r(1-p)/p$. La forme usuelle de la probabilité est :

$$f(y) = \frac{\Gamma(y+r)}{\Gamma(y+1)\Gamma(r)} p^r (1-p)^y, \quad 0 < p < 1 \text{ et } r > 0,$$

avec :

$$m = r \frac{(1-p)}{p} \Leftrightarrow p = \frac{r}{r+m},$$

en remplaçant dans l'expression de la densité on obtient :

$$\begin{aligned} \ln f(y, m) &= \ln \left[\frac{\Gamma(y+r)}{\Gamma(y+1)\Gamma(r)} \left(\frac{r}{r+m} \right)^r \left(\frac{m}{r+m} \right)^y \right] \\ &= \underbrace{r \ln \left(\frac{r}{r+m} \right)}_{A(m)} + \underbrace{\ln \frac{\Gamma(y+r)}{\Gamma(y+1)\Gamma(r)}}_{B(y)} + y \times \underbrace{\ln \left(\frac{m}{r+m} \right)}_{C(m)}. \end{aligned}$$

PROPRIÉTÉ 7.1 *Les lois de la famille exponentielle linéaire vérifient les deux propriétés suivantes :*

1. $\frac{\partial A}{\partial m} + m \frac{\partial C}{\partial m} = 0.$

2. $V(y) = \left(\frac{\partial C}{\partial m} \right)^{-1}$

PREUVE :

Pour démontrer la première propriété, on remarque que :

$$\int f(y, m) dy = 1,$$

en dérivant cette relation par rapport à m , on obtient :

$$\begin{aligned} \int \left(\frac{\partial A}{\partial m} + \frac{\partial C}{\partial m} \times y \right) f(y, m) dy &= 0 \\ \Leftrightarrow \frac{\partial A}{\partial m} \underbrace{\int f(y, m) dy}_1 + \frac{\partial C}{\partial m} \underbrace{\int y f(y, m) dy}_{E(y)} &= 0 \\ \Leftrightarrow \frac{\partial A}{\partial m} + m \frac{\partial C}{\partial m} &= 0. \end{aligned}$$

Pour démontrer la seconde propriété, on remarque que :

$$E(y) = \int y f(y, m) dy = m,$$

en dérivant cette relation par rapport à m , on obtient :

$$\begin{aligned} & \int \left(\frac{\partial A}{\partial m} + \frac{\partial C}{\partial m} \times y \right) y f(y, m) dy = 1 \\ \Leftrightarrow & \frac{\partial A}{\partial m} \underbrace{\int y f(y, m) dy}_{E(y)} + \frac{\partial C}{\partial m} \underbrace{\int y^2 f(y, m) dy}_{E(y^2)} = 1 \\ \Leftrightarrow & \frac{\partial A}{\partial m} m + \frac{\partial C}{\partial m} (V(y) + m^2) = 1 \\ \Leftrightarrow & \left(-m \frac{\partial C}{\partial m} \right) m + \frac{\partial C}{\partial m} (V(y) + m^2) = 1 \\ \Leftrightarrow & \frac{\partial C}{\partial m} V(y) = 1 \Leftrightarrow V(y) = \left(\frac{\partial C}{\partial m} \right)^{-1}. \end{aligned}$$

□

Vérifions-le sur nos exemples :

Exemple 7.6 Loi normale $N(m, \omega)$. On a $C(m) = m/\omega$ donc $\partial C/\partial m = 1/\omega$ et $V(y) = \omega$.

Exemple 7.7 Loi de Poisson $P(m)$. On a $C(m) = \ln m$ donc $\partial C/\partial m = 1/m$ et $V(y) = m$.

Exemple 7.8 Loi de Bernoulli $B(m)$. On a $C(m) = \ln(m/(1-m))$ donc $\partial C/\partial m = 1/(m(1-m))$ et $V(y) = m(1-m)$.

Exemple 7.9 Loi Gamma $\Gamma(a, b)$ avec $m = ab$. On a $C(m) = -a/m$ donc $\partial C/\partial m = a/m^2$ et $V(y) = m^2/a = ab^2$.

Exemple 7.10 Loi Binomiale Négative $BN(r, p)$ avec $m = r(1-p)/p$. On a $C(m) = \ln(m/(m+r))$ donc $\partial C/\partial m = r/(m(m+r))$ et $V(y) = m(m+r)/r = r(1-p)/p^2$.

7.1.2 Estimation

On note θ le paramètre à estimer. Ce paramètre intervient ici dans l'espérance mathématique de la distribution de y . Ainsi, dans le cas d'un modèle linéaire avec variables explicatives X_i , on aurait $m = X_i\theta$. L'espérance peut donc être différente pour chaque observation dans un modèle avec variables explicatives. Plus généralement, on considérera une espérance conditionnelle sous la forme générale :

$$m_i = m(X_i, \theta).$$

L'estimateur du pseudo maximum de vraisemblance est obtenu en maximisant la pseudo vraisemblance suivante :

$$\begin{aligned} \ell(y_i | X_i; \theta) &= \sum_{i=1}^N \ln f(y_i | m(X_i, \theta)) \\ &= \sum_{i=1}^N A[m(X_i, \theta)] + B(y_i) + C[m(X_i, \theta)] \times y_i. \end{aligned}$$

On remarque que pour la maximisation, on peut négliger les termes qui ne dépendent pas de θ . L'estimateur du pseudo maximum de vraisemblance $\tilde{\theta}_N$ se définit donc comme :

$$\begin{aligned} \tilde{\theta}_N &= \arg \max_{\theta} \sum_{i=1}^N A[m(X_i, \theta)] + B(y_i) + C[m(X_i, \theta)] \times y_i \\ &= \arg \max_{\theta} \sum_{i=1}^N A[m(X_i, \theta)] + C[m(X_i, \theta)] \times y_i. \end{aligned}$$

La condition du premier ordre est donnée par :

$$\begin{aligned} \frac{\partial \ell}{\partial \theta} (y_i | X_i; \tilde{\theta}_N) &= 0 \\ \Leftrightarrow \sum_{i=1}^N \frac{\partial m}{\partial \theta} \left[\frac{\partial A}{\partial m} + \frac{\partial C}{\partial m} \times y_i \right] &= 0 \\ \Leftrightarrow \sum_{i=1}^N \frac{\partial m}{\partial \theta} \frac{\partial C}{\partial m} \left[y_i - m(X_i, \tilde{\theta}_N) \right] &= 0 \\ \Leftrightarrow \sum_{i=1}^N \frac{\partial m}{\partial \theta} \Sigma^{-1} \left[y_i - m(X_i, \tilde{\theta}_N) \right] &= 0, \end{aligned}$$

en utilisant le fait que $\partial A / \partial m + m \times \partial C / \partial m = 0$ et que $[\partial C / \partial m]^{-1} = V(y) = \Sigma$. Remarquons ici que Σ est la matrice de covariance de la pseudo loi,

qui peut être différente de la matrice de covariance de la vraie loi, qui est inconnue dans le cas général (i.e., nous n'avons pas fait d'hypothèse sur la variance conditionnelle). La vraie matrice de covariance est notée Ω . On obtient l'équivalent des matrices $I_N(\theta)$ et $J_N(\theta)$ de la manière suivante :²

$$\begin{aligned} I_N(\theta) &= \sum_{i=1}^N E_0 \left[\frac{\partial \ln f(y_i|m(X_i;\theta))}{\partial \theta} \frac{\partial \ln f(y_i|m(X_i;\theta))}{\partial \theta'} \right] \\ &= \sum_{i=1}^N E_0 \left[\frac{\partial m}{\partial \theta} \Sigma^{-1} (y_i - m)^2 \Sigma^{-1} \frac{\partial m}{\partial \theta'} \right] \\ &= \sum_{i=1}^N \frac{\partial m}{\partial \theta} \Sigma^{-1} E_0 \left[(y_i - m)^2 \right] \Sigma^{-1} \frac{\partial m}{\partial \theta'} \\ &= \sum_{i=1}^N \frac{\partial m}{\partial \theta} \Sigma^{-1} \Omega \Sigma^{-1} \frac{\partial m}{\partial \theta'}, \end{aligned}$$

où $E_0[\cdot]$ représente l'espérance mathématique par rapport à la vraie loi de y . Pour obtenir $J_N(\theta)$ il faut se rappeler que Σ dépend de m :

$$\begin{aligned} J_N(\theta) &= \sum_{i=1}^N E_0 \left[-\frac{\partial^2 \ln f(y_i|m(X_i;\theta))}{\partial \theta \partial \theta'} \right] \\ &= -\sum_{i=1}^N E_0 \left[\frac{\partial^2 m}{\partial \theta \partial \theta'} \Sigma^{-1} (y_i - m) \right] \\ &\quad - \sum_{i=1}^N E_0 \left[\frac{\partial m}{\partial \theta} \left(\frac{\partial (\Sigma^{-1})}{\partial \theta'} (y_i - m) - \Sigma^{-1} \frac{\partial m}{\partial \theta'} \right) \right] \\ &= \sum_{i=1}^N \frac{\partial m}{\partial \theta} \Sigma^{-1} \frac{\partial m}{\partial \theta'}. \end{aligned}$$

On remarque donc que pour $\Sigma \neq \Omega$, on a :

$$I_N(\theta) \neq J_N(\theta).$$

En effectuant un développement limité de la condition du premier

²On prend les espérances mathématiques par rapport à la vraie loi car c'est vers ces quantités que convergeront nos statistiques.

ordre au voisinage de $\tilde{\theta}_N = \theta$, on obtient :

$$\begin{aligned} \frac{\partial \ell}{\partial \theta} (y_i | X_i; \tilde{\theta}_N) &\stackrel{pp}{\cong} \frac{\partial \ell}{\partial \theta} (y_i | X_i; \theta) + \frac{\partial^2 \ell}{\partial \theta \partial \theta'} (y_i | X_i; \theta) (\tilde{\theta}_N - \theta) \\ \Leftrightarrow \sqrt{N} (\tilde{\theta}_N - \theta) &\stackrel{pp}{\cong} \left[-\frac{1}{N} \frac{\partial^2 \ell}{\partial \theta \partial \theta'} (y_i | X_i; \theta) \right]^{-1} \frac{1}{\sqrt{N}} \frac{\partial \ell}{\partial \theta} (y_i | X_i; \theta), \end{aligned}$$

le premier terme converge vers $J_1(\theta)$ et le second terme est la moyenne de variables aléatoires d'espérance nulle et de variance $I_1(\theta)$. En conséquence :

$$\sqrt{N} (\tilde{\theta}_N - \theta) \stackrel{A}{\rightsquigarrow} N(0, J_1^{-1}(\theta) I_1(\theta) J_1^{-1}(\theta)),$$

dans la pratique, on utilise :

$$\tilde{\theta}_N \stackrel{A}{\rightsquigarrow} N(\theta, J_N^{-1}(\theta) I_N(\theta) J_N^{-1}(\theta)).$$

En effet, on voit que :

$$\begin{aligned} \frac{1}{N} J_1^{-1}(\theta) I_1(\theta) J_1^{-1}(\theta) &= [NJ_1(\theta)]^{-1} [NI_1(\theta)] [NJ_1(\theta)]^{-1} \\ &= J_N^{-1}(\theta) I_N(\theta) J_N^{-1}(\theta). \end{aligned}$$

Reprenons nos trois exemples. Dans tous les cas, la spécification de la moyenne est supposée juste. On a donc $E(y) = \theta$ et dans tous les cas, on trouve $\tilde{\theta}_N = 1/N \sum_{i=1}^N y_i$. Pourtant, selon que la loi est normale, de Poisson ou de Bernoulli, les variances asymptotiques sont égales à ω/N , θ/N et $\theta(1-\theta)/N$. Pour calculer les variances asymptotiques du pseudo maximum de vraisemblance, il faut recalculer les matrices $I_N(\theta)$. On suppose, dans les trois exemples suivants que la vraie variance est égale à ω . La suite montre que, lorsque l'on ne connaît pas la vraie loi avec certitude, il faut mieux utiliser l'estimateur suivant :

$$\widehat{\text{Vas}}(\tilde{\theta}_N) = \frac{\hat{\omega}}{N} \quad \text{avec} \quad \hat{\omega} = \frac{1}{N} \sum_{i=1}^N y_i^2.$$

Exemple 7.11 Pseudo loi normale $N(\theta, \omega)$. La matrice $I_N(\theta)$ est maintenant donnée par :

$$I_N(\theta) = \sum_{i=1}^N V_0[y_i - \theta] = \frac{N}{\omega},$$

donc la variance asymptotique de $\tilde{\theta}_N$ doit être estimée par :

$$J_N^{-1}(\theta) I_N(\theta) J_N^{-1}(\theta) = \left(\frac{N}{\omega}\right)^{-1} \frac{N}{\omega} \left(\frac{N}{\omega}\right)^{-1} = \frac{\omega}{N}.$$

Exemple 7.12 Pseudo loi de Poisson $P(\theta)$. La matrice $I_N(\theta)$ est maintenant donnée par :

$$I_N(\theta) = \sum_{i=1}^N V_0 \left[-1 + \frac{y_i}{\theta} \right] = \frac{N\omega}{\theta^2},$$

donc la variance asymptotique de $\tilde{\theta}_N$ doit être estimée par :

$$J_N^{-1}(\theta) I_N(\theta) J_N^{-1}(\theta) = \left(\frac{N}{\theta} \right)^{-1} \frac{N\omega}{\theta^2} \left(\frac{N}{\theta} \right)^{-1} = \frac{\omega}{N}.$$

Exemple 7.13 Pseudo loi de Bernoulli $B(\theta)$. La matrice $I_N(\theta)$ est maintenant donnée par :

$$I_N(\theta) = \sum_{i=1}^N V_0 \left[\frac{y_i}{\theta(1-\theta)} - \frac{1}{1-\theta} \right] = \frac{N\omega}{\theta^2(1-\theta)^2},$$

donc la variance asymptotique de $\tilde{\theta}_N$ doit être estimée par :

$$J_N^{-1}(\theta) I_N(\theta) J_N^{-1}(\theta) = \left(\frac{N}{\theta(1-\theta)} \right)^{-1} \frac{N\omega}{\theta^2(1-\theta)^2} \left(\frac{N}{\theta(1-\theta)} \right)^{-1} = \frac{\omega}{N}.$$

7.1.3 Matrice de covariance robuste à l'hétéroscédasticité de forme inconnue

Le modèle linéaire standard, estimé par le maximum de vraisemblance sous hypothèse de normalité, fournit un estimateur convergent de b même si la distribution de la perturbation u_i n'est pas normale $N(0, \omega)$. C'est parce que la loi normale appartient à la famille exponentielle linéaire. En effet, elle peut s'écrire sous la forme :

$$f(y, m) = \exp \{ A(m) + B(y) + y \times C(m) \}$$

où $m = Xb$ est l'espérance conditionnelle de y . Notons dès maintenant que quelle que soit la valeur du paramètre du second ordre ω , f appartient à la famille exponentielle linéaire. Nous pouvons donc nous en passer et lui donner une valeur quelconque. Ceci tient au fait que le pseudo maximum de vraisemblance à l'ordre 1 ne fait pas d'hypothèse sur la variance conditionnelle de y . Le paramètre ω est alors un paramètre de nuisance que l'on peut fixer arbitrairement; ce n'est plus nécessairement la variance des perturbations, car la vraie loi n'est pas normale dans le cas général. De plus, dans le cas du modèle linéaire, ceci n'affectera pas notre estimateur de b , puisqu'il est déterminé indépendamment de l'estimateur de ω . On peut donc poser $\omega = 1$ sans perte de généralité. On a :

$$A(m) = -\frac{1}{2} [m^2 + \ln(2\pi)], \quad B(y) = -\frac{y^2}{2}, \quad C(m) = -m.$$

La matrice de covariance de l'estimateur des moindres carrés n'est plus estimée par l'inverse de l'information de Fisher mais par :

$$\widehat{\text{Vas}}\left(\sqrt{N}\left(\widehat{b}_N - b\right)\right) = \widehat{J}_1^{-1}\left(\widehat{b}_N\right) \widehat{I}_1\left(\widehat{b}_N\right) \widehat{J}_1^{-1}\left(\widehat{b}_N\right). \quad (7.1)$$

La matrice $J_1(b)$ reste inchangée :

$$J_1(b) = \mathbb{E}_X \mathbb{E}_0 \left[-\frac{\partial^2 \ln f}{\partial b \partial b'} \right] = \mathbb{E}_X [X'X],$$

que l'on peut estimer de manière convergente par :

$$\widehat{J}_1\left(\widehat{b}_N\right) = \frac{1}{N} \sum_{i=1}^N X_i' X_i.$$

La matrice $I_1(b)$ est égale à :

$$I_1(b) = \mathbb{E}_X \mathbb{E}_0 \left[\frac{\partial \ln f}{\partial b} \frac{\partial \ln f}{\partial b'} \right] = \mathbb{E}_X \mathbb{E}_0 [X'X u^2],$$

avec

$$u = y - Xb,$$

que l'on peut estimer de manière convergente par :

$$\widehat{I}_1\left(\widehat{b}_N\right) = \frac{1}{N} \sum_{i=1}^N X_i' X_i \widehat{u}_i^2 \quad \text{avec} \quad \widehat{u}_i = y_i - X_i \widehat{b}_N.$$

On utilise donc finalement :

$$\widehat{\text{Vas}}\left(\widehat{b}_N\right) = \left(\sum_{i=1}^N X_i' X_i \right)^{-1} \sum_{i=1}^N X_i' X_i \widehat{u}_i^2 \left(\sum_{i=1}^N X_i' X_i \right)^{-1}.$$

La matrice de covariance ainsi obtenue est appelée *matrice de covariance robuste* (White, 1980). Plus précisément elle est robuste aux hypothèses de normalité et d'homoscédasticité des perturbations.³ On l'utilise donc systématiquement de nos jours.

³La variance habituelle de l'estimateur des moindres carrés ordinaires est robuste à l'hypothèse de normalité, seule l'hétéroscédasticité pose réellement problème ici.

7.2 Le pseudo maximum de vraisemblance quasi généralisé

7.2.1 La famille exponentielle quasi-généralisée

DÉFINITION 7.2 *La famille exponentielle quasi généralisée désigne une famille de distributions dont la densité admet la forme suivante :*

$$f(y, m, \eta) = \exp \{A(m, \eta) + B(y, \eta) + C(m, \eta) \times y\}$$

où $E(y) = m$ et $V(y) = [\partial C(m, \eta) / \partial m]^{-1}$.

où η est un paramètre intervenant dans la variance de la pseudo distribution. Il n'est pas forcément égal à la variance de la distribution. Les lois de Poisson et de Bernoulli n'admettent pas de paramètre spécifique intervenant dans la variance, elles n'appartiennent donc pas à la famille exponentielle quasi-généralisée. Par contre, la loi normale appartient à cette famille. Nous introduisons également la loi binomiale négative qui généralise la loi de Poisson. Le lecteur pourra vérifier par lui-même qu'elle appartient à la famille linéaire exponentielle à l'ordre 1.

Exemple 7.14 *Loi normale.* $y \rightsquigarrow N(m, \eta)$. On a $E(y) = m$ et $V(y) = \eta$. La densité s'écrit :

$$\begin{aligned} \ln f(y, m, \eta) &= \ln \left[\frac{1}{\sqrt{2\pi\eta}} \exp \left(-\frac{1}{2\eta} (y - m)^2 \right) \right] \\ &= \underbrace{-\frac{m^2}{2\eta}}_{A(m, \eta)} - \underbrace{\frac{1}{2} \ln(2\pi\eta) - \frac{y^2}{2\eta}}_{B(y, \eta)} + y \times \underbrace{\frac{m}{\eta}}_{C(m, \eta)}. \end{aligned}$$

Exemple 7.15 *Loi Gamma.* $y \rightsquigarrow \Gamma(\eta, m/\eta)$. On a $E(y) = m$ et $V(y) = m^2/\eta$. La densité s'écrit :

$$\begin{aligned} \ln f(y, m, \eta) &= \ln \left[\frac{y^{\eta-1} \exp(-y\eta/m)}{(m/\eta)^\eta \Gamma(\eta)} \right] \\ &= \underbrace{-\eta \ln(m/\eta) - \ln \Gamma(\eta)}_{A(m, \eta)} + \underbrace{(\eta - 1) \ln y}_{B(y, \eta)} + y \times \underbrace{\left(-\frac{\eta}{m} \right)}_{C(m, \eta)}. \end{aligned}$$

Exemple 7.16 *Loi binomiale négative.* $y \rightsquigarrow BN(m, \eta)$. On a $E(y) = m$ et $V(y) = m(1 + \eta m)$. On remarque que, contrairement à la loi de Poisson, la variance conditionnelle peut être différente de la moyenne conditionnelle. On retrouve la loi de Poisson en prenant la limite quand $\eta \rightarrow 0$. La densité s'applique, dans le cas de base, aux variables entières

positives et s'écrit :

$$f(y, m, \eta) = \frac{\Gamma(y + 1/\eta)}{\Gamma(y + 1)\Gamma(1/\eta)} \left(\frac{\eta m}{1 + \eta m}\right)^y \left(\frac{1}{1 + \eta m}\right)^{1/\eta},$$

ce qui implique :

$$\ln f(y, m, \eta) = \underbrace{\ln \left[\left(\frac{1}{1 + \eta m}\right)^{1/\eta} \frac{1}{\Gamma(1/\eta)} \right]}_{A(m, \eta)} + \underbrace{\ln \frac{\Gamma(y + 1/\eta)}{\Gamma(y + 1)}}_{B(y, \eta)} + \underbrace{y \ln \left(\frac{\eta m}{1 + \eta m}\right)}_{C(m, \eta)}.$$

7.2.2 Estimation

PROPRIÉTÉ 7.2 *L'estimateur du pseudo maximum de vraisemblance quasi généralisé (PMVQG) $\bar{\theta}_N$ vérifie les trois propriétés suivantes:*

1. $\sqrt{N} (\bar{\theta}_N - \theta) \overset{A}{\rightsquigarrow} N(0, J_1^{-1}(\theta))$ avec :

$$J_1(\theta) = \mathbb{E}_X \mathbb{E}_0 \left[\frac{\partial m}{\partial \theta} \Sigma^{-1} \frac{\partial m}{\partial \theta'} \right].$$

2. Il atteint la borne inférieure des matrices de covariances des estimateurs du pseudo maximum de vraisemblance à l'ordre 1.
3. Si la vraie loi de y appartient à la famille exponentielle linéaire et si les paramètres m et η sont fonctionnellement indépendants, l'estimateur du PMVQG est asymptotiquement équivalent à l'estimateur du maximum de vraisemblance sur la vraie loi.

7.2.3 Les moindres carrés pondérés

Supposons que l'on ait un modèle linéaire hétéroscédastique dont la forme de la variance est connue, donnée par :

$$V(y_i | X_i) = \omega(X_i, \eta) > 0.$$

Dans un premier temps, on applique le pseudo maximum de vraisemblance à l'ordre 1, sans tenir compte de l'hétéroscédasticité des perturbations. Cet estimateur est convergent et sa variance asymptotique est donnée par la relation (7.1). Dans un second temps, on estime le paramètre η . Pour cela, on utilise la relation :

$$\mathbb{E} \left((y_i - m_i)^2 \right) = \omega(X_i, \eta),$$

ce qui permet d'écrire la régression :

$$u_i^2 = \omega(X_i, \eta) + v_i \quad \text{avec} \quad \mathbb{E}(v_i) = 0.$$

on obtient donc un estimateur convergent de η , noté $\hat{\eta}$, en remplaçant u_i^2 par \hat{u}_i^2 . A partir de cet estimateur, on estime la variance par :

$$\hat{\omega}_i = \omega(X_i, \hat{\eta}).$$

Ensuite on maximise la pseudo vraisemblance quasi généralisée, obtenue en posant $\omega = \hat{\omega}_i$ dans la pseudo vraisemblance. On note que la variance est différente avec chaque observation et que $\hat{\omega}_i$ ne dépend que de \hat{b}_N , pas de b . Ceci donne l'estimateur :

$$\begin{aligned} \bar{b}_N &= \arg \max_b \bar{\ell}(y|X, b, \hat{\omega}) \\ &= \arg \max_b -\frac{N}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^N \ln \hat{\omega}_i - \frac{1}{2} \sum_{i=1}^N \frac{1}{\hat{\omega}_i} (y_i - X_i b)^2 \\ &= \arg \min_b \sum_{i=1}^N \left(\frac{y_i - X_i b}{\sqrt{\hat{\omega}_i}} \right)^2, \end{aligned}$$

ce qui correspond à l'estimateur des moindres carrés ordinaires avec une pondération égale à l'inverse de l'écart-type de la perturbation. La variance asymptotique de cet estimateur est donné par :

$$\widehat{\text{Vas}}(\bar{b}_N) = \left[\hat{J}_1(\bar{b}_N) \right]^{-1} = \left(\sum_{i=1}^N \frac{1}{\hat{\omega}_i} X_i' X_i \right)^{-1}.$$

CHAPITRE 8

Les variables entières

8.1 Le modèle de Poisson

8.1.1 Introduction

La loi de Poisson permet de décrire le nombre de réalisations d'un événement pendant une période de temps donnée. Soit y_i la variable que l'on souhaite expliquer (e.g., le nombre de brevets). Comme l'espérance mathématique d'une donnée de comptage est toujours strictement positive on l'écrit sous la forme :¹

$$E(y_i|X_i, b) = \exp(X_i b) > 0,$$

où X_i est le vecteur des variables explicatives et b le paramètre correspondant. On peut réécrire cette relation sous la forme :

$$\ln E(y_i|X_i, b) = X_i b,$$

de sorte que

$$\frac{\partial \ln E(y_i|X_i, b)}{\partial X_i'} = b.$$

Cette relation fait apparaître que le paramètre b sera le vecteur des élasticités de l'espérance de y par rapport à X dès lors que les variables explicatives seront prises en logarithmes. Mais, contrairement au cas du modèle linéaire, il s'agira pas ici de l'élasticité de y par rapport à X .

¹Une variable de comptage ne prend que des valeurs positives ou nulle, d'où le résultat.

8.1.1.1 Modèle homogène

Le modèle de Poisson homogène s'obtient dès lors que l'on postule que les y_i sont indépendamment et identiquement distribués selon une loi de Poisson de moyenne conditionnelle $m_i = \exp(X_i b)$. La probabilité d'observer une réalisation y_i de la variable de la variable de comptage est donc donnée par :

$$f(y_i) = \frac{\exp(-m_i) m_i^{y_i}}{y_i!}, \quad y_i \in \{0, 1, 2, \dots\}.$$

Cette hypothèse implique que la variance conditionnelle est égale à la moyenne conditionnelle :

$$V(y_i|X_i) = E(y_i|X_i) \quad \forall i.$$

On obtient ce modèle en postulant qu'il n'y a pas de perturbation dans l'expression de l'espérance conditionnelle de y_i . Les perturbations expriment généralement une forme d'hétérogénéité individuelle inobservable, de sorte qu'en leur absence on parle de modèle homogène.

8.1.1.2 Modèle hétérogène

On peut penser que le modèle précédent est insuffisant pour représenter les différences entre les individus, car celles-ci ne s'expriment que par les variables déterministes X_i . On peut penser qu'il existe également des caractéristiques individuelles inobservables, supposées sans corrélation avec les X_i , qui interviennent également dans l'espérance. La moyenne comporte alors une partie déterministe et une partie aléatoire u_i :

$$\ln E(y_i|X_i, b) = X_i b + u_i,$$

où u_i est une perturbation qui vérifie $E(\exp u_i) = 1$. La moyenne du processus de Poisson est alors elle-même aléatoire de sorte qu'il y a deux sources d'aléa dans ce modèle : d'une part, un aléa sur la moyenne et, d'autre part, un aléa lié au tirage dans une loi de Poisson de moyenne donnée. On note cette moyenne \tilde{m}_i , définie par :

$$\tilde{m}_i = \exp(X_i b + u_i) = \exp(X_i b) \exp(u_i) = m_i \exp(u_i).$$

Pour pouvoir écrire la vraisemblance de ce modèle, il faut faire une hypothèse spécifique sur la loi de $\exp(u_i)$. Nous ne suivrons pas cette approche ici car seuls quelques cas peuvent être écrits sous forme explicite; la plupart du temps, il faudrait recourir à l'intégration numérique. Nous prendrons donc une approche par le pseudo maximum de vraisemblance

à l'ordre 1, qui ne nécessite que l'expression de la moyenne conditionnelle de y_i . Ici, cette expression est tout simplement donnée par :

$$E(y_i|X_i, b) = E(m_i \exp u_i) = m_i \underbrace{E(\exp u_i)}_1 = m_i.$$

Notons ici que l'hypothèse n'implique aucune perte de généralité tant que le modèle contient un terme constant. Si on avait fait l'hypothèse que $E(\exp u_i) = k$, on aurait trouvé :

$$E(y_i|X_i, b) = km_i = k \exp(X_i b) = \exp(-\ln k + X_i b),$$

de sorte que le terme $-\ln k$ est absorbé par le terme constant du modèle.

8.1.2 Estimation

L'espérance mathématique d'un modèle hétérogène est la même que celle d'un modèle homogène. De plus, nous avons vu précédemment que la loi de Poisson appartient à la famille exponentielle linéaire. Dans ces conditions, quelle est la différence entre l'estimateur du maximum de vraisemblance et celui du pseudo maximum de vraisemblance? Réponse : la matrice de covariance asymptotique. Celle du pseudo maximum de vraisemblance est robuste.

8.1.2.1 Maximum de vraisemblance

La log-vraisemblance pour une observation est donnée par :

$$\begin{aligned} \ell_i &= \ln \left\{ \frac{\exp(-m_i) m_i^{y_i}}{y_i!} \right\} = -m_i + y_i \ln m_i - \ln y_i! \\ &= -\exp(X_i b) + y_i X_i b - \ln y_i!. \end{aligned}$$

On peut également l'écrire comme :

$$\ell(y_i, \mu_i) = -\exp(\mu_i) + y_i \mu_i - \ln y_i!,$$

avec $\mu_i = X_i b$. Les dérivées par rapport à μ_i sont égales à :

$$\frac{\partial \ell_i}{\partial \mu_i} = y_i - \exp(\mu_i) = y_i - \lambda_i, \quad \frac{\partial^2 \ell_i}{\partial \mu_i^2} = -\exp(\mu_i) = -m_i.$$

On en déduit que le score est égal à

$$\frac{\partial \ell}{\partial b} = \sum_{i=1}^N X_i' (y_i - m_i),$$

et que le hessien est égal à :

$$\frac{\partial^2 \ell}{\partial b \partial b'} = - \sum_{i=1}^N X_i' X_i m_i \ll 0.$$

La nullité du score exprime ici encore la condition d'orthogonalité entre les variables explicatives X et le résidu de la régression $\hat{u}_i = y_i - \hat{m}_i = y_i - \hat{E}(y_i | X_i, b)$. La seule différence avec le cas habituel est que l'espérance mathématique est non linéaire $\hat{m}_i = \exp(X_i \hat{b})$. Comme le hessien est défini négatif, l'estimateur du maximum de vraisemblance \hat{b} est unique et donné par la condition du premier ordre :

$$\frac{\partial \ell}{\partial b}(\hat{b}) = 0 \Leftrightarrow \sum_{i=1}^N X_i' \left[y_i - \exp(X_i \hat{b}) \right] = 0.$$

La distribution asymptotique de cet estimateur est normale :

$$\sqrt{N}(\hat{b} - b) \xrightarrow[N \rightarrow +\infty]{L} N(0, J_1^{-1}(b)),$$

avec

$$J_1(b) = \mathbb{E}_X \mathbb{E}_y \left[- \frac{\partial^2 \ln f}{\partial b \partial b'}(y|X) \right].$$

On remarque que la matrice hessienne ne dépend pas de y , ce qui implique qu'elle est égale à son espérance mathématique par rapport à la loi de y . On estime cette matrice par :

$$\hat{J}_1 = \frac{1}{N} \sum_{i=1}^N X_i' X_i \exp(X_i \hat{b}).$$

8.1.2.2 Pseudo maximum de vraisemblance

Comme l'espérance mathématique est identique dans les modèles homogène et hétérogène, et comme la loi de Poisson appartient à la famille exponentielle linéaire, seule la matrice de covariance asymptotique est changée. On a :

$$\sqrt{N}(\hat{b} - b) \xrightarrow[n \rightarrow +\infty]{L} N(0, J_1^{-1}(b) I_1(b) J_1^{-1}(b)),$$

où $J_1(b)$ a été définie dans la section précédente et

$$I_1(b) = \mathbb{E}_X \mathbb{E}_y \left[\frac{\partial \ln f}{\partial b}(y|X) \frac{\partial \ln f}{\partial b'}(y|X) \right].$$

On estime la matrice $I_1(b)$ par le moment empirique correspondant :

$$\begin{aligned}\widehat{I}_1 &= \frac{1}{N} \sum_{i=1}^N X_i' X_i \left(y_i - \exp(X_i \widehat{b}) \right)^2 \\ &= \frac{1}{N} \sum_{i=1}^N X_i' X_i \widehat{u}_i^2.\end{aligned}$$

Par rapport au cas homogène, les écarts-types sont robustes à la distribution de comptage parce que la loi de Poisson appartient à la famille exponentielle linéaire à l'ordre 1.

Examinons maintenant la variance conditionnelle de notre variable expliquée. On suppose que le terme d'hétérogénéité $\varepsilon_i = \exp(u_i)$ suit une loi d'espérance 1 et de variance $\eta_i > 0$. La variance de la variable expliquée est donnée par :

$$\begin{aligned}\mathbb{E}(y_i | X_i) &= \mathbb{E}_{\varepsilon} \mathbb{V}_{y_i} (y_i | X_i, \varepsilon_i) + \mathbb{V}_{\varepsilon} \mathbb{E}_{y_i} (y_i | X_i, \varepsilon_i) \\ &= \mathbb{E}_{\varepsilon} (\exp(X_i b + u_i) | X_i) + \mathbb{V}_{\varepsilon} (\exp(X_i b + u_i) | X_i) \\ &= \mathbb{E}_{\varepsilon} (\exp(X_i b) \varepsilon_i | X_i) + \mathbb{V}_{\varepsilon} (\exp(X_i b) \varepsilon_i | X_i) \\ &= \exp(X_i b) \underbrace{\mathbb{E}_{\varepsilon} (\varepsilon_i)}_1 + \exp(2X_i b) \underbrace{\mathbb{V}_{\varepsilon} (\varepsilon_i)}_{\eta_i} \\ &= \exp(X_i b) (1 + \eta_i \exp(X_i b)) \\ &= m_i (1 + \eta_i m_i).\end{aligned}$$

Ainsi, le modèle possède une variance supérieure à la moyenne et qui croît avec la moyenne.² Le maximum de vraisemblance du modèle homogène revient à supposer que $\eta_i = 0, \forall i$; le pseudo maximum de vraisemblance autorise n'importe quel profil de variance du terme d'hétérogénéité.

8.2 Le modèle binomial négatif

Bien que donnant un estimateur convergent, le modèle de Poisson ne donne pas forcément l'estimateur le plus efficace en présence d'hétérogénéité. Plusieurs approches sont possibles pour traiter ce problème. Premièrement, on peut postuler une loi pour le terme d'hétérogénéité et estimer le modèle par le maximum de vraisemblance; deuxièmement, on peut procéder à une estimation par le pseudo maximum de vraisemblance

²La relation précédente permet de voir que dans un modèle de Poisson homogène, la variance conditionnelle est toujours égale à l'espérance conditionnelle.

quasi généralisé à condition de choisir une loi ayant un paramètre de variance; troisièmement, on peut procéder à une estimation par le maximum de vraisemblance simulé. Le modèle binomial négatif peut être utilisé avec les deux premières approches.

8.2.1 Estimation par le maximum de vraisemblance

On dit qu'une variable aléatoire $y \in \mathbb{N}$ suit une loi binomiale négative de paramètres (r, p) quand elle admet pour distribution :

$$f(y) = \frac{\Gamma(y+r)}{\Gamma(y+1)\Gamma(r)} p^r (1-p)^y, \quad 0 < p < 1 \text{ et } r > 0,$$

ceci implique : $E(y) = \frac{r(1-p)}{p}$ et $V(y) = \frac{r(1-p)}{p^2}$.

Dans notre cas, nous souhaitons obtenir un modèle de Poisson hétérogène, ce qui impose la forme suivante pour les deux premiers moments :

$$E(y) = m \text{ et } V(y) = m(1 + \eta m),$$

de sorte qu'il faut prendre (r, p) tels que :

$$m = \frac{r(1-p)}{p} \text{ et } m(1 + \eta m) = \frac{r(1-p)}{p^2}.$$

En divisant la variance par l'espérance, on obtient :

$$1 + \eta m = \frac{1}{p} \Leftrightarrow p = \frac{1}{1 + \eta m} \in [0, 1] \text{ car } \eta > 0, m > 0.$$

En utilisant la définition de l'espérance, on a :

$$r = m \frac{p}{1-p} = \frac{1}{\eta},$$

on peut donc réécrire la densité de la manière suivante :

$$f(y, m, \eta) = \frac{\Gamma(y+1/\eta)}{\Gamma(y+1)\Gamma(1/\eta)} \left(\frac{1}{1+\eta m} \right)^{1/\eta} \left(\frac{\eta m}{1+\eta m} \right)^y, \quad \eta > 0, m > 0.$$

On retrouve donc les moments d'un modèle de Poisson hétérogène. Pour obtenir cette distribution, on fait les hypothèses suivantes :

1. Y suit une loi de Poisson d'espérance :

$$\tilde{m}_i = m_i \varepsilon_i,$$

2. ε_i suit une loi Gamma de paramètres $(1/\eta, \eta)$ dont la densité est donnée par :³

$$g(\varepsilon_i) = \frac{\varepsilon_i^{1/\eta-1} \exp(-\varepsilon_i/\eta)}{\eta^{1/\eta} \Gamma(1/\eta)},$$

et de moments

$$E(\varepsilon_i) = 1, \quad V(\varepsilon_i) = \eta.$$

Plus précisément, dans un modèle standard :

$$m_i = \exp(X_i b) \quad \text{et} \quad \varepsilon_i = \exp(u_i),$$

de sorte que :

$$\tilde{m}_i = m_i \varepsilon_i = \exp(X_i b + u_i),$$

ce qui donne un modèle de Poisson avec hétérogénéité Log-Gamma car $u_i = \ln \varepsilon_i$.

La loi binomiale négative s'obtient de la manière suivante. La densité du modèle hétérogène, noté f^* , est égale à :

$$f^*(y|X, \varepsilon) = \frac{\exp(-\tilde{m}) \tilde{m}^y}{y!} \quad \text{avec} \quad \tilde{m} = m\varepsilon, m = \exp(Xb).$$

Comme la variable aléatoire ε n'est pas observable, on intègre la densité précédente par rapport à cette distribution, afin d'obtenir la densité conditionnelle de Y par rapport à X .

$$f(y|X) = E_{\varepsilon}(f^*(y|X, \varepsilon)) = \int_0^{+\infty} f^*(y|X, u) g(\varepsilon) d\varepsilon.$$

On obtient donc l'expression suivante :⁴

$$\begin{aligned} f(y|X) &= \int_0^{+\infty} \frac{\exp(-m\varepsilon) (m\varepsilon)^y}{y!} \frac{\varepsilon^{1/\eta-1} \exp(-\varepsilon/\eta)}{\eta^{1/\eta} \Gamma(1/\eta)} d\varepsilon \\ &= \frac{m^y}{\Gamma(y+1) \eta^{1/\eta} \Gamma(1/\eta)} \int_0^{+\infty} \varepsilon^{y+1/\eta-1} \exp(-\varepsilon(m+1/\eta)) d\varepsilon \end{aligned}$$

³Dans le cas général, une variable aléatoire ε suit une loi gamma de paramètres (a, b) , noté $\varepsilon \rightsquigarrow \Gamma(a, b)$ si elle vérifie :

$$f(\varepsilon) = \frac{\varepsilon^{a-1} e^{-\varepsilon/b}}{b^a \Gamma(a)}, \quad a, b, \varepsilon > 0,$$

$$E(\varepsilon) = ab, \quad V(\varepsilon) = ab^2.$$

On rappelle également que :

$$\Gamma(a) = \int_0^{+\infty} x^{a-1} e^{-x} dx.$$

⁴On utilise $y! = \Gamma(y+1)$.

arrivé à ce stade on effectue le changement de variable :

$$z = (m + 1/\eta)\varepsilon,$$

ce qui implique :

$$\varepsilon = \frac{\eta z}{\eta m + 1}, \quad d\varepsilon = \frac{\eta dz}{\eta m + 1}, \quad \lim_{\varepsilon \rightarrow 0} z = 0, \quad \lim_{\varepsilon \rightarrow +\infty} z = +\infty,$$

d'où l'expression :

$$\begin{aligned} f(y|X) &= \frac{m^y}{\Gamma(y+1)\eta^{1/\eta}\Gamma(1/\eta)} \int_0^{+\infty} \left(\frac{\eta z}{\eta m + 1}\right)^{y+1/\eta-1} \exp(-z) \frac{\eta}{\eta m + 1} dz \\ &= \frac{m^y}{\Gamma(y+1)\eta^{1/\eta}\Gamma(1/\eta)} \left(\frac{\eta}{\eta m + 1}\right)^{y+1/\eta-1} \frac{\eta}{\eta m + 1} \underbrace{\int_0^{+\infty} z^{y+1/\eta-1} \exp(-z) dz}_{\Gamma(y+1/\eta)} \\ &= \frac{\Gamma(y+1/\eta)}{\Gamma(y+1)\Gamma(1/\eta)} \frac{1}{\eta^{1/\eta}} \left(\frac{\eta}{\eta m + 1}\right)^{1/\eta} \left(\frac{\eta m}{\eta m + 1}\right)^y \\ &= \frac{\Gamma(y+1/\eta)}{\Gamma(y+1)\Gamma(1/\eta)} \left(\frac{1}{\eta m + 1}\right)^{1/\eta} \left(\frac{\eta m}{\eta m + 1}\right)^y \end{aligned}$$

On utilise cette densité pour estimer les paramètres par le maximum de vraisemblance.

8.2.2 Estimation par le pseudo maximum de vraisemblance quasi généralisé

La loi binomiale négative appartient à la famille exponentielle linéaire à l'ordre 1, ce qui permet d'obtenir un estimateur convergent dit de première étape. Cet estimateur sera utilisé pour estimer η , et l'on pourra ensuite obtenir l'estimateur du pseudo maximum de vraisemblance quasi-généralisé. Tout d'abord, vérifions que la loi binomiale négative appartient à la famille exponentielle linéaire à l'ordre 1 :

$$\ln f(y, m) = \underbrace{-\frac{1}{\eta} \ln(\eta m + 1)}_{A(m)} + \underbrace{\ln \frac{\Gamma(y+1/\eta)}{\Gamma(y+1)\Gamma(1/\eta)}}_{B(y)} + y \underbrace{(\ln(\eta m) - \ln(\eta m + 1))}_{C(m)},$$

on vérifie facilement la variance :

$$\begin{aligned} V(y) &= \left(\frac{dC}{dm} \right)^{-1} \\ &= \left(\frac{1}{m} - \frac{\eta}{\eta m + 1} \right)^{-1} \\ &= \left(\frac{1}{m(\eta m + 1)} \right)^{-1} \\ &= m(1 + \eta m). \end{aligned}$$

8.2.2.1 Estimateur de première étape

Pour cet estimateur on fixe librement la valeur de $\eta > 0$. Par exemple, on fixe une valeur qui simplifie l'expression de la log-vraisemblance, $\eta = 1$, mais l'estimateur obtenu a peu de chance d'être de bonne qualité. Ceci importe toutefois peu, car c'est l'estimateur de seconde étape qui nous intéresse. Si on fixe, par exemple, $\eta = 1$ on obtient un estimateur convergent en optimisant la pseudo log vraisemblance :⁵

$$f(y|X) = \frac{\Gamma(y+1)}{\Gamma(y+1)\Gamma(1)} \left(\frac{m}{m+1} \right)^y \left(\frac{1}{m+1} \right)^1 = \frac{m^y}{(m+1)^{y+1}},$$

d'où une pseudo log vraisemblance à l'ordre 1 :

$$\ell_1(y|X, b) = \sum_{i=1}^N \{y_i \ln X_i b - (y_i + 1) \ln (\exp(X_i b) + 1)\}.$$

La maximisation de cette fonction fournit l'estimateur de première étape que l'on note \hat{b} .

8.2.2.2 Estimateur de seconde étape

Il s'agit de l'étape qui fournit le meilleur estimateur. Dans un premier temps, il faut trouver un estimateur convergent de δ . Pour cela, on utilise l'expression de la variance conditionnelle de y_i . On enlève le conditionnement pour ne pas alourdir les notations.

$$V(y_i) = m_i \left(1 + \frac{m_i}{\delta} \right) \Leftrightarrow E((y_i - m_i)^2) = m_i + \eta m_i^2. \quad (8.1)$$

A partir de cette relation, on peut proposer deux estimateurs convergents :

⁵A priori, rien n'empêcherait de prendre l'estimateur du maximum de vraisemblance lui-même, puisque l'estimateur du PMV1 de b est convergent pour toutes les valeurs de η . Cette approche a l'avantage d'être moins arbitraire.

1. Le premier se base sur une réécriture de la relation (8.1) :

$$\mathbb{E} \left((y_i - m_i)^2 - m_i \right) = \eta m_i^2,$$

et consiste à régresser $(y_i - \widehat{m}_i)^2 - \widehat{m}_i$ sur \widehat{m}_i^2 par les moindres carrés ordinaires sans terme constant, ce qui donne :

$$\eta_1 = \frac{\sum_{i=1}^N \left[(y_i - \widehat{m}_i)^2 - \widehat{m}_i \right] \widehat{m}_i^2}{\sum_{i=1}^N \widehat{m}_i^4}, \quad \widehat{m}_i = \exp \left(X_i \widehat{b} \right).$$

2. Un second estimateur est obtenu en réécrivant la relation (8.1) :

$$\mathbb{E} \left(\frac{1}{m_i^2} (y_i - m_i)^2 - m_i \right) = \eta \Leftrightarrow \eta = \mathbb{E} \left(\left(\frac{y_i}{m_i} - 1 \right)^2 - \frac{1}{m_i} \right),$$

ce qui donne simplement :

$$\eta_2 = \frac{1}{N} \sum_{i=1}^N \left[\left(\frac{y_i}{\widehat{m}_i} - 1 \right)^2 - \frac{1}{\widehat{m}_i} \right]$$

La pseudo vraisemblance quasi-généralisée est alors définie par :

$$\ell_{QG} = \sum_{i=1}^N \left\{ -\frac{1}{\widehat{\eta}} \ln(\widehat{\eta} m_i + 1) + \ln \frac{\Gamma(y_i + 1/\widehat{\eta})}{\Gamma(y_i + 1) \Gamma(1/\widehat{\eta})} + y_i [\ln(\widehat{\eta} m_i) - \ln(\widehat{\eta} m_i + 1)] \right\}$$

pour l'optimisation, on peut éliminer tous les termes qui ne dépendent pas de m_i , ce qui donne finalement :

$$\begin{aligned} \widetilde{\ell}_{QG} &= \sum_{i=1}^N \left\{ y_i \ln m_i - \left(y_i + \frac{1}{\widehat{\eta}} \right) \ln(\widehat{\eta} m_i + 1) \right\} \\ &= \sum_{i=1}^N \left\{ y_i X_i b - \left(y_i + \frac{1}{\widehat{\eta}} \right) \ln(\widehat{\eta} \exp(X_i b) + 1) \right\}. \end{aligned}$$

Le pseudo score, d'espérance nulle, est donné par :

$$\widetilde{s}_{QG} = \sum_{i=1}^N X_i' \left(y_i - \exp(X_i b) \frac{1 + \widehat{\eta} y_i}{1 + \widehat{\eta} \exp(X_i b)} \right),$$

et l'estimateur du PMVQG, \bar{b} , est défini par :

$$\sum_{i=1}^N X_i' \left(y_i - \exp(X_i \bar{b}) \frac{1 + \widehat{\eta} y_i}{1 + \widehat{\eta} \exp(X_i \bar{b})} \right) = 0.$$

Cet estimateur est convergent, asymptotiquement normal et l'on estime sa matrice de covariance par :

$$\begin{aligned}\widehat{\text{Vas}}(\bar{b}) &= \left[\widehat{J}_N(\bar{b}) \right]^{-1} \\ &= \left[\sum_{i=1}^N X_i' X_i \left(\frac{\bar{m}_i}{1 + \widehat{\eta} \bar{m}_i} \right) \right]^{-1} \quad \text{avec } \bar{m}_i = \exp(X_i \bar{b}),\end{aligned}$$

en utilisant $E(y_i | X_i) = m_i$.

8.3 Le modèle avec décision

Ce modèle généralise le modèle de Poisson en introduisant une forme explicite d'hétérogénéité. Tous les individus n'ont plus la même probabilité de rencontrer l'évènement étudié. L'observations d'une donnée de comptage fait donc apparaître deux types d'évènements nuls : ceux qui correspondent aux individus qui ne sont pas concernés par l'évènement étudié et ceux qui sont concernés mais qui n'ont pas rencontré l'évènement pendant la période étudiée. Le modèle comporte deux parties. La première partie est une relation de décision relative à l'évènement et se modélise par un modèle pour variable dichotomique. La seconde partie de cette relation donne le nombre d'évènements conditionnellement à la réalisation d'au moins un évènement et se modélise par un modèle de comptage. Le modèle latent qui représente la décision est donné par :

$$d_i^* = X_{1i} b_1 + u_i.$$

L'individu i entre dans le processus de comptage lorsque $d_i^* > 0$ lorsque $d_i^* > 0$. On a donc :

$$d_i = \begin{cases} 1 & \text{si } d_i^* > 0 \\ 0 & \text{sinon} \end{cases}.$$

Pour les individus qui sont entrés dans le processus de comptage le nombre de réalisations de l'évènement étudié z_i est distribué selon une loi de comptage $f(z_i)$. Cette loi est donnée par :

$$f(z_i | z_i \geq 1) = \Pr(d_i^* > 0) f(z_i).$$

On observe donc une réalisation nulle $y_i = 0$ soit lorsque $d_i = 0$ soit lorsque $z_i = 0$. Dans cette version simplifiée du modèle, on suppose que les perturbations u_i sont indépendantes du processus de comptage z_i , de sorte que :

$$\begin{aligned}\Pr[y_i = 0] &= \Pr[(d_i = 0) \cup (d_i = 1 \cap z_i = 0)] \\ &= \Pr(d_i^* \leq 0) + \underbrace{(1 - \Pr(d_i^* \leq 0))}_{\Pr(d_i^* > 0)} \Pr(z_i = 0).\end{aligned}$$

Pour procéder à une estimation par le maximum de vraisemblance, il faut préciser les distributions suivies par u_i et z_i . Dans le modèle originel, la distribution de u_i est logistique et celle de z_i est de Poisson. On a donc :

$$\Pr [d_i^* \leq 0] = \frac{1}{1 + \lambda_{1i}}, \Pr [d_i^* > 0] = \frac{\lambda_{1i}}{1 + \lambda_{1i}}, f(z_i) = \frac{\exp(-\lambda_{2i}) \lambda_{2i}^{z_i}}{z_i!},$$

avec

$$\lambda_{1i} = \exp(X_{1i}b_1).$$

La log-vraisemblance du modèle s'écrit donc :

$$\begin{aligned} \ell_i &= (1 - d_i) \ln \Pr [y_i = 0] + d_i \ln f(y_i | y_i \geq 1) \\ &= (1 - d_i) \ln [\Pr (d_i^* \leq 0) + \Pr (d_i^* > 0) \Pr (y_i = 0)] \\ &\quad + d_i \ln \Pr (d_i^* > 0) f(y_i) \\ &= (1 - d_i) \ln \left[\frac{1}{1 + \lambda_{1i}} + \frac{\lambda_{1i}}{1 + \lambda_{1i}} \exp(-\lambda_{2i}) \right] \\ &\quad + d_i \left[\ln \frac{\lambda_{1i}}{1 + \lambda_{1i}} + \ln \frac{\exp(-\lambda_{2i}) \lambda_{2i}^{y_i}}{y_i!} \right] \\ &= (1 - d_i) [\ln (1 + \lambda_{1i} \exp(-\lambda_{2i})) - \ln (1 + \lambda_{1i})] \\ &\quad + d_i [\ln \lambda_{1i} - \ln (1 + \lambda_{1i}) - \lambda_{2i} + y_i \ln \lambda_{2i} - \ln y_i!] \\ &= (1 - d_i) \ln (1 + \lambda_{1i} \exp(-\lambda_{2i})) \\ &\quad + d_i (\ln \lambda_{1i} - \lambda_{2i} + y_i \ln \lambda_{2i}) \\ &\quad - \ln (1 + \lambda_{1i}) \end{aligned}$$

8.4 Le modèle avec saut

Il s'agit d'un modèle qui permet également de s'écarter de la proportion de réalisations nulles données par la loi de Poisson simple. On considère que le processus qui génère les réalisations nulles diffère de celui qui génère les réalisations positives. La première partie des données $y_i = 0$ est générée par une loi de Poisson de paramètre $\lambda_{1i} = \exp(X_{1i}b_1)$ et que la partie des données $y_i \geq 1$ est générée par une loi de Poisson de paramètre $\lambda_{2i} = \exp(X_{2i}b_2)$. La probabilité d'une réalisation nulle est donc égale à :

$$\Pr (y_i = 0) = f_{1i}(0) = \exp(-\lambda_{1i}),$$

et celle d'une réalisation strictement positive est égale à :

$$f(y_i | y_i > 0) = \frac{1 - f_{1i}(0)}{1 - f_{2i}(0)} f_2(y_i),$$

et l'on remarque que l'on a :

$$\begin{aligned} f_1(0) + \sum_{y=1}^{+\infty} \frac{1 - f_1(0)}{1 - f_2(0)} &= f_1(0) + \frac{1 - f_1(0)}{1 - f_2(0)} \sum_{y=1}^{+\infty} f_2(y) \\ &= f_1(0) + \frac{1 - f_1(0)}{1 - f_2(0)} (1 - f_2(0)) \\ &= 1. \end{aligned}$$

d'où la log-vraisemblance :

$$\begin{aligned} \ell_i &= (1 - d_i) \ln f_{1i}(0) + d_i \ln \frac{1 - f_{1i}(0)}{1 - f_{2i}(0)} f_2(y_i) \\ &= \underbrace{(1 - d_i) \ln f_{1i}(0) + d_i \ln (1 - f_{1i}(0))}_{\text{Partie Poissit}} + \underbrace{d_i \ln \frac{f_{2i}(y_i)}{1 - f_{2i}(0)}}_{\text{Partie censurée}} \end{aligned}$$

avec $f_{ji}(0) = \exp(-\lambda_{ji})$, $j = 1, 2$.

La log-vraisemblance est séparable en deux parties indépendantes: la partie dichotomique ("Poissit") ne dépend que du paramètre b_1 ; la partie censurée ne dépend que du paramètre b_2 . On peut donc réaliser deux optimisations séparées pour obtenir b_1 et b_2 , les estimateurs du maximum de vraisemblance correspondants sont asymptotiquement indépendants.

La log-vraisemblance du modèle Poissit est donnée par :

$$\ell_{1i} = -(1 - d_i) \exp(X_{1i}b_1) + d_i \ln [1 - \exp(-\exp(X_{1i}b_1))],$$

et celle du modèle de Poisson censuré par :

$$\ell_{2i} = d_i \{-\exp(X_{2i}b_2) + y_i X_{2i}b_2 - \ln y_i! - \ln [1 - \exp(-\exp X_{2i}b_2)]\}.$$

CHAPITRE 9

Les variables de durée

On rencontre des variables de durée dans de nombreux cas. A l'origine, les modèles ont été développés pour étudier la durée de vie mais d'autres applications ont été mises en oeuvre. En économie, on étudie la durée passée au chômage, dans un emploi ou entre deux emplois, la durée d'un trajet de transport, la durée de vie d'une entreprise ou encore la durée d'un crédit de type "revolving". Or les variables de durée ont des caractéristiques particulières : elles sont strictement positives et souffrent souvent de problèmes de censure. En effet, l'arrêt de la collecte à une date donnée (date d'arrêt de l'alimentation du fichier) fait que des durées commencées n'ont pas eu le temps de se terminer et sont donc censurées. On peut juste affecter une valeur minimale à ces durées observées de manière incomplète. On parle de censure linéaire droite. Inversement, il est possible que l'on commence le fichier à une date où le processus observé a déjà commencé pour certains individus, la durée est alors censurée à gauche. Pour obtenir une bonne estimation, il faut tenir compte de toutes les observations, censurées ou non. En effet, plus une durée est longue plus elle a de chances d'être censurée, de sorte qu'enlever les durées censurées revient à causer un biais de sélection. Par exemple, si l'on étudie la durée du chômage, enlever les données censurées reviendrait à réaliser une étude sans les chômeurs de longue durée, ce qui est difficilement envisageable.

Comme pour les variables aléatoires réelles, on définit la loi d'une variable de durée par sa fonction de répartition. Toutefois, on préfère pour des raisons pratiques, utiliser d'autres concepts plus parlants que la fonction de répartition ou la densité. Cette pratique provient de la démographie et utilise donc des concepts spécifiques comme le taux de mortalité, la probabilité de survie ou l'espérance de vie à la naissance. Nous allons montrer que ces concepts sont rigoureusement équivalents à ceux utilisés dans les autres branches de l'économétrie.

9.1 Terminologie

Considérons une variable aléatoire de durée $T > 0$. Sa fonction de répartition est définie par la probabilité que cette durée soit inférieure à une valeur donnée t :

$$F(t) = \Pr [T \leq t], \quad t \in \mathbb{R}^{+*}.$$

Ce concept n'est pas toujours le plus pratique pour l'interprétation. L'économétrie des durées utilise, à la place, le concept de *fonction de survie* $S(t)$ qui donne la probabilité que la durée (de vie) soit supérieure à une valeur donnée t :

$$S(t) = \Pr [T > t] = 1 - F(t),$$

son nom vient de la démographie : elle donne la fraction d'individu d'une génération ayant survécu jusqu'à l'âge t .

La densité de la durée est donnée par :

$$f(t) = \frac{dF(t)}{dt} = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \Pr [t < T \leq t + \Delta t],$$

elle représente l'intensité d'occurrence d'une durée exactement égale à t . Cette intensité peut être supérieure à l'unité car il ne s'agit pas d'une probabilité mais d'une densité. La probabilité correspondante se calcule sur un petit intervalle de temps Δt , elle est donnée par :

$$\Pr [t < T \leq t + \Delta t] \simeq f(t) \Delta t.$$

Cette densité permet aussi de caractériser la loi de T car on en déduit la fonction de répartition de la manière suivante :

$$F(t) = \int_0^t f(x) dx.$$

Mais la densité est également reliée à la fonction de survie par la relation :

$$f(t) = \frac{dF(t)}{dt} = \frac{d}{dt} (1 - S(t)) = -\frac{dS(t)}{dt}.$$

La *fonction de hasard* représente une occurrence de mortalité instantanée. Comme pour la densité cette occurrence peut être supérieure à l'unité. Elle est définie comme la probabilité conditionnelle de sortir (i.e. décéder) à la date t sachant que l'on vécu jusqu'à cette date. En effet, le taux de mortalité à la date t se calcule sur la population survivante à

cette date, et non sur toute la population. On obtient :

$$\begin{aligned}
 h(t) &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \Pr [t < T \leq t + \Delta t | T > t] \\
 &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \frac{\Pr [(t < T \leq t + \Delta t) \cap (T > t)]}{\Pr [T > t]} \\
 &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \frac{\Pr [t < T \leq t + \Delta t]}{\Pr [T > t]} \\
 &= \frac{1}{\Pr [T > t]} \underbrace{\lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \Pr [t < T \leq t + \Delta t]}_{f(t)} \\
 &= \frac{f(t)}{S(t)}.
 \end{aligned}$$

Ceci permet également de calculer l'équivalent du taux de mortalité instantané en démographie, sur un intervalle de temps Δt , qui est égal à :

$$\Pr [t < T \leq t + \Delta t | T > t] \simeq h(t) \Delta t.$$

Comme la densité et la fonction de répartition, la fonction de hasard caractérise la loi de la durée T :

$$h(t) = \frac{f(t)}{S(t)} = -\frac{1}{S(t)} \frac{dS(t)}{dt} = -\frac{d \ln S(t)}{dt},$$

La fonction de hasard peut également être employée pour calculer la fonction de survie :

$$\begin{aligned}
 \int_0^t h(x) dx &= -\int_0^t \frac{d \ln S(x)}{dx} dx \\
 &= -[\ln S(x)]_0^t \\
 &= -\ln S(t) + \ln S(0) \\
 &= -\ln S(t),
 \end{aligned}$$

car pour une variable positive $S(0) = \Pr [T > 0] = 1$. On en déduit que :

$$S(t) = \exp \left\{ -\int_0^t h(x) dx \right\},$$

ce qui implique que l'on peut également écrire la densité en fonction du hasard :

$$\begin{aligned}
 f(t) &= h(t) S(t) \\
 &= h(t) \exp \left\{ -\int_0^t h(x) dx \right\}.
 \end{aligned}$$

Enfin, remarquons une propriété qui peut s'avérer utile pour calculer l'espérance de la durée (i.e. l'espérance de vie à la naissance) :

$$\begin{aligned} E(T) &= \int_0^{+\infty} x f(x) dx \\ &= \int_0^{+\infty} -x \frac{dS(x)}{dx} dx, \end{aligned}$$

en intégrant par partie ($u = -x$, $v' = S'(x)$), on obtient :

$$\begin{aligned} E(T) &= [-xS(x)]_0^{+\infty} + \int_0^{+\infty} S(x) dx \\ &= \int_0^{+\infty} S(x) dx, \end{aligned}$$

sous l'hypothèse que :

$$\lim_{x \rightarrow +\infty} xS(x) = 0,$$

et cette hypothèse est généralement bien vérifiée car les fonctions de survie contiennent souvent des exponentielles. Cette formule peut être utile, mais il faut garder à l'esprit que dans la plupart des cas la méthode la plus simple est celle de la fonction génératrice des moments que nous présenterons plus loin.

9.2 Lois usuelles

A priori, toutes les lois applicables aux variables réelles positives peuvent être utilisées pour modéliser les variables de durée. C'est l'approche employée pour justifier l'utilisation de la loi log-normale. Cependant, les méthodes les plus employées (Weibull, Gamma, Gamma Généralisée, Cox) se basent sur des modèles dits à hasard proportionnels que nous définirons plus loin. Ces modèles possèdent l'avantage de permettre une modélisation directe de la fonction de hasard.

9.2.1 La loi exponentielle

Cette loi, la plus simple, vérifie la propriété forte *d'indépendance temporelle* de la fonction de hasard. Le taux de mortalité (i.e. de sortie) est constant dans le temps :

$$h(t) = h, \quad \forall t,$$

cette hypothèse définit ce que l'on appelle un processus de Poisson (qui est également relié à la loi de Poisson dans le cas des données de comptage). En utilisant les propriétés de la section précédente, on retrouve les

différentes manières dont on peut caractériser la distribution :

$$\begin{aligned}
 S(t) &= \exp \left\{ - \int_0^t h dx \right\} \\
 &= \exp \left\{ -h \int_0^t dx \right\} \\
 &= \exp \left\{ -h [x]_0^t \right\} \\
 &= \exp(-ht), \\
 F(t) &= 1 - S(t) = 1 - \exp(-ht), \\
 f(t) &= h(t) S(t) = h \exp(-ht),
 \end{aligned}$$

et

$$\begin{aligned}
 E(T) &= \int_0^{+\infty} \exp(-hx) dx \\
 &= \left[-\frac{1}{h} \exp(-hx) \right]_0^{+\infty} \\
 &= \frac{1}{h}.
 \end{aligned}$$

Cette loi est surtout employée dans modèles d'économie théorique en raison de sa simplicité.

9.2.2 La loi de Weibull

Cette loi généralise la loi exponentielle en autorisant plusieurs type d'évolution de la fonction de hasard dans le temps, résumée dans le graphique 9.1. On remarque que ces évolutions restent toutefois monotones. On a :

$$h(t) = h\alpha t^{\alpha-1},$$

si $\alpha = 1$ on retrouve le modèle exponentiel mais, selon la valeur de α , le hasard peut être aussi bien croissant que décroissant avec la durée.

En utilisant les propriétés de la première section, on obtient les caractéristiques suivantes de la distribution de Weibull :

$$\begin{aligned}
 S(t) &= \exp \left\{ - \int_0^t h(x) dx \right\} \\
 &= \exp \left\{ - \int_0^t h\alpha x^{\alpha-1} dx \right\} \\
 &= \exp \left\{ -h [x^\alpha]_0^t \right\} \\
 &= \exp(-ht^\alpha),
 \end{aligned}$$

Figure 9.1: Fonction de hasard de la loi de Weibull

la fonction de répartition est donnée par :

$$F(t) = 1 - \exp(-ht^\alpha),$$

la densité par :

$$\begin{aligned} f(t) &= h(t)S(t) \\ &= h\alpha t^{\alpha-1} \exp(-ht^\alpha), \end{aligned}$$

et on peut également calculer l'espérance en utilisant la fonction de survie :

$$E(T) = \int_0^{+\infty} \exp(-hx^\alpha) dx,$$

on fait le changement de variable :

$$z = hx^\alpha \Leftrightarrow x = \left(\frac{z}{h}\right)^{1/\alpha} \Rightarrow dx = \frac{1}{\alpha h^{1/\alpha}} z^{1/\alpha-1} dz,$$

et les bornes restent inchangées :¹

$$\begin{aligned} E(T) &= \frac{1}{\alpha h^{1/\alpha}} \int_0^{+\infty} z^{1/\alpha-1} \exp(-z) dz \\ &= \frac{\Gamma(1/\alpha)}{\alpha h^{1/\alpha}} \\ &= h^{-1/\alpha} \Gamma(1 + \alpha). \end{aligned}$$

Cette loi est une des plus employées dans les applications économétriques.

9.2.3 La loi Gamma généralisée

Cette loi généralise la loi de Weibull, en introduisant un paramètre supplémentaire, qui permet d'obtenir une fonction de hasard non monotone. On la définit par sa densité :

$$f(t) = \frac{\alpha h^\beta t^{\alpha\beta-1} \exp(-ht^\alpha)}{\Gamma(\beta)}.$$

On retrouve la densité de la loi de Weibull pour $\beta = 1$ et celle de la loi exponentielle pour $\alpha = 1$ et $\beta = 1$. Pour les autres fonctions, nous aurons besoin des fonctions Gamma tronquées.² On note :

$$\begin{aligned} \underline{\Gamma}(a, x) &= \int_0^x u^{\alpha-1} e^{-u} du, \\ \overline{\Gamma}(a, x) &= \int_x^{+\infty} u^{\alpha-1} e^{-u} du, \end{aligned}$$

et l'on remarque que :

$$\begin{aligned} \lim_{x \rightarrow +\infty} \underline{\Gamma}(a, x) &= \Gamma(a), \\ \lim_{x \rightarrow 0} \overline{\Gamma}(a, x) &= \Gamma(a), \\ \underline{\Gamma}(a, x) + \overline{\Gamma}(a, x) &= \Gamma(a). \end{aligned}$$

En règle générale on évalue les fonctions Gamma tronquées numériquement. On en déduit la fonction de répartition de la variable de durée de la manière suivante :

$$F(t) = \frac{\alpha h^\beta}{\Gamma(\beta)} \int_0^t u^{\alpha\beta-1} \exp(-hu^\alpha) du,$$

¹ On rappelle que $\Gamma(x+1) = x\Gamma(x)$.

² En anglais : "incomplete Gamma functions".

on effectue le changement de variable $v = hu^\alpha$, ce qui implique $u = h^{-1/\alpha}v^{1/\alpha}$ et $du = h^{-1/\alpha}\alpha^{-1}v^{1/\alpha-1}dv$, et les bornes d'intégration deviennent 0 et ht^α :

$$\begin{aligned} F(t) &= \frac{\alpha h^\beta}{\Gamma(\beta)} \int_0^{ht^\alpha} v^{\beta-1/\alpha} h^{1/\alpha-\beta} \exp(-v) h^{-1/\alpha} \alpha^{-1} v^{1/\alpha-1} dv \\ &= \frac{1}{\Gamma(\beta)} \int_0^{ht^\alpha} v^{\beta-1} \exp(-v) dv \\ &= \frac{\underline{\Gamma}(\beta, ht^\alpha)}{\Gamma(\beta)}, \end{aligned}$$

on vérifie que lorsque $t \rightarrow 0$, $\underline{\Gamma}(\beta, 0) = \Gamma(0) = 0$ de sorte que $F(0) = 0$, et que $\lim_{t \rightarrow +\infty} \underline{\Gamma}(\beta, t) = \Gamma(\beta)$, de sorte que $F(t) \rightarrow 1$. La fonction de survie est donc donnée par :

$$S(t) = 1 - \frac{\underline{\Gamma}(\beta, ht^\alpha)}{\Gamma(\beta)} = \frac{\bar{\Gamma}(\beta, ht^\alpha)}{\Gamma(\beta)}.$$

Ceci ne permet pas d'obtenir de forme explicite pour la fonction de hasard parce que :

$$h(t) = \frac{f(t)}{S(t)} = \frac{\alpha h^\beta t^{\alpha\beta-1} \exp(-ht^\alpha)}{\bar{\Gamma}(\beta, ht^\alpha)}.$$

Le nom de cette distribution vient du fait qu'elle généralise la loi Gamma $\gamma(\beta, h)$, qui correspond au cas $\alpha = 1$. Pour l'espérance mathématique, le plus simple est de recourir à la fonction génératrice des moments, calculée plus loin.

9.2.4 La loi log-normale

Notons dès maintenant que cette distribution n'est pas reliée aux précédentes. Elle consiste à supposer directement que le logarithme de la variable de durée $\ln T$ suit une loi normale $N(m, \sigma^2)$. Sa densité est donc donnée par :

$$f(t) = \frac{1}{\sigma t} \varphi\left(\frac{\ln t - m}{\sigma}\right),$$

où $\varphi(t) = (2\pi)^{-1/2} \exp(-t^2/2)$ est la densité de loi normale centrée et réduite. Sa fonction de répartition est donc donnée par :

$$F(t) = \int_0^t f(x) dx = \int_0^t \frac{1}{\sigma x} \varphi\left(\frac{\ln x - m}{\sigma}\right) dx,$$

on remarque alors que $1/(\sigma x)$ est la dérivée par rapport à x de $(\ln x - m)/\sigma$, de sorte que :

$$F(t) = \left[\Phi \left(\frac{\ln x - m}{\sigma} \right) \right]_0^t = \Phi \left(\frac{\ln t - m}{\sigma} \right),$$

car :

$$\lim_{x \rightarrow 0} \Phi \left(\frac{\ln x - m}{\sigma} \right) = 0.$$

On en déduit :

$$S(t) = 1 - \Phi \left(\frac{\ln t - m}{\sigma} \right),$$

ainsi que :

$$h(t) = \frac{f(t)}{S(t)} = \frac{1}{\sigma t} \frac{\varphi \left(\frac{\ln t - m}{\sigma} \right)}{1 - \Phi \left(\frac{\ln t - m}{\sigma} \right)}.$$

Ceci implique que le hasard admet obligatoirement un maximum. En effet,

$$\begin{aligned} \sigma h'(t) &= -\frac{1}{t^2} \frac{\varphi(u)}{1 - \Phi(u)} + \frac{1}{t} \frac{\frac{1}{\sigma t} \varphi'(u) (1 - \Phi(u)) - \varphi(u) \left(-\frac{1}{\sigma t}\right) \varphi(u)}{(1 - \Phi(u))^2} \\ &= \frac{1}{\sigma t^2 (1 - \Phi(u))^2} \times \\ &\quad \left(-\sigma (1 - \Phi(u)) \varphi(u) - u \varphi(u) (1 - \Phi(u)) + \varphi(u)^2 \right) \\ &= \frac{\varphi(u)}{\sigma t^2 (1 - \Phi(u))^2} (-\sigma (1 - \Phi(u)) - u (1 - \Phi(u)) + \varphi(u)) \\ &= \frac{\varphi(u)}{\sigma t^2 (1 - \Phi(u))^2} (-(\sigma + u) (1 - \Phi(u)) + \varphi(u)), \end{aligned}$$

on pose :

$$\kappa(u) = \frac{\varphi(u)}{1 - \Phi(u)},$$

cette fonction est définie par analogie avec la fonction de hasard mais sur une loi normale centrée réduite. Notons qu'il ne s'agit pas à proprement parler d'une fonction de hasard parce que u peut prendre des valeurs négatives. Cette fonction est strictement croissante, comme le montre le graphique 9.2 :

On peut écrire :

$$h'(t) = \frac{1}{\sigma^2 t^2} \kappa(u) (\kappa(u) - (\sigma + u)),$$

Figure 9.2: $\varphi(x)/(1 - \Phi(x))$

le hasard atteint son maximum à un point $u = (\ln t - m)/\sigma$ tel que :

$$\kappa(u) = \sigma + u.$$

Ce type de profil est très particulier, car la présence d'un maximum est imposée, et cette hypothèse ne convient pas forcément à tous les processus de durée. Il faut donc être vigilant quand on l'emploie.

9.3 Modélisation en logarithmes

Les variables de durée peuvent toujours être prises en logarithmes, ce qui facilite l'interprétation des résultats quand les variables explicatives sont elles-mêmes en logarithmes ou sous forme d'incatrices. Mais cette modélisation peut également être utilisée pour mieux comprendre les relations entre les différentes loi usuelles, et notamment les loi exponentielles, de Weibull, Gamma et Gamma généralisée.

9.3.1 Rappels

9.3.1.1 Le changement de variable

Nous allons utiliser cette propriété dans toute la section. Supposons que l'on dispose d'une variable de durée T de densité $f_T(t)$ et que l'on effectue un changement de variable $U = g(T)$, la densité de la variable u est donnée par :

$$f_U(u) = \left| \frac{dg^{-1}(u)}{du} \right| f_T(g^{-1}(u)).$$

9.3.1.2 La loi Gamma

Une variable aléatoire X suit une loi Gamma de paramètres (a, b) , notée $\gamma(a, b)$ si sa densité s'écrit :

$$f_X(x) = \frac{b^a x^{a-1} \exp(-bx)}{\Gamma(a)}, \quad x > 0, \quad a > 0, \quad b > 0,$$

les deux premiers moments sont égaux à $E(X) = a/b$ et $V(X) = a/b^2$. Si $a = 1$, on retrouve la loi exponentielle, notée $\gamma(1, b)$, dont la densité est égale à :

$$f_X(x) = b \exp(-bx),$$

et dont les deux premiers moments sont égaux à $E(X) = 1/b$ et $V(X) = 1/b^2$. Si on prend le cas symétrique, une loi Gamma $\gamma(a, 1)$ on obtient la densité :

$$f_X(x) = \frac{x^{a-1} \exp(-x)}{\Gamma(a)},$$

dont les deux premiers moments sont égaux à $E(X) = V(X) = a$.

9.3.2 Modèle exponentiel et loi de Gumbel

Posons le modèle en logarithmes suivant :

$$\ln T = -\ln h + U,$$

où U est une variable aléatoire dont on cherche la loi. On sait seulement que la durée T suit une loi exponentielle $\gamma(1, h)$ de densité :

$$f_T(t) = h \exp(-ht).$$

Pour trouver la densité de la loi de U , on remarque que :

$$U = \ln(hT) = g(T),$$

de sorte que :

$$T = \frac{1}{h} \exp(U) = g^{-1}(U)$$

ce qui implique :

$$\frac{dg^{-1}(u)}{du} = \frac{1}{h} \exp(u),$$

d'où la densité :

$$\begin{aligned} f_U(u) &= \frac{1}{h} \exp(u) h \exp \left[-h \left(\frac{1}{h} \exp(u) \right) \right] \\ &= \exp(u) \exp(-\exp(u)), \end{aligned}$$

qui n'est autre que la densité d'une loi de Gumbel (i.e. valeur extrême de type I, minimum) d'espérance égale à l'opposée de la constante d'Euler, $E(U) = -\gamma_E$, avec $\gamma_E \simeq 0,57721$, et de variance $\pi^2/6$. Pour trouver directement ces résultats on peut utiliser les deux propriétés suivantes de la constante d'Euler :³

$$\gamma_E = - \int_0^{+\infty} (\ln x) e^{-x} dx,$$

et

$$\gamma_E^2 + \frac{\pi^2}{6} = \int_0^{+\infty} (\ln x)^2 e^{-x} dx.$$

On utilise également les propriétés suivantes de la fonction Gamma :

$$\Gamma(p) = \int_0^{+\infty} x^{p-1} e^{-x} dx \Rightarrow \Gamma'(p) = \int_0^{+\infty} (\ln x) x^{p-1} e^{-x} dx, \quad (9.1)$$

car on dérive par rapport à p et non par rapport à x .⁴ Ceci implique :

$$\Gamma'(1) = \int_0^{+\infty} (\ln x) e^{-x} dx = -\gamma_E.$$

En dérivant une nouvelle fois la relation (9.1) par rapport à p , on obtient :

$$\Gamma''(p) = \int_0^{+\infty} (\ln x)^2 x^{p-1} e^{-x} dx,$$

³Pour évaluer la constante d'Euler, on peut utiliser la définition sous forme de série donnée à l'origine par Euler lui-même :

$$\gamma_E = \sum_{k=1}^{+\infty} \left[\frac{1}{k} - \ln \left(1 + \frac{1}{k} \right) \right].$$

⁴On a :

$$\begin{aligned} \frac{d}{dp} (x^{p-1}) &= \frac{d}{dp} (e^{(p-1) \ln x}) = (\ln x) e^{(p-1) \ln x} \\ &= (\ln x) x^{p-1}. \end{aligned}$$

ce qui implique :

$$\Gamma''(1) = \gamma_E^2 + \frac{\pi^2}{6},$$

Nous retrouverons les valeurs $\Gamma'(1)$ et $\Gamma''(1)$ lors de l'étude de la fonction génératrice des moments de la loi de Gumbel. Finalement, on peut réécrire le modèle exponentiel sous la forme :

$$\begin{aligned} E(\ln T) &= -\ln h + E(U) \\ &= -(\ln h + \gamma_E), \end{aligned}$$

de sorte qu'en mettant les variables explicatives dans la fonction de hasard, on peut aboutir à un modèle log-linéaire avec une simple correction pour le terme constant du modèle. On remarque également que plus le taux de hasard h est élevé, plus l'espérance de durée est faible.

9.3.3 Modèle exponentiel et loi exponentielle

On peut également définir le modèle exponentiel directement en niveaux et non en logarithmes. C'est l'approche qui est suivie habituellement pour généraliser ce modèle vers les modèles de Weibull, Gamma et Gamma généralisé. De manière cohérente avec la section précédente, on pose :

$$T = g(V) = \frac{V}{h}, \quad (9.2)$$

où V suit une loi exponentielle de paramètre 1, notée $\gamma(1, 1)$, de densité :

$$f_V(v) = \exp(-v).$$

La densité de la variable de durée tirée de ce modèle est donnée par la transformation :

$$\begin{aligned} T = \frac{V}{h} &\Leftrightarrow V = hT = g^{-1}(T) \\ &\Rightarrow \frac{dg^{-1}(t)}{dt} = h, \end{aligned}$$

d'où la densité :

$$f_T(t) = \left| \frac{dg^{-1}(t)}{dt} \right| f_V(g^{-1}(t)) = h \exp(-ht),$$

qui correspond à la densité de la loi exponentielle $\gamma(1, h)$.

9.3.4 Modèle de Weibull

Il existe différentes manières de généraliser le modèle exponentiel (9.2). Une première manière consiste à introduire un paramètre d'échelle $\alpha > 0$ dans la définition de la variable de durée :⁵

$$T = g(V) = \left(\frac{V}{h}\right)^{1/\alpha}, \quad (9.3)$$

on retrouve le modèle exponentiel pour $\alpha = 1$. On suppose toujours que V suit une loi exponentielle $\gamma(1, 1)$. La loi suivie par T a donc changé puisque l'on a :

$$\begin{aligned} T = \left(\frac{V}{h}\right)^{1/\alpha} &\Leftrightarrow V = hT^\alpha = g^{-1}(T) \\ \Rightarrow \frac{dg^{-1}(t)}{dt} &= h\alpha t^{\alpha-1}, \end{aligned}$$

de sorte que la densité de T s'écrit :

$$\begin{aligned} f_T(t) &= \left| \frac{dg^{-1}(t)}{dt} \right| f_V(g^{-1}(t)) \\ &= \alpha h t^{\alpha-1} \exp(-ht^\alpha), \end{aligned}$$

qui correspond à la densité d'une variable de Weibull. On remarque qu'en logarithme la relation peut s'écrire :

$$\begin{aligned} \ln T &= \frac{1}{\alpha} (-\ln h + \ln V) \\ &= \frac{1}{\alpha} (-\ln h + U), \end{aligned}$$

de sorte qu'avec nos notations α est un paramètre d'échelle qui porte sur l'ensemble du modèle. L'espérance mathématique correspondante s'écrit :

$$E(\ln T) = -\frac{1}{\alpha} (\ln h + \gamma_E),$$

car U suit toujours une loi $\gamma(1, 1)$ comme dans le modèle exponentiel.

9.3.5 Modèle Gamma

Le modèle Gamma généralise le modèle exponentiel (9.2) en changeant la distribution du terme d'erreur V au lieu d'introduire un paramètre d'échelle. On suppose, comme dans le modèle exponentiel, que :

$$T = g(V) = \frac{V}{h},$$

⁵ Les notations utilisées ici sont un peu différentes de celles utilisées habituellement, c'est pour pouvoir simplifier les expressions qui apparaîtront dans la suite du chapitre.

mais cette fois-ci V suit une loi Gamma $\gamma(\beta, 1)$. La densité de V est donc donnée par :

$$f_V(v) = \frac{v^{\beta-1} e^{-v}}{\Gamma(\beta)}, \quad \beta > 0, \quad v > 0,$$

et l'on retrouve le modèle exponentiel en posant $\beta = 1$. En utilisant $g^{-1}(t) = ht$, on obtient :

$$\begin{aligned} f_T(t) &= \left| \frac{dg^{-1}(t)}{dt} \right| f_V(g^{-1}(t)) \\ &= h \frac{(ht)^{\beta-1} e^{-ht}}{\Gamma(\beta)} \\ &= \frac{h^\beta t^{\beta-1} e^{-ht}}{\Gamma(\beta)}, \end{aligned}$$

qui correspond à la densité d'une loi Gamma $\gamma(\beta, h)$. Écrit en espérance le modèle log linéaire donne :

$$E(\ln T) = -\ln h + E(\ln V),$$

et nous calculerons plus loin l'espérance de $\ln V$ à partir de sa fonction génératrice des moments.

9.3.6 Modèle Gamma généralisé

Il s'agit d'une troisième généralisation du modèle exponentiel (9.2). Cette fois-ci, nous allons combiner les deux généralisations du modèle de Weibull et du modèle Gamma. On suppose, comme dans le modèle de Weibull, que la variable de durée est définie par la relation :

$$T = g(V) = \left(\frac{V}{h} \right)^{1/\alpha},$$

et, comme dans le modèle Gamma, que V suit une loi Gamma $\gamma(\beta, 1)$ de densité :

$$f_V(v) = \frac{v^{\beta-1} e^{-v}}{\Gamma(\beta)}, \quad \beta > 0, \quad v > 0.$$

On trouve directement la nouvelle densité :

$$\begin{aligned} f_T(t) &= \left| \frac{dg^{-1}(t)}{dt} \right| f_V(g^{-1}(t)) \\ &= h\alpha t^{\alpha-1} \frac{(ht^\alpha)^{\beta-1} e^{-(ht^\alpha)}}{\Gamma(\beta)} \\ &= \frac{\alpha h^\beta t^{\alpha\beta-1} \exp(-ht^\alpha)}{\Gamma(\beta)}. \end{aligned}$$

On note cette distribution $\gamma(\beta, h, \alpha)$. On remarque que la distribution exponentielle s'obtient pour $\gamma(1, h, 1)$, la distribution de Weibull $\gamma(1, h, \alpha)$ et la distribution Gamma pour $\gamma(\beta, h, 1)$. Le nom de cette distribution vient du fait qu'elle généralise la loi Gamma puisque, pour $\alpha = 1$, on obtient :

$$f_T(t) = \frac{h^\beta t^{\beta-1} \exp(-ht)}{\Gamma(\beta)},$$

la densité de la loi Gamma à deux paramètres $\gamma(\beta, h)$. Le nom est toutefois trompeur, puisque la loi Gamma généralisée généralise également la loi de Weibull.

Plus généralement, on obtient les cas particuliers suivants :

- $\alpha = 1 : T \rightsquigarrow \gamma(\beta, h, 1)$. Loi Gamma;
- $\beta = 1 : T \rightsquigarrow \gamma(1, h, \alpha)$. Loi de Weibull;
- $\alpha = 1$ et $\beta = 1 : T \rightsquigarrow \gamma(1, h, 1)$. Loi exponentielle;
- $\alpha \neq 1$ et $\beta \neq 1 : T \rightsquigarrow \gamma(\beta, h, \alpha)$. Loi Gamma généralisée.

Le modèle log linéaire en espérance s'écrit maintenant :

$$E(\ln T) = \frac{1}{\alpha} (-\ln h + E(\ln V)),$$

où $E(\ln V)$ prend la même valeur que pour le modèle Gamma de la section précédente.

9.3.7 Modèle log-normal

On peut également utiliser la méthode du changement de variable pour le modèle log-normal, mais ici la modélisation ne porte pas sur le hasard mais sur l'espérance mathématique de la variable de durée, quantité qui est décroissante avec le taux de hasard. Dans un modèle avec des variables explicatives, un modèle basé sur l'espérance de la variable de durée implique généralement un changement de signe des coefficients par rapport à un modèle basé sur la fonction de hasard. On pose :

$$\ln T = m + \sigma U$$

où U suit une loi normale centrée et réduite, de sorte que $\ln T$ suit une loi normale $N(m, \sigma^2)$. On a donc la transformation suivante :

$$T = \exp(m + \sigma U) = g(U) \Leftrightarrow U = \frac{\ln T - m}{\sigma} = g^{-1}(T),$$

ce qui implique :

$$\frac{dg^{-1}(t)}{dt} = \frac{1}{\sigma t}.$$

D'autre part, U suit une loi normale de densité $\varphi(u)$, ce qui implique que la densité de T (i.e. de la loi log-normale) est donnée par :

$$f_T(t) = \left| \frac{dg^{-1}(t)}{dt} \right| \varphi(g^{-1}(t)) = \frac{1}{\sigma t} \varphi\left(\frac{\ln t - m}{\sigma}\right).$$

Ici, par construction, l'espérance du modèle log linéaire se réduit à sa plus simple expression :

$$E(\ln T) = m + \sigma E(U) = m.$$

9.4 Calcul des moments

Lors des estimations, on peut avoir besoin de faire une prévision de la durée moyenne passée dans l'état étudié, ainsi que de calculer la variance pour obtenir un intervalle de confiance à partir de la méthode de Slutsky. La méthode la plus simple avec les distributions qui précèdent est, souvent, d'utiliser les fonctions génératrices des moments.

9.4.1 Fonction génératrice des moments

9.4.1.1 Moments non centrés

La fonction génératrice des moments non centrés d'une variable aléatoire T est définie par :

$$M(s) = E(e^{sT}) = \int_0^{+\infty} e^{st} f_T(t) dt,$$

quand cette intégrale converge. On remarque que :

$$M(0) = E(e^0) = 1.$$

On vérifie que :

$$M'(s) = E(Te^{sT}),$$

et, par récurrence, que :

$$\frac{d^k M}{ds^k}(s) = E(T^k e^{sT}).$$

En prenant la quantité précédente en $s = 0$, on obtient :

$$\frac{d^k M}{ds^k}(0) = E(T^k).$$

On peut donc obtenir les moments non centrés par dérivation au lieu de procéder par intégration, ce qui est plus rapide.

9.4.1.2 Moments centrés

De la même manière, on peut obtenir certains moments centrés par la fonction :

$$K(s) = \ln M(s),$$

on voit que :

$$K'(s) = \frac{M'(s)}{M(s)},$$

en prenant la fonction précédente au point $s = 0$ on a :

$$K'(0) = E(T),$$

en dérivant la fonction K une deuxième fois, on obtient:

$$K''(s) = \frac{M''(s)}{M(s)} - \left(\frac{M'(s)}{M(s)} \right)^2,$$

en la prenant au point $s = 0$:

$$K''(0) = E(T^2) - E(T)^2 = V(T).$$

9.4.1.3 Moments du logarithme

Enfin, en économétrie des durées, on a souvent besoin des fonctions génératrices des moments du logarithme d'une variable de durée, parce ce que beaucoup de modèles peuvent s'écrire en logarithmes (e.g., Exponentiel, Weibull, Gamma, Gamma généralisé et log-Normal) :

$$M_{\ln T}(s) = E(e^{s \ln T}) = E(T^s),$$

dont l'intérêt est ici évident. Pour calculer la variance du logarithme d'une variable de durée, on utilisera la fonction correspondante des moments centrés :

$$K_{\ln T}(s) = \ln K_{\ln T}(s) = \ln E(T^s).$$

9.4.2 Moments des lois usuelles

9.4.2.1 Loi exponentielle

La densité est donnée par :

$$f(t) = h \exp(-ht), \quad t > 0, h > 0,$$

on a donc :

$$\begin{aligned}
 M(s) &= \mathbb{E}(e^{sT}) \\
 &= h \int_0^{+\infty} e^{-(h-s)t} dt \\
 &= h \left[-\frac{1}{h-s} e^{-(h-s)t} \right]_0^{+\infty} \\
 &= \frac{h}{h-s},
 \end{aligned}$$

remarquons bien ici que l'on a choisi de mettre dans l'exponentielle un terme en $h-s$ car on utilise cette fonction en $s=0$, ce qui garantit que $h-s > 0$ et donc la convergence de l'intégrale.

On en déduit :

$$\begin{aligned}
 M'(s) &= \frac{h}{(h-s)^2} \Rightarrow \mathbb{E}(T) = M'(0) = \frac{1}{h}, \\
 M''(s) &= \frac{2h}{(h-s)^3} \Rightarrow \mathbb{E}(T^2) = M''(0) = \frac{2}{h^2},
 \end{aligned}$$

on pourrait en déduire la variance par la formule classique :

$$\begin{aligned}
 \mathbb{V}(T) &= \mathbb{E}(T^2) - \mathbb{E}(T)^2 \\
 &= 2/h^2 - (1/h)^2 \\
 &= 1/h^2,
 \end{aligned}$$

mais on peut l'obtenir plus directement par la fonction génératrice des moments centrés. Elle est définie par :

$$K(s) = \ln h - \ln(h-s).$$

On en déduit :

$$\begin{aligned}
 K'(s) &= \frac{1}{h-s} \Rightarrow K'(0) = \mathbb{E}(T) = \frac{1}{h}, \\
 K''(s) &= \frac{1}{(h-s)^2} \Rightarrow K''(0) = \mathbb{V}(T) = \frac{1}{h^2}.
 \end{aligned}$$

Les fonctions génératrices correspondant au logarithme de la durée $\ln T$ sont celles de la loi de Gumbel données dans la section suivante.

9.4.2.2 Loi de Gumbel

On peut calculer la fonction génératrice de la loi de Gumbel en remarquant qu'une variable de ce type s'obtient comme le logarithme d'une

variable exponentielle $\gamma(1, 1)$ et en utilisant la propriété :

$$K_{\ln T}(s) = E(T^s).$$

Il suffit donc de calculer le moment d'ordre s de la loi exponentielle. En fait, dans ce cas particulier, il n'y a pas de calcul à faire, puisque l'on a :

$$E(T^s) = \int_0^{+\infty} t^s \exp(-t) dt = \Gamma(1+s),$$

on en déduit la fonction génératrice des moments centrés du logarithme de la loi exponentielle :

$$K_{\ln T}(s) = \ln \Gamma(1+s)$$

ce qui implique :

$$\begin{aligned} K'_{\ln T}(s) &= \frac{\Gamma'(1+s)}{\Gamma(1+s)} \\ \Rightarrow K'_{\ln T}(0) &= E(\ln T) = \Gamma'(1) = -\gamma_E, \end{aligned}$$

où γ_E est la constante d'Euler. De même, on voit que :

$$K''_{\ln T}(s) = \frac{\Gamma''(1+s)}{\Gamma(1+s)} - \left(\frac{\Gamma'(1+s)}{\Gamma(1+s)} \right)^2,$$

de sorte que :

$$\begin{aligned} K''_{\ln T}(0) &= V(\ln T) \\ &= \Gamma''(1) - \Gamma'(1)^2 \\ &= \gamma_E^2 + \frac{\pi^2}{6} - (-\gamma_E)^2 \\ &= \frac{\pi^2}{6}. \end{aligned}$$

9.4.2.3 Loi Gamma

La fonction génératrice de la loi Gamma $\gamma(\beta, 1)$ est définie par :

$$\begin{aligned} M_T(s) &= E(e^{sT}) \\ &= \int_0^{+\infty} e^{st} \frac{t^{\beta-1} e^{-t}}{\Gamma(\beta)} dt \\ &= \int_0^{+\infty} \frac{t^{\beta-1} e^{-(1-s)t}}{\Gamma(\beta)} dt, \end{aligned}$$

on effectue donc le changement de variable $x = (1 - s)t$, de sorte que les bornes sont inchangées et que $dt = (1 - s)^{-1}dx$. On obtient donc :

$$\begin{aligned} M_T(s) &= \frac{1}{\Gamma(\beta)} \int_0^{+\infty} \left(\frac{x}{1-s}\right)^{\beta-1} e^{-x} \frac{1}{1-s} dx \\ &= \frac{1}{(1-s)^\beta \Gamma(\beta)} \underbrace{\int_0^{+\infty} x^{\beta-1} e^{-x} dx}_{\Gamma(\beta)} \\ &= (1-s)^{-\beta}. \end{aligned}$$

Pour obtenir l'espérance et la variance de la loi Gamma, on utilise donc :

$$K_T(s) = \ln M(s) = -\beta \ln(1-s).$$

La dérivée première donne l'espérance de la distribution :

$$K'_T(s) = \frac{\beta}{1-s} \Rightarrow K'_T(0) = E(T) = \beta,$$

et la dérivée seconde donne la variance :

$$K''_T(s) = \frac{\beta}{(1-s)^2} \Rightarrow K''_T(0) = V(T) = \beta.$$

La fonction génératrice du logarithme de cette variable s'obtient par :

$$\begin{aligned} M_{\ln T}(s) &= E(T^s) \\ &= \int_0^{+\infty} \frac{t^{s+\beta-1} e^{-t}}{\Gamma(\beta)} dt \\ &= \frac{\Gamma(s+\beta)}{\Gamma(\beta)}, \end{aligned}$$

et pour trouver ses deux premiers moments on utilise :

$$K_{\ln T}(s) = \ln M_{\ln T}(s) = \ln \Gamma(s+\beta) - \ln \Gamma(\beta).$$

Pour le modèle en logarithmes, on utilise donc :

$$K'_{\ln T}(s) = \frac{\Gamma'(s+\beta)}{\Gamma(s+\beta)} \Rightarrow K'_{\ln T}(0) = E(\ln T) = \frac{\Gamma'(\beta)}{\Gamma(\beta)},$$

ainsi que :

$$\begin{aligned} K''_{\ln T}(s) &= \frac{\Gamma''(s+\beta)}{\Gamma(s+\beta)} - \left(\frac{\Gamma'(s+\beta)}{\Gamma(s+\beta)}\right)^2 \\ &\Rightarrow K''_{\ln T}(0) = V(\ln T) = \frac{\Gamma''(\beta)}{\Gamma(\beta)} - \left(\frac{\Gamma'(\beta)}{\Gamma(\beta)}\right)^2. \end{aligned}$$

9.4.2.4 Loi de Weibull

Pour trouver la fonction génératrice des moments de la loi de Weibull, il suffit de remarquer que $U = \ln hT^\alpha$ suit une loi de Gumbel de fonction génératrice des moments égale à $\Gamma(1+s)$. Ceci implique :

$$\begin{aligned} M_U(s) &= \mathbb{E}\left(e^{s \ln(hT^\alpha)}\right) = \Gamma(1+s) \\ &\Leftrightarrow \mathbb{E}(h^s T^{\alpha s}) = \Gamma(1+s) \\ &\Leftrightarrow \mathbb{E}(T^{\alpha s}) = h^{-s} \Gamma(1+s) \\ &\Leftrightarrow \mathbb{E}(T^j) = h^{-j/\alpha} \Gamma(1+j/\alpha), \end{aligned}$$

avec $j = \alpha s$. On en déduit :

$$\mathbb{E}(T) = h^{-1/\alpha} \Gamma(1+1/\alpha) \text{ et } \mathbb{E}(T^2) = h^{-2/\alpha} \Gamma(1+2/\alpha),$$

d'où la variance :

$$\begin{aligned} V(T) &= \mathbb{E}(T^2) - \mathbb{E}(T)^2 \\ &= h^{-2/\alpha} \left(\Gamma(1+2/\alpha) - \Gamma(1+1/\alpha)^2 \right). \end{aligned}$$

On remarque que l'espérance mathématique peut également se simplifier par la formule :

$$\begin{aligned} \Gamma(x) &= (x-1) \Gamma(x-1) \\ &\Rightarrow \Gamma(1+1/\alpha) = \frac{1}{\alpha} \Gamma\left(\frac{1}{\alpha}\right), \end{aligned}$$

ce qui implique :

$$\mathbb{E}(T) = h^{-1/\alpha} \frac{1}{\alpha} \Gamma\left(\frac{1}{\alpha}\right).$$

L'expression des moments du logarithme de la durée est également utilisée dans les applications. D'après ce qui précède :

$$M_{\ln T}(s) = \mathbb{E}(T^s) = h^{-s/\alpha} \Gamma(1+s/\alpha),$$

de sorte que :

$$\begin{aligned} K_{\ln T}(s) &= \ln M_{\ln T}(s) \\ &= -\frac{s}{\alpha} \ln h + \ln \Gamma(1+s/\alpha), \end{aligned}$$

d'où les dérivées :

$$\begin{aligned} K'_{\ln T}(s) &= \frac{1}{\alpha} \left(-\ln h + \frac{\Gamma'(1 + s/\alpha)}{\Gamma(1 + s/\alpha)} \right) \\ \Rightarrow K'_{\ln T}(0) &= \frac{1}{\alpha} (-\ln h + \Gamma'(1)) \\ \Leftrightarrow E(\ln T) &= (-\ln h + \Gamma'(1)) / \alpha \\ \Leftrightarrow E(\ln T) &= -\frac{1}{\alpha} (\ln h + \gamma_E), \end{aligned}$$

et que :

$$\begin{aligned} K''_{\ln T}(s) &= \frac{1}{\alpha^2} \left(\Gamma''(1) - \Gamma'(1)^2 \right) \\ \Leftrightarrow K''_{\ln T}(0) &= \frac{1}{\alpha^2} \left(\Gamma''(1) - \Gamma'(1)^2 \right) \\ \Leftrightarrow V(\ln T) &= \frac{\pi^2}{6\alpha^2}. \end{aligned}$$

9.4.2.5 Loi Gamma généralisée

Pour définir le modèle Gamma généralisé, on suppose que $hT^\alpha = V$ suit une loi Gamma $\gamma(\beta, 1)$. Les moments de cette loi peuvent être trouvés directement en utilisant :

$$E(V^j) = \frac{\Gamma(j + \beta)}{\Gamma(\beta)},$$

ce qui implique :

$$\begin{aligned} E\left((hT^\alpha)^j\right) &= \frac{\Gamma(j + \beta)}{\Gamma(\beta)} \\ \Leftrightarrow E(T^{\alpha j}) &= h^{-j} \frac{\Gamma(j + \beta)}{\Gamma(\beta)}, \end{aligned}$$

il suffit alors de poser $s = \alpha j$ ($\Leftrightarrow j = s/\alpha$) pour obtenir :

$$E(T^s) = h^{-s/\alpha} \frac{\Gamma(s/\alpha + \beta)}{\Gamma(\beta)},$$

ce qui permet d'obtenir les moments non centrés :

$$E(T) = \frac{h^{-1/\alpha} \Gamma(1/\alpha + \beta)}{\Gamma(\beta)} \text{ et } E(T^2) = \frac{h^{-2/\alpha} \Gamma(2/\alpha + \beta)}{\Gamma(\beta)},$$

dont on déduit la variance par la formule classique :

$$\begin{aligned} V(T) &= E(T^2) - E(T)^2 \\ &= h^{-2/\alpha} \left\{ \frac{\Gamma(2/\alpha + \beta)}{\Gamma(\beta)} - \left(\frac{\Gamma(1/\alpha + \beta)}{\Gamma(\beta)} \right)^2 \right\}. \end{aligned}$$

Pour obtenir les moments du logarithme de la variable de durée, il suffit de remarquer que la fonction $E(T^s)$ est identique à $M_{\ln T}(s)$ de sorte que l'on peut écrire la fonction génératrice des moments centrés :

$$\begin{aligned} K_{\ln T}(s) &= \ln M_{\ln T}(s) \\ &= -\frac{s}{\alpha} \ln h + \ln \Gamma(s/\alpha + \beta) - \ln \Gamma(\beta), \end{aligned}$$

d'où l'espérance du logarithme :

$$\begin{aligned} K'_{\ln T}(s) &= -\frac{1}{\alpha} \ln h + \frac{1}{\alpha} \frac{\Gamma'(s/\alpha + \beta)}{\Gamma(s/\alpha + \beta)} \\ \Leftrightarrow K'_{\ln T}(0) &= E(\ln T) = \frac{1}{\alpha} \left(-\ln h + \frac{\Gamma'(\beta)}{\Gamma(\beta)} \right), \end{aligned}$$

et sa variance :

$$\begin{aligned} K''_{\ln T}(s) &= \frac{1}{\alpha^2} \left[\frac{\Gamma''(s/\alpha + \beta)}{\Gamma(s/\alpha + \beta)^2} - \left(\frac{\Gamma'(s/\alpha + \beta)}{\Gamma(s/\alpha + \beta)} \right)^2 \right] \\ \Leftrightarrow K''_{\ln T}(0) &= V(\ln T) = \frac{1}{\alpha^2} \left[\frac{\Gamma''(\beta)}{\Gamma(\beta)^2} - \left(\frac{\Gamma'(\beta)}{\Gamma(\beta)} \right)^2 \right]. \end{aligned}$$

9.4.2.6 Loi normale

La loi log-normale n'admet pas de fonction génératrice des moments parce que l'intégrale qui la définit n'est pas convergente :

$$K_T(s) = +\infty,$$

mais on peut calculer tous les moments de cette loi en utilisant la fonction génératrice des moments de loi normale. C'est ce qui explique la présence de cette section. Soit une variable aléatoire X suivant une loi normale $N(m, \sigma^2)$, sa fonction génératrice des moments peut être obtenue de la manière suivante :

$$\begin{aligned} K_X(s) &= E(e^{sX}) \\ &= \int_{-\infty}^{+\infty} e^{sx} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-m}{\sigma}\right)^2} dx, \end{aligned}$$

on effectue donc le changement de variable $z = (x - m)/\sigma$, ce qui ne change pas la valeur des bornes et implique $x = m + \sigma z$ et $dx = \sigma dz$. On

a donc :

$$\begin{aligned}
 M_X(s) &= \int_{-\infty}^{+\infty} e^{s(m+\sigma z)} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \sigma dz \\
 &= e^{sm} \int_{-\infty}^{+\infty} e^{s\sigma z} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} dz \\
 &= e^{sm} \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z^2-2s\sigma z)} dz \\
 &= e^{sm+s^2\sigma^2/2} \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z^2-2s\sigma z+s^2\sigma^2)} dz \\
 &= e^{sm+s^2\sigma^2/2} \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z-s\sigma)^2} dz,
 \end{aligned}$$

et le second terme de l'expression ci-dessus est la somme sur son support de la densité d'une loi normale $N(s\sigma, 1)$, qui est égale à 1 par définition. On a donc :

$$M_X(s) = \exp(sm + s^2\sigma^2/2).$$

On peut en déduire tous les moments de la loi normale. Pour obtenir les moments centrés, on prend :

$$\begin{aligned}
 K_X(s) &= \ln M_X(s) \\
 &= sm + s^2\sigma^2/2,
 \end{aligned}$$

on en déduit :

$$\begin{aligned}
 K'_X(s) &= m + s\sigma^2 \\
 \Rightarrow K'_X(0) &= E(X) = m,
 \end{aligned}$$

ainsi que :

$$\begin{aligned}
 K''_X(s) &= \sigma^2 \\
 \Rightarrow K''_X(0) &= V(X) = \sigma^2.
 \end{aligned}$$

9.4.2.7 Loi log-normale

Par définition, la fonction génératrice du logarithme de la variable de durée $\ln T$ est identique à celle de la loi normale $N(m, \sigma^2)$ donnée par :

$$K_X(s) = K_{\ln T}(s) = E(e^{s \ln T}) = E(T^s),$$

on en déduit les moments de la loi log-normale :

$$E(T^s) = \exp(sm + s^2\sigma^2/2),$$

ce qui donne pour espérance :

$$E(T) = \exp(m + \sigma^2/2),$$

et pour variance :

$$\begin{aligned} V(T) &= E(T^2) - E(T)^2 \\ &= \exp(2m + 2\sigma^2) - \exp(2m + \sigma^2) \\ &= \exp(2m + \sigma^2) (\exp(\sigma^2) - 1). \end{aligned}$$

9.4.3 Résumé

Le tableau suivant résume les hypothèses qu'il faut effectuer pour retrouver chacun des modèles à partir de la relation suivante :

$$\ln T = k_1 + k_2 U,$$

le lecteur notera que les variables explicatives influençant la fonction de hasard ou la durée moyenne se trouvent dans la partie k_1 de ce modèle.

Modèle	Hypothèses
Exponentiel	$k_1 = -\ln h, k_2 = 1, \exp(U) \rightsquigarrow \gamma(1, 1)$
Weibull	$k_1 = -\alpha^{-1} \ln h, k_2 = \alpha^{-1}, \exp(U) \rightsquigarrow \gamma(1, 1)$
Gamma	$k_1 = -\ln h, k_2 = 1, \exp(U) \rightsquigarrow \gamma(\beta, 1)$
Gamma généralisé	$k_1 = -\alpha^{-1} \ln h, k_2 = \alpha^{-1}, \exp(U) \rightsquigarrow \gamma(\beta, 1)$
Log-Normal	$k_1 = m, k_2 = \sigma, \exp(U) \rightsquigarrow \text{LN}(0, 1)$

On peut également résumer l'espérance et la variance du terme d'erreur :

Modèle	$E(U)$	$V(U)$
Exponentiel	$-\gamma_E$	$\pi^2/6$
Weibull	$-\gamma_E$	$\pi^2/6$
Gamma	$\Gamma'(\beta)/\Gamma(\beta)$	$\Gamma''(\beta)/\Gamma(\beta) - (\Gamma'(\beta)/\Gamma(\beta))^2$
Gamma généralisé	$\Gamma'(\beta)/\Gamma(\beta)$	$\Gamma''(\beta)/\Gamma(\beta) - (\Gamma'(\beta)/\Gamma(\beta))^2$
Log-Normal	0	1

ainsi que de la perturbation complète du modèle :

Modèle	$E(k_2 U)$	$V(k_2 U)$
Exponentiel	$-\gamma_E$	$\pi^2/6$
Weibull	$-\gamma_E/\alpha$	$\pi^2/(6\alpha^2)$
Gamma	$\Gamma'(\beta)/\Gamma(\beta)$	$\Gamma''(\beta)/\Gamma(\beta) - (\Gamma'(\beta)/\Gamma(\beta))^2$
Gamma généralisé	$\alpha^{-1}\Gamma'(\beta)/\Gamma(\beta)$	$\alpha^{-2} \left\{ \Gamma''(\beta)/\Gamma(\beta) - (\Gamma'(\beta)/\Gamma(\beta))^2 \right\}$
Log-Normal	0	σ^2

9.5 Introduction des variables explicatives

9.5.1 Modèles à hasards proportionnels

Soit $X_i = (X_{1i}, \dots, X_{pi})$ un vecteur de p variables explicatives, on dit qu'un modèle est à hasard proportionnel s'il vérifie :

$$h_i(t) = h_0(t) \exp(X_i b),$$

où $h_0(t)$ est une fonction de hasard appelée hasard de base. On remarque qu'avec cette convention, le ratio des hasards de deux individus ne dépend que des variables explicatives, et non du temps :

$$\frac{h_i(t)}{h_j(t)} = \exp((X_i - X_j) b).$$

En prenant la fonction de hasard en logarithmes, on obtient :

$$\ln h_i(t) = \ln h_0(t) + X_i b,$$

de sorte que l'on peut écrire :

$$\frac{\partial \ln h_i(t)}{\partial X_{ki}} = b_k,$$

ceci implique que l'on peut interpréter b_k comme une élasticité quand la variable explicative k est en logarithmes. S'il s'agit d'une indicatrice, le coefficient b_k représente (s'il est proche de 0) l'écart de hasard en pourcentage ($100b_k$) entre la modalité 1 et la modalité 0 :

$$\begin{aligned} b_k &= \ln h_i(t|X_{ki} = 1) - \ln h_i(t|X_{ki} = 0) \\ &= \ln \left(1 + \frac{h_i(t|X_{ki} = 1) - h_i(t|X_{ki} = 0)}{h_i(t|X_{ki} = 0)} \right) \\ &\simeq \frac{h_i(t|X_{ki} = 1) - h_i(t|X_{ki} = 0)}{h_i(t|X_{ki} = 0)}. \end{aligned}$$

Pour procéder à l'estimation, on aura également besoin de la fonction de survie, à cause des données censurées :

$$\begin{aligned} S_i(t) &= \exp \left\{ - \int_0^t h(x) dx \right\} \\ &= \exp \left\{ - \exp(X_i b) \int_0^t h_0(x) dx \right\}, \end{aligned}$$

on en déduit la densité :

$$\begin{aligned} f_i(t) &= h_i(t) S_i(t) \\ &= h_0(t) \exp(X_i b) \exp \left\{ - \exp(X_i b) \int_0^t h_0(x) dx \right\}. \end{aligned}$$

Les modèles exponentiel, de Weibull et Gamma généralisés sont des modèles à hasard proportionnels. Toutefois dans le dernier cas, on ne peut pas écrire explicitement les fonction de survie et de hasard. Il faut recourir à une intégration numérique.

9.5.2 Le modèle exponentiel

Dans le cas du modèle exponentiel, on a :

$$h_0(t) = h,$$

de sorte que :

$$h_i(t) = h \exp(X_i b),$$

et

$$\int_0^t h_0(x) dx = h \int_0^t dx = ht,$$

de sorte que :

$$\begin{aligned} S_i(t) &= \exp\{-\exp(X_i b) ht\}, \\ f_i(t) &= h \exp(X_i b) \exp\{-\exp(X_i b) ht\}. \end{aligned}$$

Dans le cas du modèle de Weibull :

$$h_0(t) = h\alpha t^{\alpha-1},$$

d'où la fonction de hasard individuelle :

$$h_i(t) = h\alpha t^{\alpha-1} \exp(X_i b),$$

et

$$\begin{aligned} \int_0^t h_0(x) dx &= h \int_0^t \alpha x^{\alpha-1} dx \\ &= h [x^\alpha]_0^t \\ &= ht^\alpha, \end{aligned}$$

ce qui implique la fonction de survie :

$$S_i(t) = \exp\{-\exp(X_i b) ht^\alpha\},$$

et la densité :

$$f_i(t) = h\alpha t^{\alpha-1} \exp(X_i b) \exp\{-\exp(X_i b) ht^\alpha\}.$$

Pour la distribution Gamma, il faut évaluer numériquement les fonction de hasard $h_0(t)$ et de hasard cumulé

$$\Lambda_0(t) = \int_0^t h_0(x) dt,$$

le hasard individuel est donné par la formule habituelle alors que la fonction de survie et la densité sont égales à :

$$\begin{aligned} S_i(t) &= \exp\{-\exp(X_i b) \Lambda_0(t)\}, \\ f_i(t) &= h_0(t) \exp(X_i b) \exp\{-\exp(X_i b) \Lambda_0(t)\}. \end{aligned}$$

9.6 Ecriture de la vraisemblance

Pour écrire la log vraisemblance, on définit les variables suivantes :

- y_i est la variable de durée observable. Cette durée peut être aussi bien complète, c'est-à-dire observée jusqu'à son terme, que censurée, c'est-à-dire observée partiellement.
- $d_i \in \{0, 1\}$ est une indicatrice de censure. On observe $d_i = 1$ si l'observation i est censurée à droite et $d_i = 0$ si la durée est complète.
- Si la durée n'est pas censurée la vraisemblance de l'individu i est égale à $f_i(y_i)$, sinon elle est égale à $S_i(y_i)$.
- La log vraisemblance avec censure à droite s'écrit toujours :

$$\ell(y|X, \theta) = \sum_{i=1}^N (1 - d_i) \ln f_i(y_i) + d_i \ln S_i(y_i).$$

- En utilisant $f_i(y_i) = h_i(y_i) S_i(y_i)$ on obtient :

$$\begin{aligned} \ell(y|X, \theta) &= \sum_{i=1}^N (1 - d_i) (\ln h_i(y_i) + \ln S_i(y_i)) + d_i \ln S_i(y_i) \\ &= \sum_{i=1}^N (1 - d_i) \ln h_i(y_i) + \ln S_i(y_i) \end{aligned}$$

9.6.1 Modèle exponentiel

Les quantités dont on a besoin sont égales à :

$$\ln h_i(y_i) = \ln h + X_i b,$$

et

$$\begin{aligned}\ln S_i(y_i) &= -\exp(X_i b) h y_i \\ &= -\exp(\ln h + X_i b) y_i,\end{aligned}$$

on voit que $\ln h$ est le terme constant du modèle de sorte qu'il ne faut pas en mettre dans la liste des variables explicatives. On peut également faire un changement de paramètres :

$$Z_i = (1, X_i) \text{ et } \beta = \begin{pmatrix} \ln h \\ b \end{pmatrix}$$

de sorte que :

$$Z_i \beta = \ln h + X_i b,$$

la log-vraisemblance se réécrit donc :

$$\ell(y|X, \beta) = \sum_{i=1}^N (1 - d_i) (Z_i \beta) - \exp(Z_i \beta) y_i,$$

d'où le vecteur du score :

$$\frac{\partial \ell}{\partial \beta} (y|X, \beta) = \sum_{i=1}^N Z_i' (1 - d_i - \exp(Z_i \beta) y_i),$$

et le hessien

$$\frac{\partial^2 \ell}{\partial \beta \partial \beta'} (y|X, \beta) = -\sum_{i=1}^N Z_i' Z_i \exp(Z_i \beta) y_i \ll 0.$$

Une fois l'estimation réalisée, on obtient l'estimateur du maximum de vraisemblance de h par :

$$\hat{h} = \exp(\hat{\beta}_1) = g(\hat{\beta}_1),$$

et on estime sa variance asymptotique par :

$$\begin{aligned}\widehat{\text{Vas}}(\hat{h}) &= \frac{\partial g}{\partial \beta_1}(\hat{\beta}_1) \widehat{\text{Vas}}(\hat{\beta}_1) \frac{\partial g}{\partial \beta_1'}(\hat{\beta}_1) \\ &= \widehat{\text{Vas}}(\hat{\beta}_1) \exp(2\hat{\beta}_1).\end{aligned}$$

9.6.2 Modèle de Weibull

Les quantités dont on a besoin sont égales à :

$$\ln h_i(y_i) = \ln h + X_i b + \ln \alpha + (\alpha - 1) \ln y_i,$$

et

$$\begin{aligned} \ln S_i(y_i) &= -\exp(X_i b) h y_i^\alpha \\ &= -\exp(\ln h + X_i b) y_i^\alpha, \end{aligned}$$

on voit que $\ln h$ est le terme constant du modèle de sorte qu'il ne faut pas en mettre dans la liste des variables explicatives. On peut également faire un changement de paramètres similaire à celui du modèle exponentiel :

$$Z_i = (1, X_i), \quad \beta = \begin{pmatrix} \ln h \\ b \end{pmatrix} \quad \text{et} \quad \theta = \begin{pmatrix} \alpha \\ \beta \end{pmatrix},$$

de sorte que :

$$\ell(y|X, \theta) = \sum_{i=1}^N (1 - d_i) (Z_i \beta + \ln \alpha + (\alpha - 1) \ln y_i) - \exp(Z_i \beta) y_i^\alpha,$$

pour calculer le vecteur du score, on remarque que :

$$\frac{dy_i^\alpha}{d\alpha} = \frac{d}{d\alpha} (e^{\alpha \ln y_i}) = y_i^\alpha \ln y_i$$

d'où le vecteur du score :

$$\begin{aligned} \frac{\partial \ell}{\partial \alpha} (y|X, \theta) &= \sum_{i=1}^N (1 - d_i) \left(\frac{1}{\alpha} + \ln y_i \right) - \exp(Z_i \beta) y_i^\alpha \ln y_i \\ \frac{\partial \ell}{\partial \beta} (y|X, \theta) &= \sum_{i=1}^N Z_i' (1 - d_i - \exp(Z_i \beta) y_i^\alpha) \end{aligned}$$

et le hessien

$$\begin{aligned} \frac{\partial^2 \ell}{\partial \alpha^2} (y|X, \theta) &= - \sum_{i=1}^N \left\{ \frac{1 - d_i}{\alpha^2} + \exp(Z_i \beta) y_i^\alpha (\ln y_i)^2 \right\} \\ \frac{\partial^2 \ell}{\partial \beta \partial \alpha} (y|X, \theta) &= - \sum_{i=1}^N Z_i' \exp(Z_i \beta) y_i^\alpha \ln y_i \\ \frac{\partial^2 \ell}{\partial \beta \partial \beta'} (y|X, \theta) &= - \sum_{i=1}^N Z_i' Z_i \exp(Z_i \beta) y_i^\alpha \end{aligned}$$

Une fois l'estimation réalisée, on obtient l'estimateur du maximum de vraisemblance de h par :

$$\widehat{h} = \exp(\widehat{\theta}_1),$$

et on estime sa variance asymptotique comme dans le modèle exponentiel par :

$$\widehat{\text{Vas}}(\widehat{h}) = \widehat{\text{Vas}}(\widehat{\theta}_1) \exp(2\widehat{\theta}_1).$$

9.6.3 Modèle log-normal

Il ne s'agit pas d'un modèle à hasard proportionnel. Avec un modèle log normal, on fait directement une hypothèse sur la durée elle-même puisque l'on pose que :

$$\ln T = m + \sigma U,$$

où U suit une loi normale centrée et réduite. On peut donc voir ce modèle comme une simple extension du modèle linéaire standard normal. S'il n'y avait pas de censure des données, la méthode d'estimation adaptée serait simplement celle des moindres carrés ordinaires appliqués au logarithme de la durée. Toutefois, comme nous supposons la présence d'une censure droite, on ne peut pas appliquer les moindres carrés ordinaires. Il faut recourir à la méthode du maximum de vraisemblance. Une manière naturelle d'introduire des variables explicatives dans ce type de modèle consiste à poser simplement $m = Xb$.

La densité est alors donnée directement par :

$$f_i(y_i) = \frac{1}{\sigma y_i} \varphi\left(\frac{\ln y_i - X_i b}{\sigma}\right),$$

et la fonction de survie par :

$$S_i(y_i) = 1 - \Phi\left(\frac{\ln y_i - X_i b}{\sigma}\right),$$

ce qui donne la log-vraisemblance de l'échantillon :

$$\begin{aligned} \ell(y|X, \theta) &= \sum_{i=1}^N (1 - d_i) \ln f_i(y_i) + d_i \ln S_i(y_i) \\ &= \sum_{i=1}^N (1 - d_i) \left[\ln \varphi\left(\frac{\ln y_i - X_i b}{\sigma}\right) - \ln \sigma - \ln y_i \right] \\ &\quad + d_i \ln \left[1 - \Phi\left(\frac{\ln y_i - X_i b}{\sigma}\right) \right], \end{aligned}$$

on peut simplifier l'écriture du modèle en faisant le changement de paramètres suivant :

$$\beta = \frac{b}{\sigma}, \quad \gamma = \frac{1}{\sigma} \quad \text{et} \quad \theta = \begin{pmatrix} \beta \\ \gamma \end{pmatrix}$$

ce qui donne :

$$\begin{aligned} \ell(y|X, \theta) = \sum_{i=1}^N (1 - d_i) [\ln \varphi(\gamma \ln y_i - X_i \beta) + \ln \gamma - \ln y_i] \\ + d_i \ln [1 - \Phi(\gamma \ln y_i - X_i \beta)]. \end{aligned}$$

Pour alléger les notations, on pose :

$$u_i = \gamma \ln y_i - X_i \beta.$$

La première partie du vecteur du score est égal à :

$$\frac{\partial \ell}{\partial \beta}(y|X, \theta) = \sum_{i=1}^N X_i' \left\{ (1 - d_i) \left[-\frac{\varphi'(u_i)}{\varphi(u_i)} \right] - d_i \frac{\varphi(u_i)}{1 - \Phi(u_i)} \right\},$$

à ce stade on utilise $\varphi'(u) = -u\varphi(u)$, ce qui permet de simplifier l'expression précédente :

$$\frac{\partial \ell}{\partial \beta}(y|X, \theta) = \sum_{i=1}^N X_i' \left\{ (1 - d_i) u_i - d_i \frac{\varphi(u_i)}{1 - \Phi(u_i)} \right\},$$

si $c_i = 0 \forall i$, on retrouve les moindres carrés ordinaires, sinon on ajoute un terme pour corriger la censure droite. Pour l'autre paramètre, on trouve :

$$\frac{\partial \ell}{\partial \gamma}(y|X, \theta) = \sum_{i=1}^N \left\{ \ln y_i \left((1 - d_i) u_i - d_i \frac{\varphi(u_i)}{1 - \Phi(u_i)} \right) + \frac{1 - d_i}{\gamma} \right\}.$$

9.6.4 Généralisation

Les variables de durée peuvent également être censurées à gauche. On remarque qu'une même durée peut être censurée à la fois à gauche et à droite. Le fait d'avoir des censures à gauche ne change toutefois rien à notre analyse. En effet, si l'on observe une durée censurée y_i , on sait juste que la vraie durée est supérieure ou égale à y_i , et ce qu'elle soit censurée à gauche, à droite ou des deux côtés. Dans ce cas la contribution à la vraisemblance reste égale à $S_i(y_i)$. Il faut juste penser à définir une indicatrice de censure égale au maximum des deux indicatrices de censure

gauche et droite. Soit d_{1i} une indicatrice de censure gauche et d_{2i} une indicatrice de censure droite, on doit prendre :

$$d_i = \max(d_{1i}, d_{2i}),$$

dans les expressions de la section précédente. Cette règle reste valable si la censure a lieu avec des "trous" différents d'une observation à l'autre, car dans tous les cas la seule information disponible est que la vraie durée est supérieure à la durée observée y_i .

CHAPITRE 10

Les variables tronquées

10.1 Le modèle tronqué

On dit qu'un modèle est tronqué lorsque les variables explicatives X_i ne sont pas observables lorsque la variable expliquée z_i^* passe en dessous d'un certain seuil C_i . Ce cas peut se produire soit lorsque l'on n'interroge que les individus pour lesquels $z_i^* > C_i$ soit lorsque les réponses aux variables explicatives X_i n'ont de sens que lorsque $z_i^* > C_i$. Pour simplifier l'écriture du modèle, on pose :

$$y_i^* = z_i^* - C_i,$$

quantité qui peut toujours être calculée lorsque les seuils C_i sont connus. Avec ce changement de variable, on observe y_i^* lorsque $y_i^* > 0$ ($\Leftrightarrow z_i^* > C_i$). La variable latente est décrite par le modèle linéaire suivant :

$$y_i^* = X_i b + \sigma u_i, \quad u_i \stackrel{\text{iid}}{\rightsquigarrow} N(0, 1), \quad i = 1, \dots, N$$

avec $\sigma > 0$. La perturbation du modèle est donc égale à :

$$v_i = \sigma u_i,$$

ce qui implique que $v_i \stackrel{\text{iid}}{\rightsquigarrow} N(0, \sigma^2)$. La fonction de répartition de la loi normale centrée et réduite est notée $\Phi(z)$ et sa densité $\varphi(z)$. La variable observable, notée y_i , est définie par :

$$y_i = \begin{cases} \text{manquant} & \text{si } y_i^* \leq 0 \\ y_i^* & \text{sinon} \end{cases}$$

Pour procéder à l'estimation il nous faut l'expression de la densité de y_i^* tronquée en 0. Elle est égale, par définition, à :

$$f(y_i) = 1_{(y_i^* > 0)} \frac{f(y_i)}{\Pr[y_i^* > 0]}.$$

La probabilité d'observer la variable endogène est donnée par :

$$\begin{aligned}\Pr [y_i^* > 0] &= \Pr [X_i b + \sigma u_i > 0] \\ &= \Pr \left[u_i > -\frac{X_i b}{\sigma} \right] \\ &= 1 - \Phi \left(-\frac{X_i b}{\sigma} \right) \\ &= \Phi \left(\frac{X_i b}{\sigma} \right),\end{aligned}$$

et la vraisemblance est donnée par :

$$\ell_i = \sum_{y_i > 0} \ln \left[\frac{1}{\sigma} \varphi \left(\frac{y_i - X_i b}{\sigma} \right) \right] - \ln \Phi \left(\frac{X_i b}{\sigma} \right).$$

Il faut noter ici que seules les observations strictement positives de y_i^* sont utilisables, contrairement au modèle Tobit que nous verrons plus loin où toutes les observations sont utilisables. L'espérance mathématique de y_i est donnée par :

$$E(y_i) = \int_0^{+\infty} y_i \frac{\frac{1}{\sigma} \varphi \left(\frac{y_i - X_i b}{\sigma} \right)}{\Phi \left(\frac{X_i b}{\sigma} \right)} dy_i,$$

on effectue le changement de variable :

$$u = \frac{y_i - X_i b}{\sigma},$$

ce qui implique :

$$\lim_{y_i \rightarrow 0} u = -\frac{X_i b}{\sigma}, \quad \lim_{y_i \rightarrow +\infty} u = +\infty \quad \text{et} \quad dy_i = \sigma du,$$

d'où :

$$\begin{aligned}E[y_i | y_i > 0] &= \frac{1}{\Phi \left(\frac{X_i b}{\sigma} \right)} \int_{-\frac{X_i b}{\sigma}}^{+\infty} (X_i b + \sigma u) \varphi(u) du \\ &= \frac{1}{\Phi \left(\frac{X_i b}{\sigma} \right)} X_i b \underbrace{\int_{-\frac{X_i b}{\sigma}}^{+\infty} \varphi(u) du}_{\Phi \left(\frac{X_i b}{\sigma} \right)} + \frac{\sigma}{\Phi \left(\frac{X_i b}{\sigma} \right)} \int_{-\frac{X_i b}{\sigma}}^{+\infty} \underbrace{u \varphi(u) du}_{-\varphi'(u)} \\ &= X_i b + \frac{\sigma}{\Phi \left(\frac{X_i b}{\sigma} \right)} [-\varphi(u)]_{-\frac{X_i b}{\sigma}}^{+\infty} \\ &= X_i b + \sigma \frac{\varphi \left(\frac{X_i b}{\sigma} \right)}{\Phi \left(\frac{X_i b}{\sigma} \right)},\end{aligned}$$

cette espérance est valable sur les observations strictement positives et pourra être utilisée lors de l'estimation du modèle Tobit. On effectue les changements de paramètres suivant :

$$\beta = \frac{b}{\sigma} \quad \text{et} \quad h = \frac{1}{\sigma},$$

ce qui permet d'écrire la log-vraisemblance sous la forme :

$$\ell(y_1, \dots, y_N) = \sum_{i=1}^N \left\{ \ln h - \frac{1}{2} \ln(2\pi) - \frac{1}{2} (hy_i - X_i\beta)^2 - \ln \Phi(X_i\beta) \right\},$$

le score est donc égal à :

$$\begin{aligned} \frac{\partial \ell}{\partial \beta} &= \sum_{i=1}^N X_i' \left[hy_i - X_i\beta - \frac{\varphi(X_i\beta)}{\Phi(X_i\beta)} \right] \\ \frac{\partial \ell}{\partial h} &= \sum_{i=1}^N \left[\frac{1}{h} - hy_i^2 + X_i\beta y_i \right] \end{aligned}$$

Pour simplifier les notations, on pose :

$$m_i = X_i b \quad \text{et} \quad \lambda_i = \frac{\varphi(X_i\beta)}{\Phi(X_i\beta)},$$

où λ_i est l'inverse du ratio de Mills. La première partie du score peut donc se réécrire :

$$\frac{\partial \ell}{\partial \beta} = \sum_{i=1}^N X_i' [h(y_i - m_i) - \lambda_i]$$

On voit que le score $\partial \ell / \partial \beta$ est d'espérance nulle, puisque :

$$\mathbb{E}(y_i | y_i > 0) = m_i + \frac{1}{h} \lambda_i.$$

C'est également le cas pour $\partial \ell / \partial h$ et l'on peut donc écrire :

$$\begin{aligned} \mathbb{E} \left[\frac{1}{h} - hy_i^2 + hm_i y_i \right] &= 0 \\ \Leftrightarrow \mathbb{E} [y_i^2 | y_i > 0] &= \frac{1}{h} \left[\frac{1}{h} + hm_i \mathbb{E}(y_i | y_i > 0) \right] \\ \Leftrightarrow \mathbb{E} [y_i^2 | y_i > 0] &= \frac{1}{h^2} + m_i \mathbb{E}(y_i | y_i > 0) \\ \Leftrightarrow \mathbb{E} [y_i^2 | y_i > 0] &= \frac{1}{h^2} + m_i^2 + \frac{m_i \lambda_i}{h}. \end{aligned}$$

Cette expression nous servira pour déterminer l'algorithme du score. Les dérivées secondes sont égales à :

$$\begin{aligned}\frac{\partial^2 \ell}{\partial \beta \partial \beta'} &= - \sum_{i=1}^N X_i' X_i [1 - \lambda_i (m_i + \lambda_i)] \\ \frac{\partial^2 \ell}{\partial \beta \partial h} &= \sum_{i=1}^N X_i' y_i \\ \frac{\partial^2 \ell}{\partial h^2} &= - \sum_{i=1}^N \left[\frac{1}{h^2} + y_i^2 \right]\end{aligned}$$

d'où l'espérance de l'opposé des dérivées secondes :

$$\begin{aligned}\mathbb{E} \left[- \frac{\partial^2 \ell}{\partial \beta \partial \beta'} \right] &= \sum_{i=1}^N X_i' X_i [1 - \lambda_i (m_i + \lambda_i)] \\ \mathbb{E} \left[- \frac{\partial^2 \ell}{\partial \beta \partial h} \right] &= - \sum_{i=1}^N X_i' \left[m_i + \frac{1}{h} \lambda_i \right] \\ \mathbb{E} \left[- \frac{\partial^2 \ell}{\partial h^2} \right] &= \sum_{i=1}^N \left[\frac{2}{h^2} + m_i^2 + \frac{m_i \lambda_i}{h} \right]\end{aligned}$$

On peut alors utiliser un algorithme du score en prenant les moindres carrés ordinaires comme valeur initiale, cet estimateur n'étant pas convergent. L'algorithme de Newton-Raphson est ici plus simple que l'algorithme du score en raison de la forme particulière des espérances de la variable tronquée y_i .

10.2 Le modèle Tobit

Le modèle Tobit est un modèle censuré, ce qui signifie que l'on observe les variables explicatives X_i dans tous les cas. On peut donc utiliser cette information supplémentaire.

10.2.1 Estimation

La probabilité que la variable latente y_i^* soit négative est donnée par :

$$\Pr [y_i = 0] = \Pr [y_i^* < 0] = 1 - \Phi \left(\frac{X_i b}{\sigma} \right),$$

et la probabilité d'observer une valeur strictement positive est simplement égale à :

$$f(y_i) = \frac{1}{\sigma} \varphi \left(\frac{y_i - X_i b}{\sigma} \right), \quad y_i > 0,$$

la vraisemblance est donc égale à :

$$\ell(y_1, \dots, y_N) = (1 - d_i) \ln [1 - \Phi(X_i\beta)] + d_i \ln \left\{ \frac{h}{\sqrt{2\pi}} - \frac{1}{2} (h y_i - X_i\beta)^2 \right\},$$

avec :

$$d_i = \begin{cases} 0 & \text{si } y_i^* \leq 0 \\ 1 & \text{sinon} \end{cases}$$

où la variable dichotomique d_i suit une loi de Bernoulli de paramètre $\Phi(X_i\beta)$. Le score est égal à :

$$\frac{\partial \ell}{\partial \beta} = \sum_{i=1}^N X_i' \left[d_i (h y_i - X_i\beta) - (1 - d_i) \frac{\varphi(X_i\beta)}{1 - \Phi(X_i\beta)} \right]$$

$$\frac{\partial \ell}{\partial h} = \sum_{i=1}^N d_i \left[\frac{1}{h} - h y_i^2 + h m_i y_i \right]$$

et l'on vérifie que l'espérance du score est nulle en utilisant :

$$\begin{aligned} \mathbf{E}(d_i y_i) &= \mathbf{E}(y_i | y_i > 0) \times \Pr(y_i > 0) \\ &= \left(m_i + \frac{1}{h} \frac{\varphi_i}{\Phi_i} \right) \Phi_i \\ &= m_i \Phi_i + \frac{1}{h} \varphi_i, \end{aligned}$$

ainsi que :

$$\begin{aligned} \mathbf{E}(d_i y_i^2) &= \mathbf{E}(y_i^2 | y_i > 0) \times \Pr(y_i > 0) \\ &= \left[\frac{1}{h^2} + m_i^2 + \frac{m_i}{h} \frac{\varphi_i}{\Phi_i} \right] \Phi_i \\ &= \Phi_i \left(\frac{1}{h^2} + m_i^2 \right) + \frac{m_i \varphi_i}{h} \end{aligned}$$

Les dérivées secondes sont données par :

$$\frac{\partial^2 \ell}{\partial \beta \partial \beta'} = - \sum_{i=1}^N X_i' X_i \left\{ d_i + (1 - d_i) \frac{\varphi_i}{1 - \Phi_i} \left(\frac{\varphi_i}{1 - \Phi_i} - h m_i \right) \right\}$$

$$\frac{\partial^2 \ell}{\partial \beta \partial h} = \sum_{i=1}^N d_i X_i' y_i$$

$$\frac{\partial^2 \ell}{\partial h^2} = - \sum_{i=1}^N d_i \left(\frac{1}{h^2} + y_i^2 \right)$$

d'où les espérances mathématiques nécessaires à l'algorithme du score :

$$\begin{aligned} E \left[-\frac{\partial^2 \ell}{\partial \beta \partial \beta'} \right] &= \sum_{i=1}^N X_i' X_i \left[\Phi_i + \varphi_i \left(\frac{\varphi_i}{1 - \Phi_i} - h m_i \right) \right] \\ E \left[-\frac{\partial^2 \ell}{\partial \beta \partial h} \right] &= - \sum_{i=1}^N X_i' \left[m_i \Phi_i + \frac{1}{h} \varphi_i \right] \\ E \left[-\frac{\partial^2 \ell}{\partial h^2} \right] &= \sum_{i=1}^N \left[\Phi_i \left(\frac{2}{h^2} + m_i^2 \right) + \frac{m_i \varphi_i}{h} \right] \end{aligned}$$

10.2.2 Valeur initiale

Le fait que l'on observe toujours les variables explicatives permet de calculer facilement des valeurs initiales. Ceci provient du fait que l'on peut estimer un modèle Probit en prenant d_i comme variable expliquée. Pour les observations strictement positives, on utilise :

$$\begin{aligned} E(y_i | y_i > 0) &= X_i b + \sigma \frac{\varphi(X_i \beta)}{\Phi(X_i \beta)} \\ &= \sigma \left(X_i \beta + \frac{\varphi(X_i \beta)}{\Phi(X_i \beta)} \right) \\ &= \sigma \mu_i, \end{aligned} \tag{10.1}$$

avec $\mu_i = X_i \beta + \varphi(X_i \beta) / \Phi(X_i \beta)$. Cette quantité peut facilement être estimée à partir de l'estimateur $\hat{\beta}$ de la partie Probit du modèle :

$$\hat{\mu}_i = X_i \hat{\beta} + \frac{\varphi(X_i \hat{\beta})}{\Phi(X_i \hat{\beta})}.$$

En utilisant directement (10.1) on peut obtenir un estimateur convergent de σ en régressant y_i sur un estimateur convergent de μ_i par les moindres carrés ordinaires sans terme constant; ce qui donne :

$$\hat{\sigma} = \frac{\sum_{y_i > 0} \hat{\mu}_i y_i}{\sum_{y_i > 0} \hat{\mu}_i^2},$$

dont on déduit la valeur initiale convergente pour le paramètre h :

$$\hat{h} = \frac{1}{\hat{\sigma}} = \frac{\sum_{y_i > 0} \hat{\mu}_i^2}{\sum_{y_i > 0} \hat{\mu}_i y_i}.$$

10.2.3 Retour aux paramètres structurels

La méthode précédente permet d'obtenir des estimateurs convergents des paramètres β et h . Pour revenir aux paramètres de départ du modèle, on utilise la propriété d'invariance fonctionnelle et le théorème de Slutsky. L'invariance fonctionnelle implique que $\hat{b} = \hat{\beta}/\hat{h}$ est l'estimateur du maximum de vraisemblance de b et que $\hat{\sigma} = 1/\hat{h}$ est l'estimateur du maximum de vraisemblance de σ . Le théorème de Slutsky permet de trouver la matrice de covariance asymptotique de $(\hat{b}', \hat{\sigma})'$. Soit :

$$\theta = \begin{pmatrix} \beta \\ h \end{pmatrix},$$

la distribution asymptotique de l'estimateur du maximum de vraisemblance est normale :

$$\sqrt{N}(\hat{\theta} - \theta) \xrightarrow{L} N(0, \Omega_{\hat{\theta}}),$$

ce qui implique :

$$\sqrt{N}(g(\hat{\theta}) - g(\theta)) \xrightarrow{L} N\left(0, \frac{\partial g}{\partial \theta'} \Omega_{\hat{\theta}} \frac{\partial g}{\partial \theta'}\right),$$

avec

$$g(\theta) = (\beta/h, 1/h) = (b, \sigma) \quad \text{et} \quad \frac{\partial g}{\partial \theta}(\theta) = \begin{pmatrix} \mathbf{I}_k/h & 0 \\ -\beta/h^2 & -1/h^2 \end{pmatrix},$$

où k est le nombre de paramètres du vecteur β . On en déduit que :

$$\begin{pmatrix} \hat{b} \\ \hat{\sigma} \end{pmatrix} = \begin{pmatrix} \hat{\beta}/\hat{h} \\ 1/\hat{h} \end{pmatrix}$$

et $\hat{V} \begin{pmatrix} \hat{b} \\ \hat{\sigma} \end{pmatrix} = \begin{pmatrix} \mathbf{I}_k/\hat{h} & 0 \\ -\hat{\beta}/\hat{h}^2 & -1/\hat{h}^2 \end{pmatrix} \hat{V} \begin{pmatrix} \hat{\beta} \\ \hat{h} \end{pmatrix} \begin{pmatrix} \mathbf{I}_k/\hat{h} & -\hat{\beta}/\hat{h}^2 \\ 0 & -1/\hat{h}^2 \end{pmatrix}$

10.3 Le modèle Tobit généralisé

On présente ici le modèle développé par Heckman dans ses articles de 1976 et 1979.

10.3.1 Définition

La forme latente comporte maintenant deux équations. Une première variable latente y_1^* détermine la décision et une seconde variable latente

y_2^* détermine le montant observé quand la décision est prise. On a :

$$\begin{aligned} y_{1i}^* &= m_{1i} + \sigma_1 u_{1i} \\ y_{2i}^* &= m_{2i} + \sigma_2 u_{2i}, \end{aligned}$$

avec $m_{ji} = X_i b_j$ et :

$$\begin{pmatrix} u_{1i} \\ u_{2i} \end{pmatrix} \overset{\text{iid}}{\rightsquigarrow} \text{N} \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right].$$

On observe la décision :

$$y_{1i} = \begin{cases} 0 & \text{si } y_{1i}^* \leq 0 \\ 1 & \text{sinon} \end{cases},$$

ainsi que le montant lorsque $y_{1i} = 1$:

$$y_{2i} = \begin{cases} \text{manquant} & \text{si } y_{1i}^* \leq 0 \\ y_{2i}^* & \text{sinon} \end{cases}$$

10.3.2 Estimation

Pour écrire la log vraisemblance, on a besoin de la loi normale bivariée. Sa densité est égale à :

$$f(y_1^*, y_2^*) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} (u_1^2 - 2\rho u_1 u_2 + u_2^2) \right\} \quad (10.2)$$

Pour les observations nulles, la probabilité est simplement :

$$\Pr[y_{1i} = 0] = 1 - \Phi \left(\frac{m_{1i}}{\sigma_1} \right) = 1 - \Phi(X_{1i}\beta_1),$$

avec $\beta_1 = b_1/\sigma_1$. Pour les observations positives, il faut calculer :

$$\begin{aligned} f(y_2^* \cap (y_1^* > 0)) &= \int_{-X_{1i}\beta_1}^{+\infty} f(u_1, u_2) du_1 \\ &= \int_{-X_{1i}\beta_1}^{+\infty} f(u_1|u_2) \varphi(u_2) du_1 \\ &= \varphi(u_2) \int_{-X_{1i}\beta_1}^{+\infty} f(u_1|u_2) du_1. \end{aligned}$$

D'après la densité (10.2) on a :

$$\begin{aligned} f(u_1|u_2) &= \frac{\frac{1}{2\pi\sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} (u_1^2 - 2\rho u_1 u_2 + u_2^2) \right\}}{\frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{u_2^2}{2} \right\}} \\ &= \frac{1}{\sqrt{1-\rho^2}} \varphi \left(\frac{u_1 - \rho u_2}{\sqrt{1-\rho^2}} \right) \end{aligned}$$

ce qui permet d'écrire :

$$I = \int_{-X_{1i}\beta_1}^{+\infty} f(u_1|u_2) du_1 = \int_{-X_{1i}\beta_1}^{+\infty} \frac{1}{\sqrt{1-\rho^2}} \varphi\left(\frac{u_1 - \rho u_2}{\sqrt{1-\rho^2}}\right) du_1,$$

en effectuant le changement de variable :

$$z = \frac{u_1 - \rho u_2}{\sqrt{1-\rho^2}},$$

on a :

$$du_1 = \sqrt{1-\rho^2} dz, \quad \lim_{u_1 \rightarrow -X_{1i}\beta_1} z = \frac{-X_{1i}\beta_1 - \rho u_2}{\sqrt{1-\rho^2}}, \quad \lim_{u_1 \rightarrow +\infty} z = +\infty,$$

ce qui implique :

$$\begin{aligned} I &= \int_{\frac{-X_{1i}\beta_1 - \rho u_2}{\sqrt{1-\rho^2}}}^{+\infty} \varphi(z) dz \\ &= 1 - \int_{-\infty}^{\frac{-X_{1i}\beta_1 - \rho u_2}{\sqrt{1-\rho^2}}} \varphi(z) dz \\ &= 1 - \Phi\left(\frac{-X_{1i}\beta_1 - \rho u_2}{\sqrt{1-\rho^2}}\right) \\ &= \Phi\left(\frac{X_{1i}\beta_1 + \rho u_2}{\sqrt{1-\rho^2}}\right) \\ &= \Phi\left(\frac{X_{1i}\beta_1 + \rho(y_{2i} - X_{2i}\beta_2)/\sigma_2}{\sqrt{1-\rho^2}}\right), \end{aligned}$$

en posant :

$$h = \frac{1}{\sigma_2} \quad \text{et} \quad \beta_2 = \frac{b_2}{\sigma_2},$$

on obtient finalement :

$$I = \Phi\left(\frac{X_{1i}\beta_1 + \rho(h y_{2i} - X_{2i}\beta_2)}{\sqrt{1-\rho^2}}\right).$$

La log vraisemblance du modèle Tobit généralisé est donc égale à :

$$\begin{aligned} \ell(\beta_1, \beta_2, h) &= \sum_{y_{1i}=0} \ln(1 - \Phi(X_{1i}\beta_1)) \\ &\quad + \sum_{y_{1i}=1} \ln \Phi\left(\frac{X_{1i}\beta_1 + \rho(h y_{2i} - X_{2i}\beta_2)}{\sqrt{1-\rho^2}}\right) \\ &\quad + \sum_{y_{1i}=1} \ln \left\{ \frac{h}{\sqrt{2\pi}} \exp -\frac{1}{2} (h y_{2i} - X_{2i}\beta_2)^2 \right\} \end{aligned}$$

10.3.3 Valeur initiale

Comme dans le modèle Tobit simple, il est possible de trouver une valeur initiale à partir d'une méthode en deux étapes. Pour cela, on utilise l'espérance conditionnelle suivante :

$$\begin{aligned} E(y_{2i}^* | y_{1i}^* > 0) &= X_{2i}b_2 + \sigma_2 E(u_{2i} | u_{1i} > -X_{1i}\beta_1) \\ &= X_{2i}b_2 + \rho\sigma_2 \frac{\varphi(X_{1i}\beta_1)}{\Phi(X_{1i}\beta_1)}. \end{aligned}$$

Dans un premier temps, on estime donc un modèle Probit sur toutes les observations, ce qui permet d'obtenir un estimateur de β_1 noté $\hat{\beta}_1$. On estime ensuite l'inverse du ratio de Mills noté $\hat{\lambda}_i$:

$$\hat{\lambda}_i = \frac{\varphi(X_{1i}\hat{\beta}_1)}{\Phi(X_{1i}\hat{\beta}_1)},$$

en régressant les observations positives de y_2^* sur X_2 et \hat{M}_i on obtient un estimateur convergent de b_2 et de $\rho\sigma_2$. On peut ensuite soit estimer le modèle en faisant un balayage sur ρ soit utiliser une expression similaire sur la variance conditionnelle de y_2^* pour estimer ρ . Le lecteur intéressé est invité à se reporter à l'ouvrage de C. Gouriéroux.

10.3.4 Amélioration de l'estimation

La plupart des logiciels n'ont besoin que de la log-vraisemblance pour déterminer le maximum de la fonction précédente, surtout si elle prend pour valeur initiale l'estimateur en deux étapes de Heckman, parce qu'il est convergent. Toutefois, pour pouvoir estimer un système d'équation incluant une variable modélisée par un Tobit généralisé, il faut disposer des dérivées premières analytiques. Elles permettent de calculer la matrice de covariance de la forme réduite du modèle. Ceci permet également d'accélérer les procédures d'optimisation numériques. On pose les notations suivantes :

$$\mu_{1i} = X_{1i}\beta_1, \quad \mu_{2i} = X_{2i}\beta_{2i}$$

et

$$\lambda_{2i} = \varphi\left(\frac{\mu_{1i} + \rho(h y_{2i} - \mu_{2i})}{\sqrt{1 - \rho^2}}\right) \Phi\left(\frac{\mu_{1i} + \rho(h y_{2i} - \mu_{2i})}{\sqrt{1 - \rho^2}}\right)^{-1}.$$

La log vraisemblance pour une observation s'écrit :

$$\begin{aligned} \ell_i &= (1 - y_{1i}) \ln(1 - \Phi(\mu_{1i})) \\ &\quad + y_{1i} \ln \Phi\left(\frac{\mu_{1i} + \rho(h y_{2i} - \mu_{2i})}{\sqrt{1 - \rho^2}}\right) \\ &\quad + y_{1i} \left(\ln h - \frac{1}{2} \ln(2\pi) - \frac{1}{2} (h y_{2i} - \mu_{2i})^2\right). \end{aligned}$$

Les dérivées pour chaque observation s'écrivent donc :

$$\begin{aligned} \frac{\partial \ell_i}{\partial \mu_{1i}} &= \frac{y_{1i} \lambda_{2i}}{\sqrt{1 - \rho^2}} - \frac{(1 - y_{1i}) \varphi_{1i}}{1 - \Phi_{1i}}, \\ \frac{\partial \ell_i}{\partial \mu_{2i}} &= y_{1i} \left\{ h y_{2i} - \mu_{2i} - \frac{\rho y_{1i} \lambda_{2i}}{\sqrt{1 - \rho^2}} \right\}, \\ \frac{\partial \ell_i}{\partial h} &= y_{1i} \left\{ \frac{\rho y_{2i} \lambda_{2i}}{\sqrt{1 - \rho^2}} + \frac{1}{h} - y_{2i} (h y_{2i} - \mu_{2i}) \right\}, \\ \frac{\partial \ell_i}{\partial \rho} &= y_{1i} \lambda_{2i} (1 - \rho^2)^{-3/2} \{ \rho \mu_{1i} + h y_{2i} - \mu_{2i} \}. \end{aligned}$$

On en déduit les dérivées par rapport aux paramètres :

$$\frac{\partial \ell}{\partial \beta_1} = \sum_{i=1}^N X'_{1i} \frac{\partial \ell_i}{\partial \mu_{1i}}, \quad \frac{\partial \ell}{\partial \beta_2} = \sum_{i=1}^N X'_{2i} \frac{\partial \ell_i}{\partial \mu_{2i}}, \quad \frac{\partial \ell}{\partial h} = \sum_{i=1}^N \frac{\partial \ell_i}{\partial h} \quad \text{et} \quad \frac{\partial \ell}{\partial \rho} = \sum_{i=1}^N \frac{\partial \ell_i}{\partial \rho}.$$

10.3.5 Programmation

Pour procéder à l'optimisation de la log vraisemblance, il est pratique de procéder aux changements de paramètres suivants :

$$h = \exp(\gamma_1) > 0, \quad \rho = \sin(\gamma_2) \in [-1, +1],$$

on obtient alors les dérivées de la manière suivante :

$$\frac{\partial \ell_i}{\partial \gamma_1} = \frac{\partial \ell_i}{\partial h} \frac{\partial h}{\partial \gamma_1} = \frac{\partial \ell_i}{\partial h} \exp(\gamma_1),$$

et

$$\frac{\partial \ell_i}{\partial \gamma_2} = \frac{\partial \ell_i}{\partial \rho} \frac{\partial \rho}{\partial \gamma_2} = \frac{\partial \ell_i}{\partial \rho} \cos(\gamma_2).$$

Pour les valeurs initiales, on utilise la régression sur les données quantitatives observables :

$$\{\forall i | y_1 = 1\}, \quad y_{2i} = X_{2i} b_2 + c \times \hat{\lambda}_{1i} + v_{2i},$$

avec :

$$\hat{\lambda}_{1i} = \frac{\varphi\left(X_{1i}\hat{\beta}_1\right)}{\Phi\left(X_{1i}\hat{\beta}_1\right)}, \quad c = \rho \times \sigma_2.$$

A la suite de Gouriéroux (1989), on prend comme valeur initiale

$$\hat{\sigma}_2^2 = \frac{1}{N_1} \sum_{\forall i|y_i=1} \left(\hat{v}_{2i}^2 + \hat{c}^2 \hat{\lambda}_{1i} \left(X_{1i}\hat{\beta}_1 + \hat{\lambda}_{1i} \right) \right),$$

ce qui implique :

$$\hat{h} = \frac{1}{\sqrt{\hat{\sigma}_2^2}} \quad \text{donc} \quad \hat{\gamma}_1 = \ln \hat{h},$$

ainsi que :

$$\hat{\rho} = \frac{\hat{c}}{\hat{\sigma}_2} = \hat{h} \times \hat{c} \quad \text{donc} \quad \hat{\gamma}_2 = \sin^{-1}(\hat{\rho}).$$

CHAPITRE 11

Estimation de modèles à plusieurs équations

11.1 Estimation de la forme réduite

Pour fixer les idées, on cherche à estimer le système à deux équations suivant :

$$\begin{cases} y_1^* &= a_{12} y_2^* + X_1 b_1 + u_1 \\ y_2^* &= a_{21} y_1^* + X_2 b_2 + u_2 \end{cases} \quad (11.1)$$

en résolvant ce système par rapport aux variables expliquées (y_1^*, y_2^*) en fonction des variables explicatives et des perturbations, on obtient la forme réduite du modèle :

$$\begin{cases} y_1^* &= X \pi_1 + v_1 \\ y_2^* &= X \pi_2 + v_2 \end{cases}$$

où X est la matrice de toutes les variables explicatives et :

$$v_1 = \frac{u_1 + a_{12} u_2}{1 - a_{12} a_{21}}, \quad v_2 = \frac{a_{21} u_1 + u_2}{1 - a_{12} a_{21}}. \quad (11.2)$$

On voit que ce système peut être estimé très simplement, équation par équation, puisqu'il n'y a plus de variable qualitative endogène dans les membres de droite des équations de la forme réduite. Le seul problème consiste à estimer la matrice de covariance globale des ces estimateurs obtenus séparément.

Il est facile de voir que les estimateurs de la forme réduite $(\hat{\pi}_1, \hat{\pi}_2)$ peuvent être obtenus par la maximisation d'un objectif de la forme suivante :

$$\hat{\pi} = \arg \max_{(\pi_1, \pi_2)} \sum_{i=1}^N \Psi_1(\pi_1; y_{1i}, X_i) + \sum_{i=1}^N \Psi_2(\pi_2; y_{2i}, X_i) \quad (11.3)$$

En effet, la dérivée par rapport à π_1 ne fait intervenir que la première partie de l'objectif, qui réalise l'estimation par le (pseudo) maximum de vraisemblance, alors que la dérivée par rapport à π_2 ne fait intervenir que la dérivée par rapport à la seconde partie de l'objectif. En conséquence, les conditions du premier ordre de ce problème sont identiques à celles des estimations séparées ce qui implique que les estimateurs obtenus en maximisant l'objectif (11.3) sont numériquement identiques à ceux obtenus par les estimations séparées. Il nous reste à voir comment calculer la matrice de covariance de $\hat{\pi}$ pour résoudre notre problème. Le problème d'optimisation se réécrit :

$$\hat{\pi} = \arg \max_{\pi} \sum_{i=1}^N \Psi(\pi; y_i, X_i)$$

avec $y = (y_1, y_2)$ et $\Psi(\pi; y_i, X_i) = \Psi_1(\pi_1; y_{1i}, X_i) + \Psi_2(\pi_2; y_{2i}, X_i)$. L'estimateur de la forme réduite est défini par :

$$\sum_{i=1}^N \frac{\partial \Psi}{\partial \pi}(\hat{\pi}; y_i, X_i) = 0.$$

La matrice de covariance est donnée par le résultat suivant, qui s'applique aux M-estimateurs en général (Gouriéroux et Monfort, 1989) :

$$\sqrt{N}(\hat{\pi} - \pi) \xrightarrow{L} N(0, \Sigma),$$

avec

$$\begin{aligned} \Sigma &= J^{-1} I J^{-1}, \\ J &= E \left[-\frac{\partial^2 \Psi}{\partial \pi \partial \pi'}(\pi; y, X) \right], \\ I &= E \left[\frac{\partial \Psi}{\partial \pi}(\pi; y, X) \frac{\partial \Psi}{\partial \pi}(\pi; y, X) \right]. \end{aligned}$$

Dans la pratique, on estimera ces quantités par :

$$\begin{aligned} \hat{J} &= -\frac{1}{N} \sum_{i=1}^N \frac{\partial^2 \Psi}{\partial \pi \partial \pi'}(\hat{\pi}; y_i, X_i) \\ \text{et } \hat{I} &= \frac{1}{N} \sum_{i=1}^N \frac{\partial \Psi}{\partial \pi}(\hat{\pi}; y_i, X_i) \frac{\partial \Psi}{\partial \pi'}(\hat{\pi}; y_i, X_i). \end{aligned}$$

On en déduit les remarques importantes suivantes :

1. Les dérivées secondes croisées entre équations sont toutes nulles puisque le paramètre d'une équation de la forme réduite n'apparaît que dans cette équation.
2. De la première remarque, on déduit le résultat suivant :

$$\begin{aligned} \Omega_{\hat{\pi}} &= \begin{pmatrix} J_{11}^{-1} & 0 \\ 0 & J_{22}^{-1} \end{pmatrix} \begin{pmatrix} I_{11} & I_{12} \\ I'_{12} & I_{22} \end{pmatrix} \begin{pmatrix} J_{11}^{-1} & 0 \\ 0 & J_{22}^{-1} \end{pmatrix} \\ &= \begin{pmatrix} J_{11}^{-1} I_{11} J_{11}^{-1} & J_{11}^{-1} I_{12} J_{22}^{-1} \\ J_{22}^{-1} I'_{12} J_{11}^{-1} & J_{22}^{-1} I_{22} J_{22}^{-1} \end{pmatrix}. \end{aligned}$$

3. Les estimateurs des matrices de covariance asymptotiques de chaque équation pris séparément, qui sont situées sur la diagonale, sont identiques à ceux du pseudo maximum de vraisemblance.
4. La covariance asymptotique entre les estimateurs des deux équations est donnée par :

$$\text{Covas} \left(\sqrt{N} (\hat{\pi}_1 - \pi_1), \sqrt{N} (\hat{\pi}_2 - \pi_2) \right) = J_{11}^{-1} I_{12} J_{22}^{-1}.$$

5. Il faut donc sauvegarder les dérivées premières individu par individu pour pouvoir estimer la matrice de covariance de l'estimateur de la forme réduite. La seule nouveauté est donc la matrice I_{12} , que l'on estimera par :

$$\hat{I}_{12} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \Psi_1}{\partial \pi_1} (\hat{\pi}_1; y_{1i}, X_{1i}) \frac{\partial \Psi_2}{\partial \pi_2'} (\hat{\pi}_2; y_{2i}, X_{2i}),$$

et que l'on obtient par un simple produit croisé des matrices des dérivées premières individuelles.

11.2 Estimation de la forme structurelle

En reportant les formes réduites de (y_1^*, y_2^*) dans le système (11.1) on obtient les identités suivantes :

$$\begin{aligned} X\pi_1 + v_1 &= a_{12} (X\pi_2 + v_2) + X_1 b_1 + u_1 \\ X\pi_2 + v_2 &= a_{21} (X\pi_1 + v_1) + X_2 b_2 + u_2 \end{aligned}$$

en prenant l'espérance mathématique du système précédent, on obtient les égalités :

$$\begin{cases} X\pi_1 &= X\pi_2 a_{12} + X_1 b_1 \\ X\pi_2 &= X\pi_1 a_{21} + X_2 b_2 \end{cases} \quad (11.4)$$

Pour obtenir une relation entre les paramètres du modèle, on introduit les matrices d'exclusion E_1 et E_2 , définies de la manière suivante :

$$X_1 = X E_1, \quad X = X E_2.$$

Ces matrices résument les contraintes qui permettent d'identifier le modèle, c'est-à-dire de remonter de la forme réduite du modèle à sa forme structurelle. On obtient la propriété suivante :

$$\begin{cases} X(\pi_1 - \pi_2 a_{12} - E_1 b_1) = 0 \\ X(\pi_2 - \pi_1 a_{21} - E_2 b_2) = 0 \end{cases} \Rightarrow \begin{cases} \pi_1 - \pi_2 a_{12} - E_1 b_1 = 0 \\ \pi_2 - \pi_1 a_{21} - E_2 b_2 = 0 \end{cases} \quad (11.5)$$

car X est de plein rang colonne. Les relations (11.5) s'appellent les contraintes identifiantes. La méthode des moindres carrés asymptotiques permet d'estimer la forme structurelle du modèle à partir d'un estimateur convergent et asymptotiquement normal (CAN) de la forme réduite. On note cet estimateur :

$$\sqrt{N}(\hat{\pi} - \pi) \xrightarrow{L} N(0, \Omega_{\hat{\pi}}).$$

Les équations auxiliaires sont définies par :

$$\begin{cases} \hat{\pi}_1 = \hat{\pi}_2 a_{12} + E_1 b_1 + \omega_1 \\ \hat{\pi}_2 = \hat{\pi}_1 a_{21} + E_2 b_2 + \omega_2 \end{cases} \quad (11.6)$$

où $\omega = (\omega_1, \omega_2)'$ est un terme d'erreur qui vérifie $\text{Plim} \sqrt{N}\omega = 0$. Ce système peut être estimé en deux étapes. Une première étape sert à estimer la matrice de covariance asymptotique de ω ; la seconde étape sert à obtenir l'estimateur optimal. On estime d'abord la relation suivante par les moindres carrés ordinaires :

$$\underbrace{\begin{pmatrix} \hat{\pi}_1 \\ \hat{\pi}_2 \end{pmatrix}}_{\hat{\pi}} = \underbrace{\begin{pmatrix} \hat{\pi}_2 & E_1 & 0 & 0 \\ 0 & 0 & \hat{\pi}_1 & E_2 \end{pmatrix}}_{\hat{H}} \underbrace{\begin{pmatrix} a_{12} \\ b_1 \\ a_{21} \\ b_2 \end{pmatrix}}_{\beta} + \underbrace{\begin{pmatrix} \omega_1 \\ \omega_2 \end{pmatrix}}_{\omega},$$

on obtient l'estimateur suivant :

$$\hat{\beta} = \left(\hat{H}' \hat{H} \right)^{-1} \hat{H}' \hat{\pi}.$$

Cet estimateur de première étape permet de calculer la variance de ω de la manière suivante :

$$\omega = \underbrace{\left[\begin{pmatrix} 1 & -a_{12} \\ -a_{21} & 1 \end{pmatrix} \otimes I_k \right]}_A \hat{\pi} + \begin{pmatrix} E_1 b_1 \\ E_2 b_2 \end{pmatrix},$$

ce qui implique :

$$V(\omega) = A V(\hat{\pi}) A',$$

on obtient un estimateur de cette variance en remplaçant (a_{12}, a_{21}) par leurs estimations :

$$\hat{A} = \begin{pmatrix} 1 & -\hat{a}_{12} \\ -\hat{a}_{21} & 1 \end{pmatrix} \otimes I_k \quad \text{et} \quad \hat{V}(\omega) = \hat{A} V(\hat{\pi}) \hat{A}'.$$

L'estimateur optimal β^* est obtenu en appliquant les moindres carrés généralisés à la relation (11.6) :

$$\beta^* = \left(\hat{H}' \hat{V}(\omega)^{-1} \hat{H} \right)^{-1} \hat{H}' \hat{V}(\omega)^{-1} \hat{\pi}, \quad (11.7)$$

et sa matrice de covariance peut être estimée par :

$$V(\beta^*) = \left(\hat{H}' V^*(\omega)^{-1} \hat{H} \right)^{-1}$$

avec

$$V^*(\omega) = A^* V(\hat{\pi}) A^{*'} \quad \text{et} \quad A^* = \begin{pmatrix} 1 & -a_{12}^* \\ -a_{21}^* & 1 \end{pmatrix} \otimes I_k.$$

On peut effectuer une troisième itération en remplaçant $\hat{V}(\omega)$ par $V^*(\omega)$ dans la relation (11.7).

Annexe A

Moments empiriques et moments théoriques

A.1 Moments empiriques des vecteurs

Le but de cette section est de se familiariser avec les notations de calcul matriciel, car c'est sous cette forme qu'apparaissent le plus souvent les moments empiriques. Il faut donc savoir les simplifier quand on les rencontre dans une expression.

A.1.1 Moyenne arithmétique

La moyenne arithmétique d'un vecteur colonne $z = (z_1, z_2, \dots, z_N)'$ peut se trouver sous les formes équivalentes suivantes :

$$\bar{z} = \frac{z'e}{e'e} = \frac{z'e}{N} = \frac{1}{N} \sum_{i=1}^N z_i,$$

car on a :

$$z'e = (z_1, z_2, \dots, z_N) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = z_1 + z_2 + \dots + z_N = \sum_{i=1}^N z_i,$$

et :

$$e'e = (1, 1, \dots, 1) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \underbrace{1 + 1 + \dots + 1}_{N \text{ fois}} = N.$$

A.1.2 Variance empirique

La variance empirique de la série z , notée $V_e(z)$, peut se trouver sous les formes équivalentes :

$$\begin{aligned} V_e(z) &= \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z})^2 \\ &= \frac{1}{N} \sum_{i=1}^N z_i^2 - (\bar{z})^2 \\ &= \frac{1}{N} (z - \bar{z}e)' (z - \bar{z}e), \\ &= \frac{z'z}{N} - (\bar{z})^2 \end{aligned}$$

car

$$z - \bar{z}e = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_N \end{pmatrix} - \begin{pmatrix} \bar{z} \\ \bar{z} \\ \vdots \\ \bar{z} \end{pmatrix} = \begin{pmatrix} z_1 - \bar{z} \\ z_2 - \bar{z} \\ \vdots \\ z_N - \bar{z} \end{pmatrix},$$

ce qui implique :

$$\begin{aligned} (z - \bar{z}e)' (z - \bar{z}e) &= (z_1 - \bar{z}, z_2 - \bar{z}, \dots, z_N - \bar{z}) \begin{pmatrix} z_1 - \bar{z} \\ z_2 - \bar{z} \\ \vdots \\ z_N - \bar{z} \end{pmatrix} \\ &= (z_1 - \bar{z})^2 + (z_2 - \bar{z})^2 + \dots + (z_N - \bar{z})^2 \\ &= \sum_{i=1}^N (z_i - \bar{z})^2. \end{aligned}$$

En posant $\bar{z} = 0$, on trouve :

$$z'z = \sum_{i=1}^N z_i^2.$$

A.1.3 Ecart-type empirique

Il s'agit simplement de la racine carrée de la variance empirique. On le note :

$$\sigma_e(x) = \sqrt{V_e(x)}.$$

A.1.4 Covariance empirique

La covariance empirique entre le vecteur $z = (z_1, z_2, \dots, z_N)'$ et le vecteur $x = (x_1, x_2, \dots, x_N)'$, $\text{Cov}_e(z, x)$, s'écrit :

$$\begin{aligned}\text{Cov}_e(x, z) &= \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z})(x_i - \bar{x}) \\ &= \frac{1}{N} \sum_{i=1}^N z_i x_i - \bar{z} \bar{x} \\ &= \frac{1}{N} (z - \bar{z}e)' (x - \bar{x}e) \\ &= \frac{z'x}{N} - \bar{z} \bar{x}.\end{aligned}$$

En effet :

$$\begin{aligned}(z - \bar{z}e)' (x - \bar{x}e) &= (z_1 - \bar{z}, z_2 - \bar{z}, \dots, z_N - \bar{z}) \begin{pmatrix} x_1 - \bar{x} \\ x_2 - \bar{x} \\ \vdots \\ x_N - \bar{x} \end{pmatrix} \\ &= (z_1 - \bar{z})(x_1 - \bar{x}) + \dots + (z_N - \bar{z})(x_N - \bar{x}) \\ &= \sum_{i=1}^N (z_i - \bar{z})(x_i - \bar{x}).\end{aligned}$$

En posant $\bar{z} = 0 = \bar{x}$ dans l'expression précédente, on a :

$$z'x = \sum_{i=1}^N z_i x_i.$$

On remarque de plus que lorsque $z = x$:

$$\begin{aligned}\text{Cov}_e(x, x) &= \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x}) \\ &= \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \\ &= V_e(x).\end{aligned}$$

A.1.5 Corrélation empirique

Le coefficient de corrélation linéaire empirique entre les séries z et x , noté $\rho_e(x, z)$ est défini par :

$$\rho_e(x, z) = \frac{\text{Cov}_e(x, z)}{\sqrt{V_e(x) V_e(z)}} = \frac{\text{Cov}_e(x, z)}{\sigma_e(x) \sigma_e(z)}.$$

Il peut donc prendre différentes formes en fonction des expressions que nous avons vu plus haut. On peut faire apparaître son expression dans la définition des différents estimateurs.

A.2 Moments empiriques des matrices

A.2.1 Moyenne arithmétique

On considère maintenant une matrice X de dimension (N, p) . Chaque ligne de X correspond à une observation et chaque colonne de X correspond à une variable. On note ces variables $X = (X^{(1)} | X^{(2)} | \dots | X^{(p)})$. On a :

$$\bar{X} = \underbrace{\frac{X'e}{N}}_{(p,1)} = \frac{1}{N} \begin{pmatrix} X^{(1)'} \\ X^{(2)'} \\ \vdots \\ X^{(p)'} \end{pmatrix} e = \frac{1}{N} \begin{pmatrix} X^{(1)'}e \\ X^{(2)'}e \\ \vdots \\ X^{(p)'}e \end{pmatrix} = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_p \end{pmatrix}.$$

A.2.2 Matrice de covariance empirique

Contrairement au cas univarié, on définit une matrice qui contient à la fois les variances et les covariances des variables. Les variances sont sur la diagonale de la matrice de covariance. On a :

$$V_e(X) = \frac{X'X}{N} - \bar{X}\bar{X}'$$

On peut définir la matrice des produits croisés des variables explicatives $X'X$ à partir du modèle écrit par observations ou par variables. Selon le contexte une expression peut s'avérer plus pratique que l'autre, et il faut pouvoir passer facilement entre les différentes expressions.

Par rapport aux variables, on a:

$$\begin{aligned}
\begin{matrix} X' \\ (N,p)(N,p) \end{matrix} X &= \begin{pmatrix} X^{(1)'} \\ X^{(2)'} \\ \vdots \\ X^{(p)'} \end{pmatrix} (X^{(1)} | X^{(2)} | \dots | X^{(p)}) \\
&= \begin{pmatrix} X^{(1)'} X^{(1)} & X^{(1)'} X^{(2)} & \dots & X^{(1)'} X^{(p)} \\ X^{(1)'} X^{(2)} & X^{(2)'} X^{(2)} & \dots & X^{(2)'} X^{(p)} \\ \vdots & \vdots & \ddots & \vdots \\ X^{(p)'} X^{(1)} & X^{(p)'} X^{(2)} & \dots & X^{(p)'} X^{(p)} \end{pmatrix} \\
&= \begin{pmatrix} \sum_{i=1}^N x_{i1}^2 & \sum_{i=1}^N x_{i1}x_{i2} & \dots & \sum_{i=1}^N x_{i1}x_{ip} \\ \sum_{i=1}^N x_{i1}x_{i2} & \sum_{i=1}^N x_{i2}^2 & \dots & \sum_{i=1}^N x_{i2}x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^N x_{i1}x_{ip} & \sum_{i=1}^N x_{i2}x_{ip} & \dots & \sum_{i=1}^N x_{ip}^2 \end{pmatrix}
\end{aligned}$$

La matrice des moments empiriques non centrés de X est définie par :

$$\frac{X'X}{N} = \begin{pmatrix} N^{-1} \sum_{i=1}^N x_{i1}^2 & \dots & N^{-1} \sum_{i=1}^N x_{i1}x_{ip} \\ N^{-1} \sum_{i=1}^N x_{i1}x_{i2} & \dots & N^{-1} \sum_{i=1}^N x_{i2}x_{ip} \\ \vdots & \ddots & \vdots \\ N^{-1} \sum_{i=1}^N x_{i1}x_{ip} & \dots & N^{-1} \sum_{i=1}^N x_{ip}^2 \end{pmatrix}$$

On en déduit la matrice de covariance empirique :

$$\begin{aligned}
V_e(X) &= \begin{pmatrix} N^{-1} \sum_{i=1}^N x_{i1}^2 & \dots & N^{-1} \sum_{i=1}^N x_{i1}x_{ip} \\ N^{-1} \sum_{i=1}^N x_{i1}x_{i2} & \dots & N^{-1} \sum_{i=1}^N x_{i2}x_{ip} \\ \vdots & \ddots & \vdots \\ N^{-1} \sum_{i=1}^N x_{i1}x_{ip} & \dots & N^{-1} \sum_{i=1}^N x_{ip}^2 \end{pmatrix} \\
&\quad - \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_p \end{pmatrix} \begin{pmatrix} \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_p \end{pmatrix}
\end{aligned}$$

$$\begin{aligned}
&= \begin{pmatrix} N^{-1} \sum_{i=1}^N x_{i1}^2 & \dots & N^{-1} \sum_{i=1}^N x_{i1}x_{ip} \\ N^{-1} \sum_{i=1}^N x_{i1}x_{i2} & \dots & N^{-1} \sum_{i=1}^N x_{i2}x_{ip} \\ \vdots & \ddots & \vdots \\ N^{-1} \sum_{i=1}^N x_{i1}x_{ip} & \dots & N^{-1} \sum_{i=1}^N x_{ip}^2 \end{pmatrix} - \begin{pmatrix} \bar{x}_1^2 & \dots & \bar{x}_1\bar{x}_p \\ \bar{x}_1\bar{x}_2 & \dots & \bar{x}_2\bar{x}_p \\ \vdots & \ddots & \vdots \\ \bar{x}_1\bar{x}_p & \dots & \bar{x}_p^2 \end{pmatrix} \\
&= \begin{pmatrix} N^{-1} \sum_{i=1}^N x_{i1}^2 - \bar{x}_1^2 & \dots & N^{-1} \sum_{i=1}^N x_{i1}x_{ip} - \bar{x}_1\bar{x}_p \\ N^{-1} \sum_{i=1}^N x_{i1}x_{i2} - \bar{x}_1\bar{x}_2 & \dots & N^{-1} \sum_{i=1}^N x_{i2}x_{ip} - \bar{x}_2\bar{x}_p \\ \vdots & \ddots & \vdots \\ N^{-1} \sum_{i=1}^N x_{i1}x_{ip} - \bar{x}_1\bar{x}_p & \dots & N^{-1} \sum_{i=1}^N x_{ip}^2 - \bar{x}_p^2 \end{pmatrix}
\end{aligned}$$

On obtient donc finalement :

$$V_e(X) = \begin{pmatrix} V_e(x_1) & \text{Cov}_e(x_1, x_2) & \dots & \text{Cov}_e(x_1, x_p) \\ \text{Cov}_e(x_1, x_2) & V_e(x_2) & \dots & \text{Cov}_e(x_2, x_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}_e(x_1, x_p) & \text{Cov}_e(x_2, x_p) & \dots & V_e(x_p) \end{pmatrix}$$

Par rapport aux observations. La matrice de covariance empirique peut s'écrire :

$$V_e(X) = \frac{1}{N} \sum_{i=1}^N X_i' X_i - \bar{X} \bar{X}'$$

on a :

$$\begin{aligned}
\sum_{i=1}^N X_i' X_i &= \sum_{i=1}^N (x_{i1}, x_{i2}, \dots, x_{ip}) \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} \\
&= \sum_{i=1}^N \begin{pmatrix} x_{i1}^2 & x_{i1}x_{i2} & \dots & x_{i1}x_{ip} \\ x_{i1}x_{i2} & x_{i2}^2 & \dots & x_{i2}x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ x_{i1}x_{ip} & x_{i2}x_{ip} & \dots & x_{ip}^2 \end{pmatrix} \\
&= \begin{pmatrix} \sum_{i=1}^N x_{i1}^2 & \sum_{i=1}^N x_{i1}x_{i2} & \dots & \sum_{i=1}^N x_{i1}x_{ip} \\ \sum_{i=1}^N x_{i1}x_{i2} & \sum_{i=1}^N x_{i2}^2 & \dots & \sum_{i=1}^N x_{i2}x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^N x_{i1}x_{ip} & \sum_{i=1}^N x_{i2}x_{ip} & \dots & \sum_{i=1}^N x_{ip}^2 \end{pmatrix} \\
&= X'X
\end{aligned}$$

On retrouve donc le même résultat que précédemment. De même pour les produits croisés entre les variables explicatives et la variable expliquée, on a :

$$\begin{pmatrix} X' \\ (N,p) \end{pmatrix} \begin{pmatrix} y \\ (N,1) \end{pmatrix} = \begin{pmatrix} X^{(1)'} \\ X^{(2)'} \\ \vdots \\ X^{(p)'} \end{pmatrix} y = \begin{pmatrix} X^{(1)'}y \\ X^{(2)'}y \\ \vdots \\ X^{(p)'}y \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N x_{i1}y_i \\ \sum_{i=1}^N x_{i2}y_i \\ \vdots \\ \sum_{i=1}^N x_{ip}y_i \end{pmatrix} = \sum_{i=1}^N X_i' y_i.$$

Les moments centrés donnent donc :

$$\text{Cov}_e(X, y) = \frac{1}{N} \sum_{i=1}^N X_i' y_i - \bar{X} \bar{y} = \frac{X'y}{N} - \bar{X} \bar{y}.$$

Le vecteur correspondant est égal à :

$$\text{Cov}_e(X, y) = \begin{pmatrix} N^{-1} \sum_{i=1}^N x_{i1}y_i \\ N^{-1} \sum_{i=1}^N x_{i2}y_i \\ \vdots \\ N^{-1} \sum_{i=1}^N x_{ip}y_i \end{pmatrix} - \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_p \end{pmatrix} \bar{y}$$

$$\begin{aligned}
&= \begin{pmatrix} N^{-1} \sum_{i=1}^N x_{i1}y_i \\ N^{-1} \sum_{i=1}^N x_{i2}y_i \\ \vdots \\ N^{-1} \sum_{i=1}^N x_{ip}y_i \end{pmatrix} - \begin{pmatrix} \bar{x}_1\bar{y} \\ \bar{x}_2\bar{y} \\ \vdots \\ \bar{x}_p\bar{y} \end{pmatrix} \\
&= \begin{pmatrix} N^{-1} \sum_{i=1}^N x_{i1}y_i - \bar{x}_1\bar{y} \\ N^{-1} \sum_{i=1}^N x_{i2}y_i - \bar{x}_2\bar{y} \\ \vdots \\ N^{-1} \sum_{i=1}^N x_{ip}y_i - \bar{x}_p\bar{y} \end{pmatrix} \\
&= \begin{pmatrix} \text{Cov}_e(x_1, y) \\ \text{Cov}_e(x_2, y) \\ \vdots \\ \text{Cov}_e(x_p, y) \end{pmatrix}.
\end{aligned}$$

Sous certaines conditions, les moments empiriques que nous venons de voir convergent en probabilité vers les moments théoriques correspondants. Ce point est examiné dans la section suivante.

A.3 Convergence en probabilité

DÉFINITION A.1 Soit \hat{b}_N une variable aléatoire dont la réalisation dépend du nombre d'observations disponibles dans un échantillon (noté N). On dit que cette suite de variables aléatoires \hat{b}_N converge en probabilité vers une valeur b lorsque le nombre d'observations N tend vers l'infini, si elle vérifie la propriété suivante :

$$\forall \varepsilon > 0, \Pr \left[\left| \hat{b}_N - b \right| > \varepsilon \right] \xrightarrow{N \rightarrow +\infty} 0.$$

La convergence en probabilité de \hat{b}_N vers b est notée de manière abrégée par l'expression :

$$\text{Plim } \hat{b}_N = b,$$

où Plim est l'abréviation de "probability limit" (i.e., limite en probabilité). Elle s'écrit également :

$$\hat{b}_N \xrightarrow[N \rightarrow +\infty]{P} b.$$

Cette définition signifie que l'évènement " \widehat{b}_N s'écarte de b d'une distance supérieure à ε " est de probabilité nulle (i.e., impossible) lorsque $N \rightarrow +\infty$. Cette propriété s'étend à certaines fonctions de \widehat{b}_N , comme le montre le théorème suivant.

THÉORÈME A.1 [Slutsky]

Soit \widehat{b}_N une suite de variables aléatoires qui converge en probabilité vers b :

$$\text{Plim } \widehat{b}_N = b,$$

et soit $g(\cdot)$ une fonction continue définie au point b . On a :

$$\text{Plim } g(\widehat{b}_N) = g(\text{Plim } \widehat{b}_N) = g(b).$$

Les définitions précédentes et le théorème de Slutsky s'étendent au cas vectoriel en raisonnant composante par composante. En particulier le théorème de Slutsky permet de simplifier considérablement le calcul des limites en probabilités. Prenons deux estimateurs convergents, \widehat{a} d'un paramètre a et \widehat{b} d'un paramètre b . On a :

$$\text{Plim } \widehat{a} + \widehat{b} = \text{Plim } \widehat{a} + \text{Plim } \widehat{b} = a + b,$$

car la fonction $g(a, b) = a + b$ est continue et les estimateurs convergent en probabilité. De même, en utilisant les fonctions $g(a, b) = ab$, $g(a, b) = a/b$ (pour $b \neq 0$) on obtient les propriétés :

$$\begin{aligned} \text{Plim } \widehat{a} \widehat{b} &= \text{Plim } \widehat{a} \text{Plim } \widehat{b} = ab, \\ \text{Plim } \frac{\widehat{a}}{\widehat{b}} &= \frac{\text{Plim } \widehat{a}}{\text{Plim } \widehat{b}} = \frac{a}{b}, \quad b \neq 0. \end{aligned}$$

A.4 Inégalité de Bienaymé-Chebichev

Le théorème suivant est très important. Il nous permet notamment de démontrer la loi des grands nombres et le fait que la convergence en moyenne quadratique implique la convergence en probabilité...en une seule ligne.

THÉORÈME A.2 [Inégalité de Bienaymé-Chebichev]

Soit Z une variable de carré intégrable, on a :

$$\forall \delta > 0, \text{Pr } [|Z| \geq \delta] \leq \frac{1}{\delta^2} \text{E}(Z^2).$$

PREUVE :

Soit la variable de Bernoulli :

$$D = \begin{cases} 1 & \text{si } |Z| \geq \delta \\ 0 & \text{sinon} \end{cases}$$

son espérance mathématique est égale à :

$$\mathbb{E}(D) = 1 \times \Pr[|Z| \geq \delta] + 0 \times \Pr[|Z| < \delta] = \Pr[|Z| \geq \delta].$$

D'autre part :

1. Si $|Z| \geq \delta$ on a $D = 1$ donc :

$$\frac{|Z|}{\delta} \geq 1 \Rightarrow \frac{Z^2}{\delta^2} \geq D = 1.$$

2. Si $|Z| < \delta$ on a $D = 0$ donc :

$$\frac{|Z|}{\delta} \geq 0 \Rightarrow \frac{Z^2}{\delta^2} \geq D = 0.$$

donc dans tous les cas on a :

$$\begin{aligned} \frac{Z^2}{\delta^2} \geq D &\Rightarrow \mathbb{E}\left(\frac{Z^2}{\delta^2}\right) \geq \mathbb{E}(D) \\ &\Leftrightarrow \frac{1}{\delta^2} \mathbb{E}(Z^2) \geq \Pr[|Z| \geq \delta]. \end{aligned}$$

□

Remarque A.1 En posant $Z = X - \mathbb{E}(X)$, on obtient l'expression :

$$\forall \delta > 0, \Pr[|X - \mathbb{E}(X)| \geq \delta] \leq \frac{1}{\delta^2} \mathbb{V}(X),$$

car $\mathbb{V}(X) = \mathbb{E}\left[(X - \mathbb{E}(X))^2\right]$.

La convergence en probabilité est parfois difficile à vérifier directement, on utilise alors une conditions suffisante, qui correspond en fait à la convergence en moyenne quadratique.

DÉFINITION A.2 Soit \widehat{b}_N une variable aléatoire dont la réalisation dépend du nombre d'observations disponibles dans un échantillon (noté N). On dit que cette suite de variables aléatoires \widehat{b}_N converge en moyenne quadratique vers une valeur b lorsque le nombre d'observations N tend vers l'infini, si elle vérifie une des deux propriétés équivalentes suivantes :

1. $\mathbb{E}\left[\left(\widehat{b}_N - b\right)^2\right] \rightarrow 0$ lorsque $N \rightarrow +\infty$.
2. $\mathbb{E}\left(\widehat{b}_N\right) \rightarrow b$ et $\mathbb{V}\left(\widehat{b}_N\right) \rightarrow 0$ lorsque $N \rightarrow +\infty$.

On note ce résultat :

$$\widehat{b}_N \xrightarrow[N \rightarrow +\infty]{m.q.} b.$$

Cette définition porte directement sur la distance entre \widehat{b}_N et b . Elle impose que cette distance s'annule quand le nombre d'observations devient suffisamment grand. L'équivalence entre les deux définitions vient du développement suivant :¹

$$\begin{aligned} \mathbb{E} \left[\left(\widehat{b}_N - b \right)^2 \right] &= \mathbb{V} \left[\widehat{b}_N - b \right] + \left[\mathbb{E} \left(\widehat{b}_N - b \right) \right]^2 \\ &= \mathbb{V} \left(\widehat{b}_N \right) + \left(\mathbb{E} \left(\widehat{b}_N \right) - b \right)^2 \geq 0. \end{aligned}$$

Les deux termes précédents sont positifs ou nuls donc pour que l'expression s'annule lorsque $N \rightarrow +\infty$, il faut que l'on ait simultanément $\mathbb{V} \left(\widehat{b}_N \right) \rightarrow 0$ et $\mathbb{E} \left(\widehat{b}_N \right) \rightarrow b$.

PROPRIÉTÉ A.1 Soit \widehat{b}_N une suite de variables aléatoires, on a :

$$\widehat{b}_N \xrightarrow[N \rightarrow +\infty]{m.q.} b \quad \Rightarrow \quad \text{Plim } \widehat{b}_N = b,$$

la convergence en moyenne quadratique implique la convergence en probabilité.

PREUVE :

C'est une conséquence de l'inégalité de Bienaymé-Chebichev. En posant $Z = \widehat{b}_N - b$ et $\delta = \varepsilon > 0$ dans le théorème [A.2], on obtient :

$$\forall \varepsilon > 0, \quad 0 \leq \Pr \left[\left| \widehat{b}_N - b \right| \geq \varepsilon \right] \leq \frac{1}{\varepsilon^2} \mathbb{E} \left[\left(\widehat{b}_N - b \right)^2 \right] \xrightarrow[N \rightarrow +\infty]{} 0.$$

□

A.5 La loi faible des grands nombres

Cette section permet de faire le lien entre les moments empiriques que nous avons vu plus haut et la convergence en probabilité que nous venons de voir. Elle signifie que sous certaines conditions, les moments empiriques convergent en probabilité vers les moments théoriques correspondants. On l'appelle loi faible des grands nombres, car la convergence en probabilité est également appelée convergence faible. La version de cette loi que nous utilisons est due à Markov (cf. Petrov 1995, p.134).

¹On rappelle que : $\mathbb{V}(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2 \Leftrightarrow \mathbb{E}(X^2) = \mathbb{V}(X) + \mathbb{E}(X)^2$. Ici on pose $X = \widehat{b}_n - b$.

THÉORÈME A.3 [Markov]

Soit (X_1, \dots, X_N) une suite de variables aléatoires qui admettent une espérance mathématique $E(X_k) = m_k$ pour toute valeur de $k \in \{1, \dots, N\}$, et qui vérifient la propriété suivante :

$$\frac{1}{N^2} V \left[\sum_{k=1}^N X_k \right] \rightarrow 0 \quad \text{lorsque } N \rightarrow +\infty,$$

alors

$$\text{Plim} \left[\frac{1}{N} \sum_{k=1}^N X_k - \frac{1}{N} \sum_{k=1}^N m_k \right] = 0.$$

PREUVE :

Il suffit de poser $Z = N^{-1} \sum_{k=1}^N (X_k - m_k)$ dans l'inégalité de Bienaymé-Chebichev (théorème [A.2]) :

$$\forall \delta > 0, \Pr \left[\left| \frac{1}{N} \sum_{k=1}^N X_k - \frac{1}{N} \sum_{k=1}^N m_k \right| \geq \delta \right] \leq \frac{1}{\delta^2 N^2} V \left[\sum_{k=1}^N X_k \right] \xrightarrow{N \rightarrow +\infty} 0.$$

En effet, on a :

$$E(Z) = \frac{1}{N} \sum_{k=1}^N [E(X_k) - m_k] = 0$$

$$V(Z) = V \left[\frac{1}{N} \sum_{k=1}^N X_k - \frac{1}{N} \sum_{k=1}^N m_k \right] = V \left[\frac{1}{N} \sum_{k=1}^N X_k \right]$$

car $N^{-1} \sum_{k=1}^N m_k$ est une quantité certaine et que l'on a : $V \left[N^{-1} \sum_{k=1}^N X_k \right] = N^{-2} V \left[\sum_{k=1}^N X_k \right]$.

□

Une moyenne arithmétique de variable aléatoires converge donc vers la moyenne des espérances mathématiques des variables aléatoires, à condition que la variance de leur moyenne $V \left[N^{-1} \sum_{k=1}^N X_k \right]$ tende vers 0 lorsque $N \rightarrow +\infty$.

Exemple A.1 On considère un échantillon de variables (X_1, \dots, X_k) indépendantes, d'espérance et de variance constantes : $\forall k, m_k = m$ et $V(X_k) = \sigma^2$. Sous hypothèse d'indépendance, on obtient la condition suivante :

$$\frac{1}{N^2} V \left[\sum_{k=1}^N X_k \right] = \frac{1}{N^2} \sum_{k=1}^N V(X_k) = \frac{N\sigma^2}{N^2} = \frac{\sigma^2}{N} \rightarrow 0 \quad \text{quand } N \rightarrow +\infty.$$

D'autre part $N^{-1} \sum_{k=1}^N m_k = N^{-1} (N \times m) = m$. On a donc le résultat de convergence suivant :

$$\text{Plim } \bar{X} = m,$$

la moyenne empirique converge vers l'espérance mathématique commune des variables (X_1, \dots, X_k) .

Exemple A.2 On considère un échantillon de variables (X_1, \dots, X_k) indépendantes de variances différentes et finies : $V(X_k) = \sigma_k^2$. La moyenne arithmétique de ces variances $N^{-1} \sum_{k=1}^N \sigma_k^2 = \bar{\sigma}$ est également finie. En effet :

$$\bar{\sigma} \leq \max_{k=1, \dots, N} \sigma_k^2 \text{ qui est finie.}$$

ce qui implique :

$$\frac{1}{N^2} V \left[\sum_{k=1}^N X_k \right] = \frac{1}{N^2} \sum_{k=1}^N \sigma_k^2 = \frac{\bar{\sigma}}{N} \rightarrow 0 \quad \text{quand } N \rightarrow +\infty.$$

On en déduit que :

$$\text{Plim } \bar{X} = \text{Plim } \frac{1}{N} \sum_{k=1}^N E(X_k).$$

A.6 Théorème de la limite centrale

Le théorème suivant nous permet de déterminer la loi asymptotique de la plupart de nos estimateurs.

THÉORÈME A.4 (Liapunov) Soit u_1, u_2, \dots, u_N une suite de variables aléatoires indépendantes d'espérances mathématiques $E(u_i) = \mu_i$ et de variances respectives $V(u_i) = E(u_i - \mu_i)^2 = \sigma_i^2 \neq 0$, $i = 1, \dots, n$. On suppose également que le moment absolu d'ordre trois existe $E|u_i - \mu_i|^3 = \beta_i \forall i$. Soient :

$$B_N = \left(\sum_{i=1}^N \beta_i \right)^{1/3}, D_N = \left(\sum_{i=1}^N \sigma_i^2 \right)^{1/2},$$

alors, si $\lim B_N/D_N = 0$ lorsque $N \rightarrow +\infty$, on a :

$$\sum_{i=1}^N \frac{u_i - \mu_i}{D_N} \xrightarrow{N \rightarrow +\infty} N(0, 1).$$

Annexe B

Algèbre linéaire

B.1 Calcul matriciel

On considère une matrice $A = [A_{ij}]$ de format (m, n) .

1. La transposée de A , notée A' , est définie par $A' = [A_{ji}]$, on intervertit donc les lignes et les colonnes.
2. A est de plein rang colonne si ses colonnes sont linéairement indépendantes. C'est-à-dire si :

$$\forall \alpha \in \mathbb{R}^n, \quad A\alpha = 0 \Rightarrow \alpha = 0.$$

3. A est de plein rang ligne si ses lignes sont linéairement indépendantes (i.e., si A' est de plein rang colonne).

On considère maintenant deux matrices A de format (m, n) et B de format (r, p) .

1. Le produit matriciel de A par B n'existe que si le nombre de colonnes de A est égal au nombre de lignes de B : $n = r$. Dans ce cas, on le note $F = AB$ et il est de format (m, p) .
2. Les éléments de la matrice produit $F = [F_{ij}] = AB$ sont définis comme les produits scalaires de la i -ème ligne de A et de la j -ième colonne de B .
3. AB n'est généralement pas égal à BA , le produit matriciel n'est pas commutatif.
4. $A(B + C) = AB + AC$.

$$5. (A + B)C = AC + BC.$$

$$6. (AB)' = B'A'.$$

On considère maintenant deux matrices carrées A de format (m, m) et B de format (r, r) .

1. Une matrice est carrée si elle a autant de lignes que de colonnes.
2. Une matrice carrée A est symétrique si $A' = A$.
3. La trace d'une matrice carrée A est définie par la somme de ses éléments diagonaux. On la note $\text{tr}(A) = \sum_{i=1}^m A_{ii}$.
4. $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$.
5. Si ABC est une matrice carrée et si les formats sont compatibles : $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$. Il n'est pas nécessaire que chaque matrice soit carrée à l'intérieur des produits précédents.
6. Si A est une matrice carrée de plein rang (ligne ou colonne), elle admet une inverse notée A^{-1} telle que $AA^{-1} = A^{-1}A = I$.
7. Si les matrices A et B sont inversibles : $(AB)^{-1} = B^{-1}A^{-1}$.
8. Une matrice carrée A est idempotente si $A^2 = A$.

B.2 Matrices définies positives

DÉFINITION B.1 Une matrice A de format (m, m) est semi définie positive lorsque :

$$\forall \alpha \in \mathbb{R}^m, s(\alpha, A) = \alpha' A \alpha \geq 0.$$

DÉFINITION B.2 Une matrice A de format (m, m) est définie positive lorsque :

$$\forall \alpha \in \mathbb{R}^m, \alpha \neq 0, s(\alpha, A) = \alpha' A \alpha > 0.$$

La propriété suivante est utile pour comparer les variances des différents estimateurs.

PROPRIÉTÉ B.1 Soit $X_{(n,p)}$ une matrice quelconque, alors $X'X$ est semi définie positive.

PREUVE :

En posant $A = X'X$, on obtient :

$$s(\alpha, X'X) = \alpha' X'X \alpha = \underbrace{(X\alpha)'}_{(1,n)} \underbrace{(X\alpha)}_{(n,1)} = \|X\alpha\|^2 \geq 0.$$

□

La propriété suivante est utile pour montrer l'existence de certains estimateurs.

PROPRIÉTÉ B.2 Soit $X_{(n,p)}$ une matrice de plein rang colonne, $\text{rang}(X) = p$, alors $X'X$ est définie positive (donc de rang égal à p).

PREUVE :

La matrice X est de plein rang colonne :

$$\forall \alpha \in \mathbb{R}^p, \quad X'\alpha = 0 \Rightarrow \alpha = 0$$

donc $\|X\alpha\|^2$ ne peut être nul que dans le cas $\alpha = 0$. En conséquence :

$$\forall \alpha \in \mathbb{R}^p, \alpha \neq 0, \quad \|X\alpha\|^2 > 0.$$

□

B.3 Produits de Kronecker

Soient deux matrices $A = [A_{ij}]$ de format (a, b) et $B = [B_{ij}]$ de format (c, d) . Le produit de Kronecker de la matrice A par la matrice B , noté $A \otimes B$, donne une matrice $F = [F_{ij}]$ de format (ac, bd) . Cette matrice est définie par :

$$F = [A_{ij}B] = \begin{pmatrix} A_{1,1}B & A_{1,2}B & \cdots & A_{1,b}B \\ A_{2,1}B & A_{2,2}B & \cdots & A_{2,b}B \\ \vdots & \vdots & \ddots & \vdots \\ A_{a,1}B & A_{a,2}B & \cdots & A_{a,b}B \end{pmatrix},$$

chaque élément originel de la matrice A se voit multiplié par la totalité de la matrice B . Chacun des éléments de la matrice ci-dessus est donc de dimensions égales à celles de B , et C est de format (ac, bd) . Les propriétés suivantes sont valables sous réserve que les formats des matrices autorisent les multiplications matricielles indiquées.

1. Dans le cas général $(A \otimes B)$ n'est pas égal à $(B \otimes A)$, le produit de Kronecker n'est donc pas commutatif.
2. $0 \otimes A = 0$.
3. $A \otimes 0 = 0$, mais attention, le format de ce 0 n'est pas nécessairement le même que celui de la propriété précédente.
4. $A \otimes (B + C) = A \otimes B + A \otimes C$.

5. $(A + B) \otimes C = A \otimes B + B \otimes C.$
6. $\forall (x, y) \in \mathbb{R}^2, (xA) \otimes (yB) = xy(A \otimes B).$
7. $(A \otimes B)(C \otimes D) = (AC \otimes BD).$
8. $(A \otimes B)' = (A' \otimes B').$
9. Si A et B sont inversibles : $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}.$
10. $\text{tr}(A \otimes B) = \text{tr } A. \text{tr } B.$

Annexe C

La loi normale

La loi normale centrée réduite admet pour densité :

$$\phi(u) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{u^2}{2}\right\},$$

elle est d'espérance nulle et variance unitaire. Plus généralement, on peut définir une loi normale d'espérance m et de variance σ^2 en définissant la variable suivante :

$$Y = g(U) = m + \sigma U, \quad U \rightsquigarrow N(0, 1),$$

la réciproque de la fonction est :

$$g^{-1}(y) = \frac{y - m}{\sigma},$$

et la densité de Y est donnée par :

$$\begin{aligned} f(y) &= \phi(g^{-1}(y)) \left| \frac{dg^{-1}(y)}{dy} \right| \\ &= \frac{1}{\sigma} \phi\left(\frac{y - m}{\sigma}\right). \end{aligned}$$

La fonction génératrice des moments de la loi normale centrée réduite est définie par :

$$\begin{aligned} M(s) &= E(e^{sU}) \\ &= \int_{-\infty}^{+\infty} \exp(su) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du \\ &= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(u^2 - 2su)\right) du \end{aligned}$$

C.1 Loi normale univariée tronquée

On considère une loi normale de densité :

$$f(y) = \frac{1}{\sigma} \phi\left(\frac{y-m}{\sigma}\right),$$

et on cherche la densité de la loi tronquée en un seuil c . La densité de la loi tronquée est définie par :

$$f(y|y > c) = \frac{f(y) 1_{(y>c)}}{\Pr(y > c)}.$$

Pour calculer l'espérance mathématique de la loi tronquée, il nous faut la quantité :

$$\begin{aligned} I &= \int_{-\infty}^{+\infty} y f(y) 1_{(y>c)} dy \\ &= \int_c^{+\infty} y \frac{1}{\sigma} \phi\left(\frac{y-m}{\sigma}\right) dy, \end{aligned}$$

on fait le changement de variable :

$$z = \frac{y-m}{\sigma},$$

ce qui implique :

$$\lim_{y \rightarrow c} z = \frac{y-c}{\sigma}, \quad \lim_{y \rightarrow +\infty} z = +\infty, \quad y = m + \sigma z \quad \text{et} \quad dy = \sigma dz,$$

on obtient donc :

$$\begin{aligned} I &= \int_{(c-m)/\sigma}^{+\infty} (m + \sigma z) \phi(z) dz \\ &= m \int_{(c-m)/\sigma}^{+\infty} \phi(z) dz + \sigma \int_{(c-m)/\sigma}^{+\infty} z \phi(z) dz \\ &= m(1 - \Phi((c-m)/\sigma)) + \sigma \int_{(c-m)/\sigma}^{+\infty} -\phi'(z) dz \\ &= m\Phi((m-c)/\sigma) + \sigma\phi((m-c)/\sigma), \end{aligned}$$

d'autre part :

$$\begin{aligned}
 \Pr(y > c) &= \Pr\left(\frac{y-m}{\sigma} > \frac{c-m}{\sigma}\right) \\
 &= 1 - \Pr\left(\frac{y-m}{\sigma} \leq \frac{c-m}{\sigma}\right) \\
 &= 1 - \Phi\left(\frac{c-m}{\sigma}\right) \\
 &= \Phi\left(\frac{m-c}{\sigma}\right),
 \end{aligned}$$

ce qui implique :

$$E(y|y > m) = \frac{I}{\Pr(y > c)} = m + \sigma \frac{\phi\left(\frac{m-c}{\sigma}\right)}{\Phi\left(\frac{m-c}{\sigma}\right)}, \quad (\text{C.1})$$

la quantité :

$$\lambda\left(\frac{m-c}{\sigma}\right) = \frac{\phi\left(\frac{m-c}{\sigma}\right)}{\Phi\left(\frac{m-c}{\sigma}\right)},$$

est égale à l'inverse du ratio de Mills.

C.2 Loi normale bivariée

Cette annexe présente la loi normale bivariée ainsi que les distributions conditionnelles qui y sont associées. On se limite ici à deux variables, mais l'extension à un nombre quelconque est possible. On considère deux variables (y_1, y_2) d'espérance (m_1, m_2) et de matrice de covariance :

$$V \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix},$$

La loi normale bivariée est définie par la densité :

$$\begin{aligned}
 \phi_2(y_1, y_2) &= \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \times \\
 \exp\left\{-\frac{1}{2(1-\rho^2)} \left(\left[\frac{y_1-m_1}{\sigma_1}\right]^2 + \left[\frac{y_2-m_2}{\sigma_2}\right]^2 - 2\rho\frac{y_1-m_1}{\sigma_1}\frac{y_2-m_2}{\sigma_2} \right)\right\}
 \end{aligned}$$

C.3 Loi normale conditionnelle

La densité de la loi conditionnelle de y_1 sachant y_2 est définie par :

$$\phi_c(y_1|y_2) = \frac{\phi_2(y_1, y_2)}{f_2(y_2)}, \quad (\text{C.2})$$

où $f_2(y_2)$ est la densité marginale de y_2 . On a :

$$f_2(y_2) = \frac{1}{\sigma_2} \phi\left(\frac{y_2 - m_2}{\sigma_2}\right),$$

où $\phi(\cdot)$ est la densité de la loi normale centrée réduite. En prenant le ratio (C.2) on obtient :

$$\begin{aligned} \phi_c(y_1|y_2) &= \frac{1}{\sigma_1\sqrt{1-\rho^2}} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2(1-\rho^2)}\left(\frac{y_1 - m_1}{\sigma_1} - \rho\frac{y_2 - m_2}{\sigma_2}\right)^2\right\} \\ &= \frac{1}{\sigma_1\sqrt{1-\rho^2}} \phi\left(\frac{1}{\sigma_1\sqrt{1-\rho^2}}\left[y_1 - m_1 - \rho\frac{\sigma_1}{\sigma_2}(y_2 - m_2)\right]\right) \\ &= \frac{1}{\sigma_1\sqrt{1-\rho^2}} \phi\left(\frac{1}{\sigma_1\sqrt{1-\rho^2}}\left[y_1 - \left(m_1 + \rho\frac{\sigma_1}{\sigma_2}(y_2 - m_2)\right)\right]\right) \end{aligned}$$

il s'agit de la densité d'une loi normale d'espérance :

$$E(y_1|y_2) = m_1 + \rho\frac{\sigma_1}{\sigma_2}(y_2 - m_2).$$

et de variance :

$$V(y_1|y_2) = \sigma_1^2(1 - \rho^2).$$

Plus généralement, on peut montrer directement que :

$$y_1 = m_1 + \rho\frac{\sigma_1}{\sigma_2}(y_2 - m_2) + \varepsilon_1, \quad (\text{C.3})$$

où ε_1 est une perturbation normale indépendante de y_2 . Pour voir cela, il suffit de remarquer que ε_1 est une combinaison linéaire de lois normales et est donc une variable normale. Pour l'indépendance, il suffit de calculer la covariance entre ε_1 et y_2 , puisque l'indépendance est équivalente à l'absence de corrélation pour cette loi. On a :

$$\begin{aligned} \text{Cov}(\varepsilon_1, y_2) &= \text{Cov}\left(y_1 - m_1 - \rho\frac{\sigma_1}{\sigma_2}(y_2 - m_2), y_2\right) \\ &= \text{Cov}(y_1, y_2) - \rho\frac{\sigma_1}{\sigma_2}V(y_2) \\ &= \rho\sigma_1\sigma_2 - \rho\frac{\sigma_1}{\sigma_2}\sigma_2^2 \\ &= 0, \end{aligned}$$

de sorte que ε_1 et y_2 sont indépendantes. L'espérance de ε_1 est nulle :

$$E(\varepsilon_1) = E\left(y_1 - m_1 - \rho\frac{\sigma_1}{\sigma_2}(y_2 - m_2)\right) = 0,$$

et sa variance est égale à :

$$\begin{aligned}
 V(\varepsilon_1) &= V\left(y_1 - m_1 - \rho \frac{\sigma_1}{\sigma_2} (y_2 - m_2)\right) \\
 &= \sigma_1^2 V\left(\frac{y_1 - m_1}{\sigma_1} - \rho \frac{y_2 - m_2}{\sigma_2}\right) \\
 &= \sigma_1^2 (1 + \rho^2 - 2\rho^2) \\
 &= \sigma_1^2 (1 - \rho^2).
 \end{aligned}$$

La propriété (C.3) est très pratique lorsque l'on étudie la troncature d'une variable normale par une autre variable normale.

C.4 Loi normale bivariée tronquée

Ici, on recherche l'espérance conditionnelle d'une première variable tronquée par la valeur d'une seconde variable avec laquelle elle est corrélée. Ce cas se retrouve lorsque l'on estime une équation de salaire en tenant compte de la participation. On cherche donc la valeur de l'espérance conditionnelle suivante :

$$\begin{aligned}
 E(y_1 | y_2 > c) &= E\left(m_1 + \rho \frac{\sigma_1}{\sigma_2} (y_2 - m_2) + \varepsilon_1 \mid y_2 > c\right) \\
 &= m_1 + \rho \frac{\sigma_1}{\sigma_2} (E(y_2 | y_2 > c) - m_2) + E(\varepsilon_1 | y_2 > c) \\
 &= m_1 + \rho \frac{\sigma_1}{\sigma_2} \left(m_2 + \sigma_2 \frac{\phi\left(\frac{m_2 - c}{\sigma_2}\right)}{\Phi\left(\frac{m_2 - c}{\sigma_2}\right)} - m_2\right) + E(\varepsilon_1),
 \end{aligned}$$

en utilisant la propriété (C.1) et l'indépendance entre y_2 et ε_1 . Après simplification, on obtient :

$$E(y_1 | y_2 > c) = m_1 + \rho \sigma_1 \lambda\left(\frac{m_2 - c}{\sigma_2}\right). \quad (\text{C.4})$$

Annexe D

Simplification du calcul des dérivées

La plupart des modèles font intervenir dans la log-vraisemblance des termes linéaires :

$$m(X, b) = Xb$$

où X est un vecteur ligne $1 \times p$ et b un vecteur colonne $1 \times p$. Cette propriété vient de la forme latente linéaire de la plupart des modèles à variables qualitatives. Pour estimer le modèle, on a besoin des dérivées de la fonction à maximiser par rapport à un vecteur b . En fait, nous allons voir un certain nombre de simplifications qui permettent de se limiter à des dérivées par rapport à une variable réelle non vectorielle.

Tout d'abord, la fonction à maximiser est la somme de N fonctions qui ne diffèrent que par les valeurs que prennent les variables expliquée et explicatives. La forme fonctionnelle reste la même quelle que soit l'observation. Dans le cas le plus simple :

$$\ell(y, X, b) = \sum_{i=1}^N g(y_i, X_i, b),$$

où la fonction g est identique pour tous les individus. C'est la forme que l'on obtient systématiquement sous hypothèse d'indépendance, où la fonction g est le logarithme de la densité de probabilité. Elle se simplifie souvent comme :

$$\ell(y, X, b) = \sum_{i=1}^N g(y_i, m_i, \xi) \text{ avec } m_i = X_i b,$$

où ξ est un paramètre, indépendant de b , généralement du second ordre (i.e., de variance ou de corrélation). Les observations des variables explicatives pour le i -ème individu sont rangées dans un vecteur ligne $X_i =$

(X_{1i}, \dots, X_{pi}) et b est le vecteur colonne correspondant $b = (b_1, \dots, b_p)'$. On traite les paramètres de ξ séparément comme des paramètres réels car ils sont en petit nombre dans les cas usuels. La dérivation ne pose donc pas de problème particulier par rapport à ξ . Nous sommes donc ramenés au calcul de la dérivée par rapport au *vecteur* b . Il est clair qu'il suffit de dériver :

$$g(y_i, m_i, \xi),$$

et de faire la somme des dérivées ensuite. Ceci est valable aussi bien pour les dérivées premières que pour les dérivées secondes. En dérivant en chaîne, on a :

$$\frac{\partial g}{\partial b}(y_i, m_i, \xi) = \frac{\partial g}{\partial m_i}(y_i, m_i, \xi) \frac{\partial m_i}{\partial b}.$$

La première dérivée est celle d'une fonction réelle et s'effectue comme d'habitude. La deuxième dérivée est obtenue en empilant les dérivées dans le même sens que le vecteur par rapport auquel on dérive (i.e., en ligne ou en colonne). En effet, par convention :

$$\frac{\partial m_i}{\partial b} = \frac{\partial m_i}{\partial \begin{pmatrix} b_1 \\ \vdots \\ b_p \end{pmatrix}} = \begin{pmatrix} \partial m_i / b_1 \\ \vdots \\ \partial m_i / b_p \end{pmatrix},$$

en remarquant que $m_i = X_i b = X_{1i} b_1 + \dots + X_{pi} b_p$, on a $\partial m_i / b_j = X_{ji}$ pour $j = 1, \dots, p$. Donc, en empilant ces dérivées :

$$\frac{\partial m_i}{\partial b} = \begin{pmatrix} X_{1i} \\ \vdots \\ X_{pi} \end{pmatrix} = X_i'.$$

En conséquence le gradient, pour une observation i , est égal à :

$$\frac{\partial g}{\partial b}(y_i, m_i, \xi) = X_i' \underbrace{\frac{\partial g}{\partial m_i}(y_i, m_i, \xi)}_{\in R}.$$

Le calcul du hessien est simplifié du fait de la linéarité de m_i en b , qui implique que la dérivée seconde de m_i par rapport à b est nulle :

$$\begin{aligned}
\frac{\partial^2 g}{\partial b \partial b'}(y_i, m_i, \xi) &= \frac{\partial^2 g}{\partial m_i^2}(y_i, m_i, \xi) \frac{\partial m_i}{\partial b} \frac{\partial m_i}{\partial b'} + \frac{\partial g}{\partial m_i}(y_i, m_i, \xi) \underbrace{\frac{\partial^2 m_i}{\partial b \partial b'}}_{=0} \\
&= \frac{\partial^2 g}{\partial m_i^2}(y_i, m_i, \xi) \frac{\partial m_i}{\partial b} \frac{\partial m_i}{\partial b'} \\
&= X_i' X_i \underbrace{\frac{\partial^2 g}{\partial m_i^2}(y_i, m_i, \xi)}_{\in R}.
\end{aligned}$$

Le calcul se fait donc en trois étapes :

1. Calcul des dérivées première et seconde par rapport à une variable réelle $m = Xb$.
2. Multiplication par X_i' pour le gradient et par $X_i' X_i$ pour le hessien.
3. Addition des dérivées sur l'ensemble des observations.

Exemple D.1 *Nous verrons plus loin que la log-vraisemblance du modèle Logit pour une observation i peut s'écrire sous la forme :*

$$\ell_i = y_i \ln p(m_i) + (1 - y_i) \ln(1 - p(m_i)),$$

où $y_i \in \{0, 1\}$ est la réponse qualitative que l'on étudie (0 pour "non" et 1 pour "oui") et p la fonction de répartition de la loi logistique. Comme précédemment $m_i = X_i b$ résume l'influence des variables explicatives X_i sur le choix y_i qui a été effectué par l'individu i . On remarque dès maintenant que la fonction de répartition de la loi logistique est égale à :

$$p(m) = \frac{1}{1 + \exp(-m)}, \quad X \in R,$$

ce qui entraîne que

$$p'(m) = \frac{\exp(-m)}{1 + \exp(-m)} = p(m)(1 - p(m)).$$

Nous pouvons donc écrire la log-vraisemblance sous la forme :

$$\ell(b) = \sum_{i=1}^N g(y_i, m_i) \text{ avec } g(y_i, m) = y_i \ln p(m) + (1 - y_i) \ln(1 - p(m)),$$

$\forall m \in R.$

La dérivée première de g par rapport à m est égale à :

$$\frac{\partial g}{\partial m}(y_i, m) = y_i \frac{p'(m)}{p(m)} - (1 - y_i) \frac{p'(m)}{1 - p(m)} = y_i - p(m),$$

après simplification. La dérivée seconde est égale à :

$$\frac{\partial^2 g}{\partial m^2}(y_i, X) = -p'(m) = -p(m)(1 - p(m)).$$

Le score s'écrit donc :

$$s(b) = \sum_{i=1}^N X_i' \frac{\partial g}{\partial m}(y_i, m_i) = \sum_{i=1}^N X_i'(y_i - p(m_i)),$$

et le hessien est égal à :

$$H(b) = \sum_{i=1}^N X_i' X_i \frac{\partial^2 g}{\partial m^2}(y_i, m_i) = - \sum_{i=1}^N X_i' X_i p(m_i)(1 - p(m_i)).$$