

MODELISATION DE DONNÉES QUALITATIVES

REGRESSION LOGISTIQUE MULTIPLE

Pierre-Louis Gonzalez

RÉGRESSION LOGISTIQUE MULTIPLE

(cas Y binaire)

La régression logistique multiple généralise la régression logistique simple au cas où il y a plusieurs variables explicatives $\mathbf{X}_1 \dots \mathbf{X}_k$.

I - LE MODÈLE

On note $\mathbf{x} = (\mathbf{x}_1 \dots, \mathbf{x}_k)$ une valeur de $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_k)$.

$$\Pi(\mathbf{x}) = \text{Pr ob}(Y = 1 \mid \mathbf{X} = \mathbf{x}) = \frac{e^{\beta_0 + \beta_1 \mathbf{x}_1 + \dots + \beta_k \mathbf{x}_k}}{1 + e^{\beta_0 + \beta_1 \mathbf{x}_1 + \dots + \beta_k \mathbf{x}_k}}$$

1 Estimation par le maximum de vraisemblance

On a :

$$\text{Log } L(\beta) = \sum_{i=1}^n [y_i \text{Log } \Pi_i + (1 - y_i) \text{Log } (1 - \Pi_i)]$$

$$\text{où } \Pi_i = \Pi(\mathbf{x}_i) = \Pi(\mathbf{x}_{i1}, \dots, \mathbf{x}_{ik})$$

On obtient $\hat{\beta}$ en annulant les dérivées partielles :

$$\frac{\partial \text{Log } L(\beta)}{\partial \beta_j} = \sum_{i=1}^n x_{ij} (y_i - \Pi_i) = 0 \quad j=0, \dots, k$$

$$(\text{où } \mathbf{x}_{i0} = \mathbf{1}, \text{ pour tout } \mathbf{i})$$

Matrice Variances-Covariances de $\hat{\beta}$:

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{1k} \\ \vdots & & \\ \vdots & & \\ 1 & x_{n1} & x_{nk} \end{pmatrix} \quad \mathbf{V} = \begin{pmatrix} \hat{\Pi}_1(1-\hat{\Pi}_1) & & \\ & \ddots & 0 \\ 0 & & \ddots \\ & & & \hat{\Pi}_n(1-\hat{\Pi}_n) \end{pmatrix}$$

$$\text{On a } \mathbf{Var}(\hat{\beta}) = (\mathbf{X}'\mathbf{V}\mathbf{X})^{-1}$$

2 Influence globale des variables

Modèle $\Pi(\mathbf{x}) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k}}$

Test
$$\begin{cases} \mathbf{H}_0 : \beta_1 = \dots = \beta_k = \mathbf{0} \\ \mathbf{H}_1 : \text{au moins un } \beta_j \neq \mathbf{0} \end{cases}$$

Trois statistiques

▶ Statistique de Wald

▶ Rapport des vraisemblances

$$\Lambda = -2 \text{ Log} \left[\frac{L(\text{constante})}{L(\text{toutes les variables})} \right] \rightarrow \chi^2(\mathbf{k}) \text{ sous } \mathbf{H}_0$$

▶ Score $\hat{\beta}_{H_0} = \begin{pmatrix} \text{Log} \frac{\mathbf{n}_1}{\mathbf{n}_0} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}$

$$\text{SCORE} = U(\hat{\beta}_{H_0})' \hat{I}(\hat{\beta}_{H_0})^{-1} U(\hat{\beta}_{H_0}) \rightarrow \chi^2(k) \text{ sous } H_0$$

3 Influence d'un groupe de variables

Test $H_0 : \beta_{r+1} = \dots = \beta_k = 0$

H_1 : au moins un de ces $\beta_j \neq 0$

Statistiques utilisées

► Rapport des vraisemblances

$$\Lambda = -2 \text{Log} \left[\frac{L(\mathbf{X}_1, \dots, \mathbf{X}_r)}{L(\mathbf{X}_1, \dots, \mathbf{X}_k)} \right] \longrightarrow \chi_{k-r}^2 \text{ sous } H_0$$

► Score

$$U(\hat{\beta}_{H_0})' \hat{I}(\hat{\beta}_{H_0})^{-1} U(\hat{\beta}_{H_0}) \longrightarrow \chi_{k-r}^2 \text{ sous } H_0$$

II SÉLECTION DE VARIABLES : RÉGRESSION LOGISTIQUE PAS À PAS

1 Pas à pas ascendant

À chaque étape, on sélectionne la variable qui aura le «score» le plus élevé une fois introduite dans le modèle.

Soit $\mathbf{X}_1, \dots, \mathbf{X}_t$ les variables déjà introduites. On choisit la variable hors modèle \mathbf{X}_j maximisant le score (\mathbf{X}_j) dans le modèle contenant les variables $\mathbf{X}_1, \dots, \mathbf{X}_t, \mathbf{X}_j$.

L'influence des variables hors modèle $X_{t+1} \dots X_k$ en plus des variables $X_1 \dots X_t$ est testée globalement à l'aide de la statistique Score (notée «Residual chi-square dans SAS).

Si les variables hors modèle X_{t+1}, \dots, X_k sont sans influence marginale sur Y , alors :

$$\text{"Residual chi - square"} \longrightarrow \chi^2_{(k-t)}$$

2. Pas à pas descendant

À chaque étape, on enlève la variable ayant le plus petit «Wald chi-square» à condition qu'il ne soit pas significatif.

3. Meilleur sous-ensemble.

Cette méthode permet de déterminer les meilleurs modèles à une variable explicative, à deux variables explicatives, à trois variables explicatives.....

Méthode basée sur le score de vraisemblance et utilisant l'algorithme de Furnival et Wilson.

Choix parmi les $2^K - 1$ modèles possibles.

Cette option n'est pas disponible en présence de variables qualitatives identifiées par l'instruction CLASS.

III RÉSUMÉ DES TESTS DE VALIDITÉ GÉNÉRALE DU MODÈLE

1 Existe-t-il des statistiques permettant de juger de la bonne adéquation du modèle, en jouant un rôle analogue à celui du R^2 classique ?

→ Déviance

→ Rapport de vraisemblance

→ Statistique du score

→ Critère d' Akaike

$$\text{AIC} = -2 \text{ Log } \mathbf{L} + 2\mathbf{k}$$

(\mathbf{k} = nombre de paramètres à estimer)

→ Critère de Schwartz

$$\text{SC} = -2 \text{ Log } \mathbf{L} + \mathbf{k} \text{ Log } \mathbf{n}$$

(\mathbf{n} = nombre total d'observations)

Les critères de Schwartz et d' Akaike sont utiles pour comparer des modèles différents portant sur les mêmes données (nombre de d.d.l. différents).

On préférera le modèle pour lequel ces statistiques ont la valeur la plus faible.

2 Tables de classement

Objectif Évaluer la capacité prédictive du modèle

$$\begin{cases} \hat{Y}_i = 1 \text{ si } \hat{\Pi}_1 \geq 1/2 \\ \hat{Y}_i = 0 \text{ si } \hat{\Pi}_1 < 1/2 \end{cases}$$

		PRÉDICTION	
		$\hat{Y}_i = 1$	$\hat{Y}_i = 0$
OBSERVATION	$Y_i = 1$	a	b
	$Y_i = 0$	c	d

SENSIBILITE

$$\frac{a}{a + b}$$

SPÉCIFICITÉ

$$\frac{d}{c + d}$$

TAUX D'ERREUR PAR EXCÈS

$$\frac{c}{a + c}$$

TAUX D'ERREUR PAR DÉFAUT

$$\frac{b}{b + d}$$

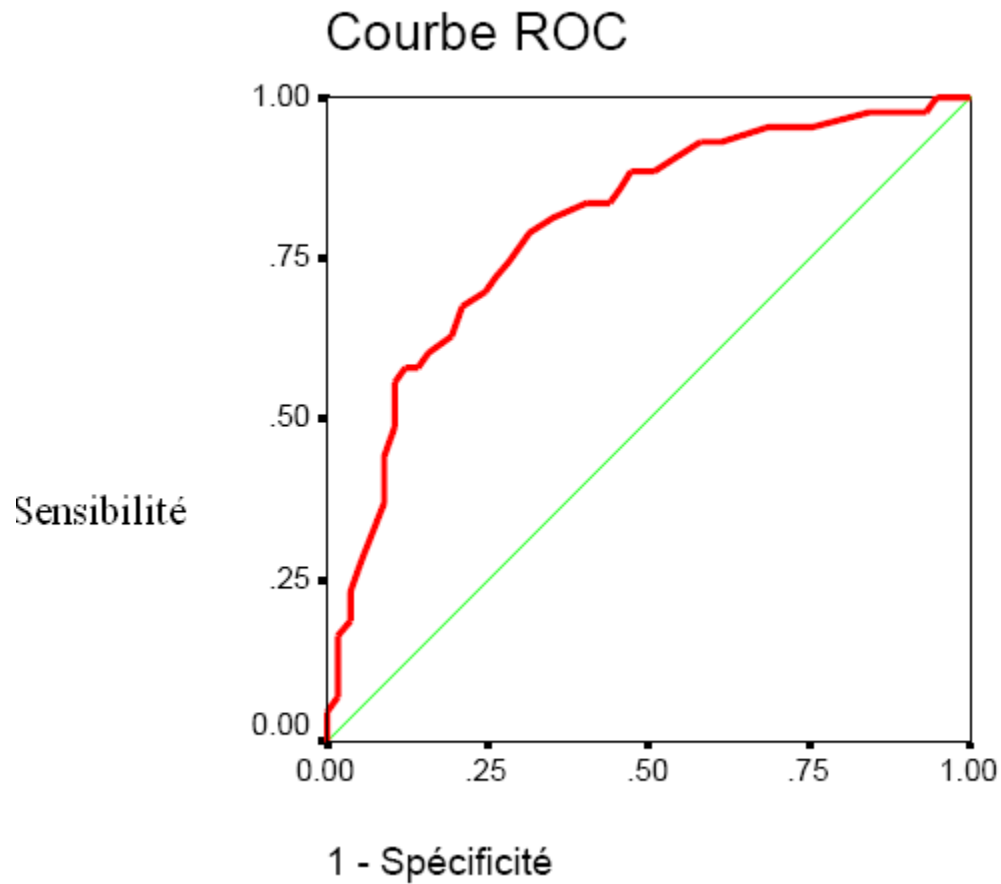
BIEN CLASSÉS

$$\frac{a + d}{n}$$

La capacité prédictive du modèle est d'autant meilleure que :

**{ SENSIBILITE
SPECIFICITÉ
BIEN CLASSÉS**

sont élevés.



Site consacré à la courbe ROC

<http://www.anaesthetist.com/mnm/stats/roc/>

3 Indices de corrélation entre «PRÉDICTION» et «OBSERVATION»

$t = \text{Card} \left\{ (Y_j, Y_k) / Y_j \neq Y_k \right\}$ nombre de paires différentes

$n_c = \text{Card} \left\{ (Y_j = 1, Y_k = 0) , \hat{Y}_j > \hat{Y}_k \right\}$ nombre de paires concordantes

$n_d = \text{Card} \left\{ (Y_j = 1, Y_k = 0) / \hat{Y}_j < \hat{Y}_k \right\}$ nombre de paires discordantes

$t = \text{nombre de concordances} + \text{nombre de discordances} + \text{nombre d'ex aequo}$

1.
$$c = \frac{n_c + 0,5(t - n_c - n_d)}{t} \in [0,1]$$

En l'absence d'ex aequo, c représente l'aire sous la courbe ROC égale à n_c/t souvent noté AUC

2. Somer's
$$D = \frac{n_c - n_d}{t} \in [-1,+1]$$

En l'absence d'ex aequo le D de Somer's est l'indice de Gini: $2AUC-1$

3. Gamma
(Goodman-Kruskal)
$$= \frac{n_c - n_d}{n_c + n_d} \in [-1,+1]$$

4. Kendall's tau-a
$$= \frac{n_c - n_d}{0,5 n(n-1)} \in [-1,+1]$$

La capacité prédictive du modèle est d'autant meilleure que les indices de corrélation sont élevés.

IV INTERPRÉTATION DES COEFFICIENTS DE LA RÉGRESSION LOGISTIQUE

1 Cas d'une variable indépendante binaire

Y Variable dépendante	Variable indépendante X	
	x = 1	x = 0
y = 1	$\Pi(1) = \frac{e^{\beta_0 + \beta_1}}{1 + e^{\beta_0 + \beta_1}}$	$\Pi(0) = \frac{e^{\beta_0}}{1 + e^{\beta_0}}$
y = 0	$1 - \Pi(1) = \frac{1}{1 + e^{\beta_0 + \beta_1}}$	$1 - \Pi(0) = \frac{1}{1 + e^{\beta_0}}$

La cote (**odds**) d'obtenir $Y=1$ pour les personnes ayant $X=1$ est définie par :

$$\frac{\Pi_1}{1 - \Pi_1}$$

La cote (**odds**) d'obtenir $Y=1$ pour les personnes ayant $X=0$ est définie par :

$$\frac{\Pi_0}{1 - \Pi_0}$$

Le rapport de cotes (odds ratio) noté Ψ est défini par :

$$\Psi = \frac{\Pi(1) / (1 - \Pi(1))}{\Pi(0) / (1 - \Pi(0))}$$

En notant : $\mathbf{g}(\mathbf{1}) = \mathbf{Log} \frac{\Pi(\mathbf{1})}{1 - \Pi(\mathbf{1})}$

$$\mathbf{g}(\mathbf{0}) = \mathbf{Log} \frac{\Pi(\mathbf{0})}{1 - \Pi(\mathbf{0})}$$

le logarithme de l'odds ratio ou log odds est :

$$\mathbf{Log} (\psi) = \mathbf{Log} \left[\frac{\Pi(\mathbf{1}) / 1 - \Pi(\mathbf{1})}{\Pi(\mathbf{0}) / 1 - \Pi(\mathbf{0})} \right] = \mathbf{g}(\mathbf{1}) - \mathbf{g}(\mathbf{0})$$

soit la différence de logits.

Le logit $\mathbf{g}(\mathbf{x})$ valant $\beta_0 + \beta_1 \mathbf{x}$, il vient : $\mathbf{g}(\mathbf{1}) - \mathbf{g}(\mathbf{0}) = \beta_1$

$$\text{d'où } \underline{\mathbf{Log}(\psi)} = \beta_1 \Rightarrow \underline{\psi} = e^{\beta_1}$$

Interprétation de l'odds ratio

Lorsque $\Pi(1)$ et $\Pi(0)$ sont petits.

$$\psi = \frac{\Pi(1) / 1 - \Pi(1)}{\Pi(0) / 1 - \Pi(0)} \cong \frac{\Pi(1)}{\Pi(0)} = \eta = \text{risque relatif}$$

On communique souvent les résultats sous cette forme dans un contexte médical.

Exemple

CHD «présence d'une maladie cardiaque»

CHD (y)	≥ 55 (1)	< 55 (0)	TOTAL
Présent (1)	21	22	43
Absent (0)	6	51	57
TOTAL	27	73	100

$$\hat{\psi} = \frac{21/6}{22/51} = 8,11$$

Le quotient de malades par rapport aux non malades est huit fois plus important chez les plus de 55 ans que chez les moins de 55 ans.

Résultats de la régression logistique de Y sur l'âge dichotomisé

Variable	Coefficient estimé	Erreur-type	$\frac{\text{coefficient}}{\text{erreur-type}}$	$\hat{\psi}$
Âge	2,094	0,529	3,96	8,1
Constante	- 0,841	0,255	- 3,30	

$$\hat{\psi} = e^{\beta_1} = e^{2,094} = \underline{8,1}$$

De plus, on obtient un intervalle de confiance de l'odds ratio (ici à 95 %)

$$\exp\left[\hat{\beta}_1 \pm 1,96 s(\hat{\beta}_1)\right] = \exp\left[2,094 \pm 1,96*0,529\right] = [2,9 ; 22,9]$$

Remarque

Les résultats présentés dans cette section dépendent du codage utilisé.

Si on avait codé \mathbf{X} par $\mathbf{1} = (\geq 55)$ et $-\mathbf{1} = (< 55)$

$$\begin{aligned}\text{alors : } \text{Log}(\hat{\psi}) &= \hat{\mathbf{g}}(\mathbf{1}) - \hat{\mathbf{g}}(-\mathbf{1}) \\ &= \hat{\beta}_0 + \hat{\beta}_1 - (\hat{\beta}_0 - \hat{\beta}_1) \\ &= 2 \hat{\beta}_1\end{aligned}$$

$$\text{et } \hat{\psi} = \exp(2 \hat{\beta}_1)$$

2 Variable indépendante polytomique

On souhaite étudier la relation entre la maladie cardiaque et le groupe.

CHD	Blanc	Noir	Hispanique	Autre	TOTAL
Présent	5	20	15	10	50
Absent	20	10	10	10	50
Total	25	30	25	20	100
Odds ratio $\hat{\Psi}$	1	8	6	4	
I.C. 95 % $\hat{\Psi}$		[2,3 ; 27,6]	[1,7 ; 21,3]	[1,1 ; 14,9]	
Log $\hat{\Psi}$	0	2,08	1,79	1,39	

Le groupe «blanc» est utilisé comme groupe de référence.

Le résultats du tableau ci-dessus peuvent être retrouvés à l'aide de la régression logistique de $Y = \text{CHD}$ sur les variables indicatrices des groupes «Noir» «Hispanique» «Autre».

Régression logistique de CHD sur Noir, Hispanique, Autre

Variable	Coefficient estimé	Erreur-type	$\frac{\text{coefficient}}{\text{erreur-type}}$	$\hat{\psi}$
Groupe (1)	2,079	0,633	3,29	8,0
Groupe (2)	1,792	0,646	2,78	6.0
Groupe (3)	1,386	0,671	2,07	4.0
Constante	- 1,386	0,500	- 2,77	

PROGRAMME SAS

```
options ls=75 ps=50 nodate nocenter nopage;
data CHDMULT;
input groupe $ chd $ effectif;
cards;
ref present 5
ref absent 20
gr1 present 20
gr1 absent 10
gr2 present 15
gr2 absent 10
gr3 present 10
gr3 absent 10
;
proc logistic order=data ;
freq effectif ;
class groupe / param= glm;
model chd = groupe / ctable;
run;
```

Estimations par l'analyse du maximum de vraisemblance

Paramètre		DDL	Valeur estimée	Erreur type	Khi-2 de Wald	Pr > Khi-2
Intercept		1	-1.3863	0.5000	7.6871	0.0056
groupe	gr1	1	2.0794	0.6325	10.8100	0.0010
groupe	gr2	1	1.7917	0.6455	7.7048	0.0055
groupe	gr3	1	1.3863	0.6708	4.2706	0.0388
groupe	ref	0	0	.	.	.

Estimations des rapports de cotes			
Effet	Valeur estimée du point	Intervalle de confiance de Wald à 95 %	
groupe gr1 vs ref	8.000	2.316	27.633
groupe gr2 vs ref	6.000	1.693	21.261
groupe gr3 vs ref	4.000	1.074	14.895