

MODELISATION DE DONNÉES QUALITATIVES

PREMIÈRE PARTIE

Pierre-Louis Gonzalez

I INTRODUCTION

■ 1 variable qualitative

- Tri à plat

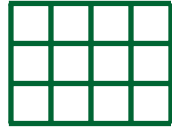
- Représentations graphiques



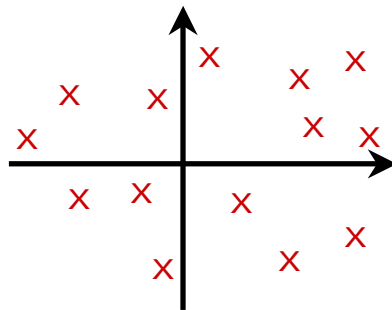
- Modélisation : loi binomiale
loi multinomiale

■ 2 variables qualitatives

- Tri croisé



- Indépendance ?
- Khi-deux ...
- Description du tableau de contingence par analyse des correspondances simples.



■ Plus de deux variables qualitatives

- Tris croisés pour tous les couples de variables (tableau de Burt)
- Analyse des correspondances multiples

But de l'étude ?

- Modélisation
- Expliquer une variable à l'aide d'autres variables ...

▪ Effets de structure

Le recours à l'utilisation de modèles (linéaires, logistiques) est nécessaire pour isoler les effets propres.

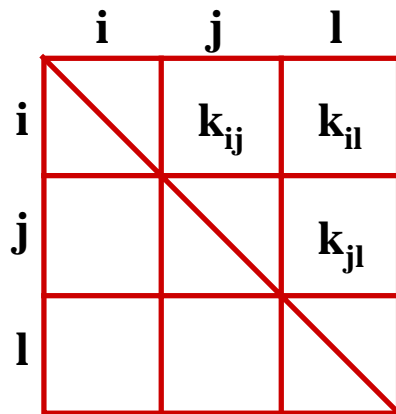
- Séparation des effets
- Effet d'une variable toutes choses égales par ailleurs
- Effet d'une variable conditionnellement aux variables introduites dans le modèle

Exemple Vocations spécifiques de deux approches : description
modélisation

Correspondances multiples (DESCRIPTION)	Modèle log linéaire (EXPLORATION DE L'UNIVERS DES MODÈLES)
Description des liaisons entre les variables prises deux à deux sous forme essentiellement graphique.	Description des interactions entre plus de deux variables dans un cadre inférentiel.
N'impose aucune hypothèse sur les liaisons, mais impose une certaine homogénéité de l'ensemble des variables actives.	Des hypothèses sur les liaisons doivent être formulées au préalable.
N'est pas limitée dans le nombre de variables.	Est limité à peu de variables (en pratique moins de 5).

Correspondances multiples (DESCRIPTION)

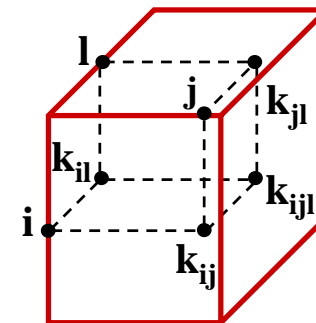
Met seulement en jeu les faces de l'hypercube représentées par le tableau de Burt.



Les individus peuvent jouer un rôle central. L'analyse sert à produire des typologies d'individus.

Modèle log linéaire (EXPLORATION DE L'UNIVERS DES MODELES)

Met en jeu toutes les cases d'un hypercube de contingence.



Les individus n'apparaissent pas.

II LES MÉTHODES EXPLICATIVES

VARIABLE À EXPLIQUER	VARIABLES EXPLICATIVES X_1, \dots, X_K		
Y	Numériques	Nominales	Mixte
Numérique	Régression multiple REG GLM	Analyse de la variance GLM ANOVA	Analyse de la covariance GLM
Qualitative	Analyse discriminante CANDISC STEPPDISC DISCRIM	DISQUAL Analyse discriminante sur variables qualitatives	

VARIABLE À EXPLIQUER	VARIABLES EXPLICATIVES X_1, \dots, X_K		
Y	Numériques	Nominales	Mixte
Nominale à deux modalités	RÉGRESSION LOGISTIQUE LOGISTIC GENMOD		
Nominale	MODÈLE LINÉAIRE GÉNÉRALISÉ LOGISTIC CATMOD GENMOD		
Ordinale	RÉGRESSION LOGISTIQUE ou MODÈLE LINÉAIRE GÉNÉRALISÉ LOGISTIC CATMOD		

III VARIABLE QUALITATIVE À EXPLIQUER

1 Variable dichotomique : $Y \in \{0,1\}$

Exemple 1 Soit P la population des ménages :

$$Y_i = \begin{cases} \mathbf{1} & \text{si le ménage } \mathbf{i} \in P \text{ , possède un bien durable} \\ \mathbf{0} & \text{sinon} \end{cases}$$

$X_i = (\text{AGE, CSP, SALAIRE, HABITAT, ...})$ régresseurs

Exemple 2 Soit P la population des clients potentiels d'une banque :
«CREDIT SCORING»

$$Y_i = \begin{cases} \mathbf{1} & \text{si un crédit est accordé au client } \mathbf{i} \\ \mathbf{0} & \text{sinon} \end{cases}$$

$$X_i = (\text{AGE, REVENU, PRODUIT BANCAIRE, LIEU DE NAISSANCE, ...})$$

Exemple 3 Soit P la population des sujets testés à une dose «DOSAGE LEVEL»

$$Y_i = \begin{cases} \mathbf{1} & \text{si le sujet } \mathbf{i} \in P \text{ réagit au stimulus} \\ \mathbf{0} & \text{sinon} \end{cases}$$

$$\mathbf{X}_i = (\text{NIVEAU DE LA DOSE, POIDS, AGE, ...})$$

La variable réponse à expliquer \mathbf{Y} est une variable de Bernoulli de paramètre \mathbf{p}_i .

$$\mathbf{p}_i = \Pr(\mathbf{Y}_i = \mathbf{1} \mid \mathbf{X}_i) = \mathbf{E}(\mathbf{Y}_i \mid \mathbf{X}_i)$$

$$Y_i \mid X_i \rightarrow B(1, p_i)$$

OBJECTIF

Exprimer \mathbf{p}_i en fonction de \mathbf{X}_i

2 Variable polytomique

→ Polytomique ordonnée

Exemple 1 Soit P la population d'étudiants :

$$Y_i = \begin{cases} 1 & \text{si l'étudiant } i \in P \text{ pratique du sport tous les jours} \\ 2 & \text{si l'étudiant } i \in P \text{ pratique du sport une ou plusieurs fois par semaine} \\ 3 & \text{si l'étudiant } i \in P \text{ pratique du sport plus rarement} \end{cases}$$

La variable réponse Y : «pratique du sport» est codée

$$X_i = (\text{AGE, SEXE, TYPE D'ETUDES, ...})$$

Exemple 2 Soit P la population de chômeurs à la date t :

$$Y_i = \begin{cases} \mathbf{1} & \text{si l'individu } \mathbf{i} \in P \text{ est toujours au chômage à la date } \mathbf{t} + \delta \\ \mathbf{2} & \text{si l'individu } \mathbf{i} \in P \text{ est en formation (stage)} \\ \mathbf{3} & \text{si l'individu } \mathbf{i} \in P \text{ a un contrat CDD} \\ \mathbf{4} & \text{si l'individu } \mathbf{i} \in P \text{ a un contrat CDI} \end{cases}$$

$$X_i = (\text{AGE, SEXE, DIPLÔME, QUALIFICATION ...})$$

→ Polytomique non ordonnée

Y_i «distraction du samedi soir»

$$Y_i = \begin{cases} 1 = \text{télévision} \\ 2 = \text{théâtre} \\ 3 = \text{cinéma} \\ 4 = \text{visite amis} \end{cases}$$

$X_i = (\text{AGE, SEXE, CSP, HABITAT, ...})$

OBJECTIF

Exprimer $p_{ij} = P(Y_i = j \mid X_i)$

en fonction de X_i pour $j = 1, 2, 3 \dots$

IV POURQUOI DES MODÈLES PARTICULIERS ?

1 Cas de la régression linéaire classique

$$Y_i = x_i \beta + \varepsilon_i$$

↑
variable aléatoire quantitative

↑
prédicteur linéaire (élément déterminé)

↑
variable aléatoire

1

- $\mathbf{E}(\boldsymbol{\varepsilon}_i \mid \mathbf{X}_i = \mathbf{x}_i) = \mathbf{0}$

Par la suite, on notera les espérances sans conditionnement $\mathbf{X}_i = \mathbf{x}_i$ ce qui revient à considérer \mathbf{X}_i est non aléatoire. On notera indifféremment \mathbf{X}_i ou \mathbf{x}_i .

- $\mathbf{V}(\boldsymbol{\varepsilon}_i) = \sigma^2$

Si de plus $\boldsymbol{\varepsilon}_i$ est supposée gaussienne, l'estimateur des moindres carrés ordinaire :

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$$

est l'estimateur du maximum de vraisemblance.

2 Cas de la régression d'une variable dichotomique

Si on modélise par $\diamond 1$, on obtient un résidu qui est une v.a.r. discrète prenant deux valeurs :

$$\diamond 1 \Rightarrow \begin{cases} \varepsilon_i = 1 - \mathbf{x}_i\beta & \text{avec la probabilité } \mathbf{p}_i = \mathbf{P}(Y_i = 1) \\ \varepsilon_i = -\mathbf{x}_i\beta & \text{avec la probabilité } 1 - \mathbf{p}_i \end{cases}$$

Si on modélise par $\diamond 1$ l'estimateur $\hat{\beta}$ n'est plus efficace.

$$\begin{array}{l} \diamond 1 \Rightarrow \mathbf{E}(Y_i) = \mathbf{x}_i\beta \\ \text{Or } Y_i \rightarrow B(1, p_i) \Rightarrow \mathbf{E}(Y_i) = p_i \end{array} \left. \vphantom{\begin{array}{l} \diamond 1 \Rightarrow \mathbf{E}(Y_i) = \mathbf{x}_i\beta \\ \text{Or } Y_i \rightarrow B(1, p_i) \Rightarrow \mathbf{E}(Y_i) = p_i \end{array}} \right\} \Rightarrow \mathbf{p}_i = \mathbf{x}_i\beta$$

T

une valeur qui n'est pas forcément entre 0 et 1

Le modèle $\diamond 1$ est donc inapproprié !

V NIVEAU D'UTILITÉ, VARIABLE LATENTE

1 Cas de variable latente

Z_i «intensité du désir de posséder le bien»
↓
pour le ménage i caractérisé par \mathbf{X}_i
 Z_i non observable

$$\begin{cases} Y_i = 0 & \Leftrightarrow Z_i < s \rightarrow (\text{seuil théorique}) \\ Y_i = 1 & \Leftrightarrow Z_i \geq s \end{cases}$$

c'est-à-dire $Y_i = \mathbf{1}_{(Z_i \geq s)}$

2 Fonction d'utilité

Soit $u(\mathbf{1}, \mathbf{x}_i)$ le niveau d'utilité procuré par la possession du bien

$u(\mathbf{0}, \mathbf{x}_i)$ le niveau d'utilité procuré par la non possession du bien

$$\begin{cases} Y_i = 0 & \Leftrightarrow u(\mathbf{0}, \mathbf{x}_i) > u(\mathbf{1}, \mathbf{x}_i) \\ Y_i = 1 & \Leftrightarrow u(\mathbf{1}, \mathbf{x}_i) \geq u(\mathbf{0}, \mathbf{x}_i) \end{cases}$$

c'est-à-dire :

$$Z_i = u(\mathbf{1}, \mathbf{x}_i) - u(\mathbf{0}, \mathbf{x}_i)$$

$$Y_i = \mathbf{1}_{(Z_i \geq 0)}$$



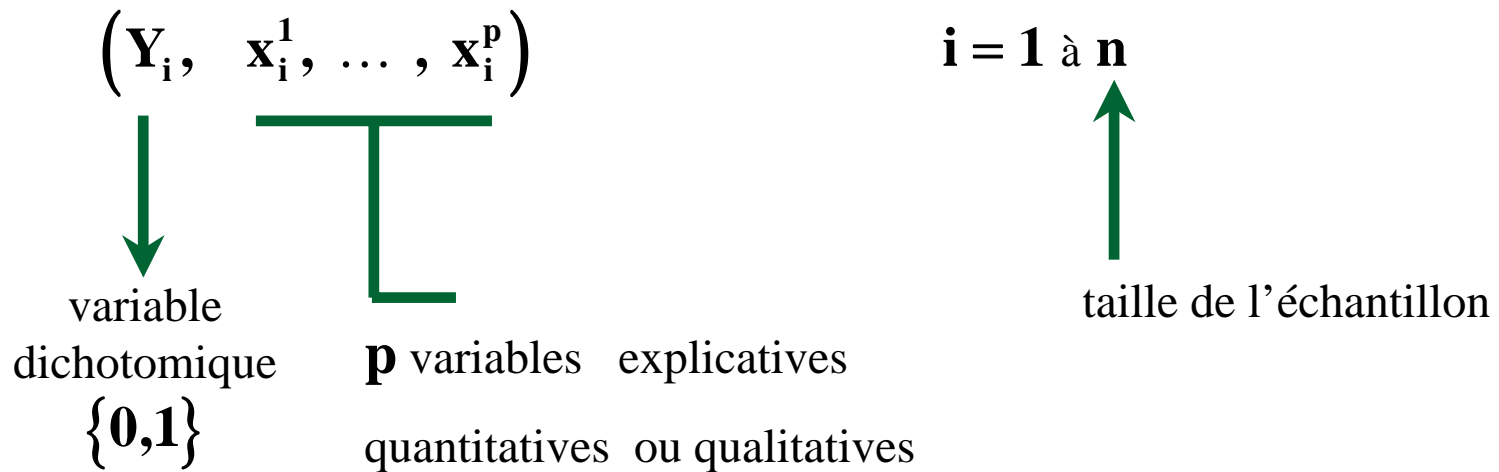
Dans ces deux cas, on peut exprimer la probabilité

$$\mathbf{p}_i = \mathbf{P}(\mathbf{Y}_i = \mathbf{1} \mid \mathbf{x}_i) \text{ comme :}$$

$$\mathbf{p}_i = \mathbf{P}(\mathbf{Z}_i \geq \mathbf{s})$$

VI MODÈLE THÉORIQUE

1 Données statistiques



\mathbf{X}_i vecteur de \mathbb{R}^p

(On supposera $(x_i^1 = 1 \forall i)$)

de façon à définir un modèle avec constante)

En introduisant \mathbf{Z} variable latente non observable telle que :

$$\left. \begin{array}{l} \mathbf{Z}_i = \mathbf{X}_i\boldsymbol{\beta} + \boldsymbol{\varepsilon}_i \\ \mathbf{Y}_i = \mathbf{1}_{\mathbf{Z}_i > \mathbf{0}} \end{array} \right\} \mathbf{p}_i = \mathbf{P}(\mathbf{Y}_i = 1) = \mathbf{P}(-\boldsymbol{\varepsilon}_i < \mathbf{X}_i\boldsymbol{\beta})$$

$$= \mathbf{F}(\mathbf{X}_i\boldsymbol{\beta})$$



fonction de répartition
de $-\boldsymbol{\varepsilon}_i$

2 Modèle stochastique général

$(Y_i, X_i)_{i=1..n}$ i.i.d tel que :

$$H_1 : Y_i | X_i \rightarrow B(1, p_i)$$

$$H_2 : p_i = P(Y_i = 1 | X_i) = F(X_i \beta)$$

où $F : \mathbf{R} \rightarrow [0,1]$ fonction de répartition

Le paramètre β , vecteur de \mathbf{R}^p formé des coefficients de régression est inconnu.

3 Modèles PROBIT, LOGIT, ...

L'hypothèse \mathbf{H}_2 dépend du choix de la fonction \mathbf{F} . Les modèles paramétriques usuels sont :

3.1 Le modèle probit

$$F(\omega) = \Phi(\omega) = \int_{-\infty}^{\omega} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right) dt \quad \forall \omega \in \mathbb{R}$$

Fonction de répartition de la loi normale centrée réduite $N(0;1)$

$$\mathbf{F}^{-1} = \Phi^{-1} \quad \text{probit}$$

3.2 Le modèle logit

$$F(\omega) = \frac{e^{\omega}}{1 + e^{\omega}} = \frac{1}{1 + e^{-\omega}} \quad \forall \omega \in \mathbb{R}$$

Fonction de répartition de la loi logistique
de moyenne **0** et de variance $\frac{\pi^2}{3}$

$$F^{-1}(t) = \ln \frac{t}{1-t} \quad \underline{\text{logit}}$$

3.3 Le modèle complémentaire log-log (ou modèle Gompit)

$$F(\omega) = 1 - \exp(-e^\omega) \quad \forall \omega \in \mathbb{R}$$

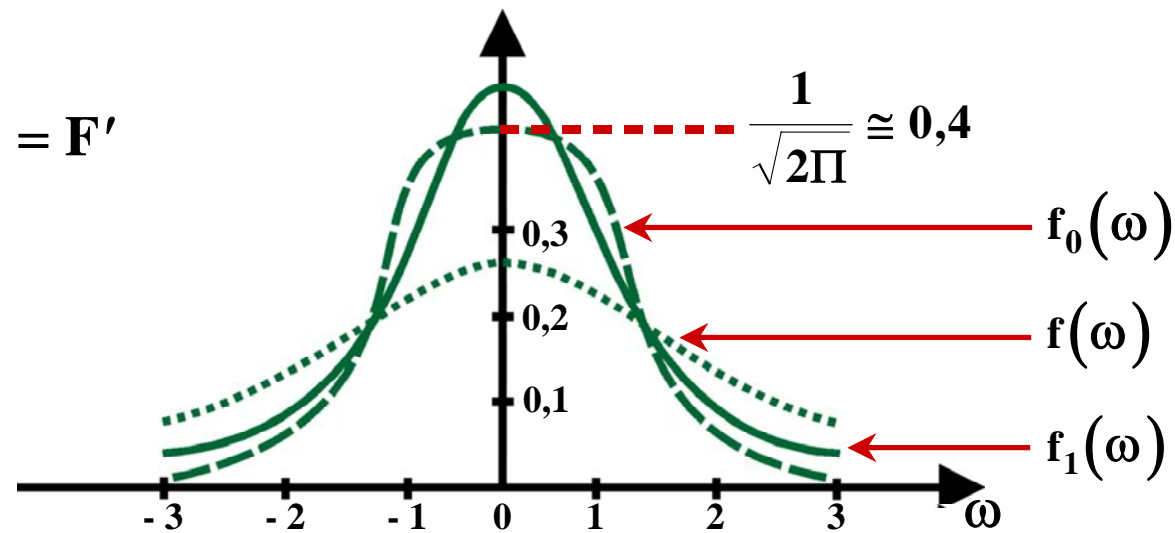
Fonction de répartition de la loi de Gompertz
de moyenne 0,577 (constante d'Euler) et de
variance $\frac{\pi^2}{6}$

$$F^{-1}(t) = \ln(-\ln(1-t))$$

Remarque Cette loi est dissymétrique.

4 Comparaison des modèles LOGIT et PROBIT

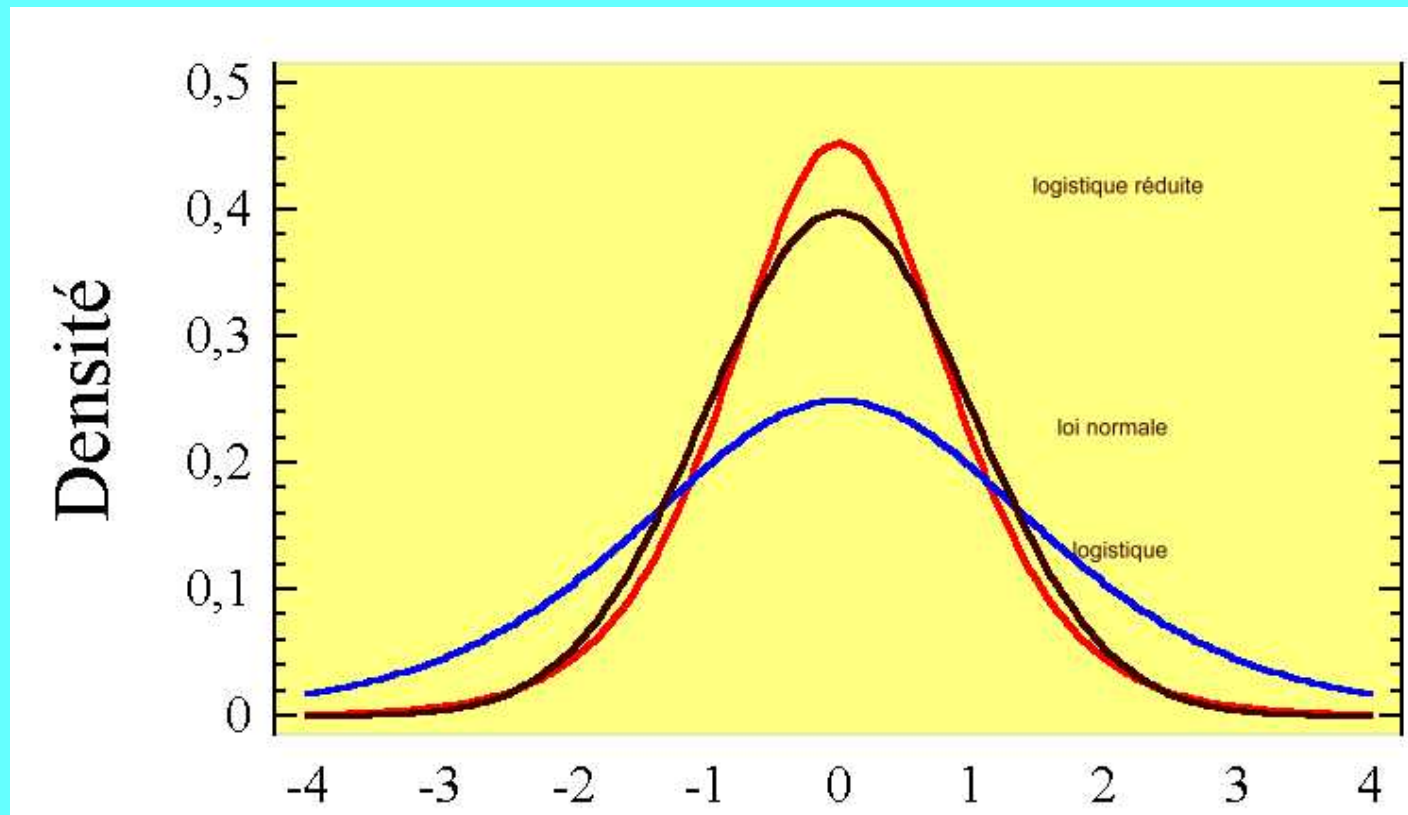
DENSITÉS $f = F'$



Modèle PROBIT $\Phi(\omega) = \int_{-\infty}^{\omega} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \Rightarrow f_0(\omega) = \frac{1}{\sqrt{2\pi}} e^{-\omega^2/2}$

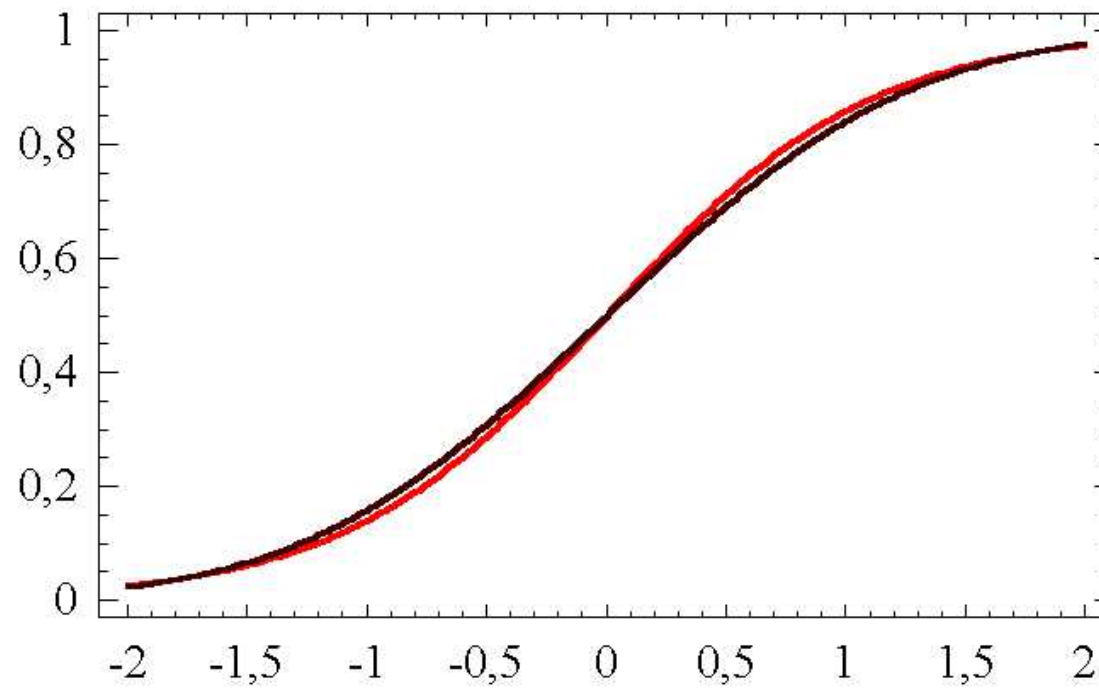
Modèle LOGIT $F(\omega) = \frac{1}{1 + e^{-\omega}} \Rightarrow f(\omega) = \frac{e^{\omega}}{(1 + e^{\omega})^2}$

Modèle LOGIT réduit $F_1(\omega) = \frac{1}{1 + e^{-\Pi\omega/\sqrt{3}}} \quad f_1(\omega) = \frac{\Pi}{\sqrt{3}} \frac{e^{\Pi\omega/\sqrt{3}}}{(1 + e^{\Pi\omega/\sqrt{3}})^2}$

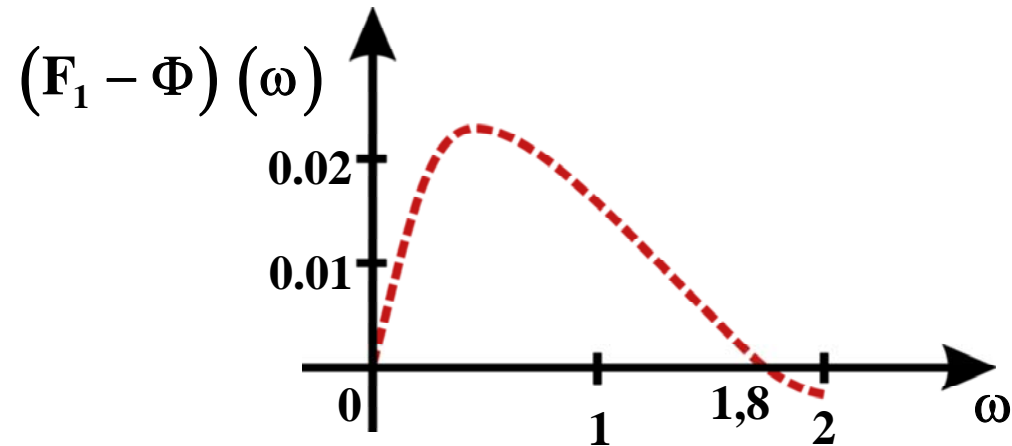


Sur les extrêmes la loi logistique s'approche un peu plus lentement de 0 ou de 1.

Comparaison des fonctions de répartition



DIFFÉRENCE



CONCLUSION

- Les lois F_1 et Φ sont proches \Rightarrow modèles équivalents.
- La procédure de SAS, «LOGISTIC» utilise F ou Φ .
 \Rightarrow Les estimateurs obtenus avec F (logit) seront $\Pi / \sqrt{3}$ fois plus grands qu'avec Φ (probit).
- Le modèle LOGIT est préférable car les calculs sont plus simples.

Dans la plupart des cas pratiques, on peut donc choisir indifféremment l'un ou l'autre modèle.

Le **modèle LOGIT** a l'avantage d'une plus grande simplicité numérique.

Le **modèle PROBIT** est en revanche plus proche du modèle habituel de régression par les moindres carrés.

Avantages du modèle LOGIT

Les coefficients du modèle LOGIT sont interprétables en termes d'odds-ratio.

Un échantillonnage ne respectant pas les proportions réelles dans la population des deux modalités de la variable à expliquer Y ne change que la constante dans le modèle.

VII PRINCIPES GÉNÉRAUX : ESTIMATION DU VECTEUR β

1 La méthode du maximum de vraisemblance

Échantillon $\mathbf{x}_1 \dots \mathbf{x}_n$

Loi de probabilité de \mathbf{x}_i $f(\mathbf{x}_i, \theta)$ où $\theta = (\theta_1, \dots, \theta_k)'$ $\in \Omega$

Vraisemblance

$$L(\theta) = \prod_{i=1}^n f(\mathbf{x}_i, \theta) \quad \text{dépend des } \mathbf{x}_i \text{ et des } \theta$$

Estimation du maximum de vraisemblance

$$\hat{\theta} \quad L(\hat{\theta}) = \text{MAX}_{\theta \in \Omega} L(\theta)$$

On obtient en général $\hat{\theta}$ en annulant les dérivées premières $\frac{\partial \text{Log } L(\theta)}{\partial \theta_i}$

Scores

$$\mathbf{u}_i(\boldsymbol{\theta}) = \frac{\partial \text{Log } \mathbf{L}(\boldsymbol{\theta})}{\partial \theta_i}$$

$$\mathbf{u}(\boldsymbol{\theta}) = (\mathbf{u}_1(\boldsymbol{\theta}), \dots, \mathbf{u}_k(\boldsymbol{\theta}))' = \text{vecteur score}$$

$$\text{On a : } \mathbf{u}(\hat{\boldsymbol{\theta}}) = \mathbf{0}$$

Matrice d'information de Fisher

$$\mathbf{I}(\boldsymbol{\theta}) = \mathbf{E} \left[\frac{-\partial^2 \text{Log } \mathbf{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^2} \right]$$

$$\text{estimée par : } \hat{\mathbf{I}}(\hat{\boldsymbol{\theta}}) = \left(\frac{-\partial^2 \text{Log } \mathbf{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^2} \right)_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}$$

Résultats

θ = vecteur des paramètres

$\hat{\theta}$ = estimation du maximum de vraisemblance

1. $\hat{\theta} \rightarrow N \left(\theta ; I(\theta)^{-1} \right)$
2. $u(\theta) \rightarrow N(0 ; I(\theta))$
3. $(\hat{\theta} - \theta)' I(\theta) (\hat{\theta} - \theta) \rightarrow \chi_{(k)}^2$
4. $u(\theta)' I(\theta)^{-1} u(\theta) \rightarrow \chi_{(k)}^2$
5. $\Lambda = -2 \text{Log} \frac{L(\theta)}{L(\hat{\theta})} \rightarrow \chi_{(k)}^2$

2 Test global

$$\mathbf{H}_0 : \theta = \theta_0$$

Statistiques

1. Statistique de Wald

$$(\hat{\theta} - \theta_0)' \mathbf{I}(\theta_0) (\hat{\theta} - \theta_0) \rightarrow \chi_{(k)}^2 \text{ sous } \mathbf{H}_0$$

2. Statistique du score

$$\mathbf{u}(\theta_0)' \mathbf{I}(\theta_0)^{-1} \mathbf{u}(\theta_0) \rightarrow \chi_{(k)}^2 \text{ sous } \mathbf{H}_0 \text{ Avantage : pas de calcul de } \hat{\theta}$$

3. Statistique des vraisemblances

$$\Lambda = -2 \text{Log} \frac{\mathbf{L}(\theta_0)}{\mathbf{L}(\hat{\theta})} \rightarrow \chi_{(k)}^2 \text{ sous } \mathbf{H}_0$$

3 Test partiel

$$\theta = (\theta_1, \theta_2) \quad \theta_1 \text{ a } \mathbf{p} \text{ coordonnées}$$

$$\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2) = \text{estimation du M. V.}$$

$$\text{Test } \mathbf{H}_0 : \theta_1 = \theta_{10}$$

$$\text{On calcule } \hat{\theta}_{\mathbf{H}_0} = (\theta_{10}, \tilde{\theta}_2)$$

$$\text{avec } \mathbf{L}(\theta_{10}, \tilde{\theta}_2) = \max_{\theta_2} \mathbf{L}(\theta_{10}, \theta_2)$$

$$\hat{\theta}_{\mathbf{H}_0} = \text{estimation de } \theta \text{ sous } \mathbf{H}_0$$

Statistiques utilisées

1. Wald $(\hat{\theta}_1 - \theta_{10})' \text{Var}(\hat{\theta}_1)^{-1} (\hat{\theta}_1 - \theta_{10}) \rightarrow \chi_{(p)}^2$ sous \mathbf{H}_0

$$\text{Var}(\hat{\theta}_1) \text{ est extrait de } \text{Var}(\hat{\theta}) = \hat{\mathbf{I}}(\hat{\theta})^{-1}$$

2. Score

$$\mathbf{u}(\hat{\theta}_{\mathbf{H}_0})' \hat{\mathbf{I}}(\hat{\theta}_{\mathbf{H}_0})^{-1} \mathbf{u}(\hat{\theta}_{\mathbf{H}_0}) \rightarrow \chi_{(p)}^2 \text{ sous } \mathbf{H}_0$$

3. Rapport de vraisemblance

$$\Lambda = -2 \text{Log} \frac{\mathbf{L}(\theta_{10}, \tilde{\theta}_2)}{\mathbf{L}(\hat{\theta}_1, \hat{\theta}_2)} \rightarrow \chi_{(p)}^2 \text{ sous } \mathbf{H}_0$$