

## La régression logistique

Frédéric Bertrand et Myriam Maumy<sup>1</sup>

<sup>1</sup>IRMA, Université Louis Pasteur  
Strasbourg, France

Master 2ème Année 21-11-2007

Frédéric Bertrand et Myriam Maumy	La régression logistique	La régression logistique
Introduction		
Régression logistique : variable explicative qualitative		
Régression logistique : variable explicative continue		
Régression logistique : variables explicatives mixtes		

Frédéric Bertrand et Myriam Maumy	La régression logistique	La régression logistique
Introduction		
Régression logistique : variable explicative qualitative		
Régression logistique : variable explicative continue		
Régression logistique : variables explicatives mixtes		

- Pour **tester** l'**existence** de ce lien il serait possible de procéder à un **test du Khi-deux** (étudié en L3) :

Les dénombresments attendus sont imprimés sous les dénombresments observés

Succès Echec Total

1 21 2 23

16,73 6,27

2 19 13 32  
23,27 8,73

**Question : Existe-t-il une corrélation entre le développement de la maladie et l'apparition du cancer ?**

Total 40 15 55  
Khi deux = 1,091 + 2,910 + 0,784 + 2,092 = 6,878

Df = 1, P = 0,009

Groupe	Tumeur présente	Tumeur absente	Total
Contrôle	19	13	32
Traitement	21	2	23

Ce test ne permet pas de déterminer la **nature** de ce lien, c'est-à-dire comment sont liées les variations des deux variables.

- **Pour parer à cet inconvenient :** On utilise *la régression logistique* qui permet de **modéliser** la probabilité de succès à l'aide des variables explicatives dont nous disposons. Ceci nous permettra de tester si ces changements sont significatifs à un niveau  $\alpha$  donné.

De même que la régression linéaire (simple ou multiple) est un prolongement de l'étude du coefficient de corrélation linéaire de deux variables quantitatives, de même la régression logistique est une généralisation d'un coefficient servant à évaluer la corrélation de deux variables qualitatives : *le rapport des cotes* ou *odds-ratio*.

### Définition

On appelle **côte du succès** *le rapport*

$$\exp(\theta) = \frac{\pi}{1 - \pi}$$

où  $\pi$  est *la probabilité de succès*.

### Définition

*La probabilité de succès s'exprime à partir de la côte de succès de la manière suivante :*

$$\pi = \frac{\exp(\theta)}{1 + \exp(\theta)}.$$

Pour fixer les idées voici quelques valeurs de la côte du succès en fonction la probabilité de succès. (Le logarithme de) cette côte :

- est ( $< 0$ )  $< 1$  lorsque  $\pi < 0.5$ .
- est ( $= 0$ )  $= 1$  lorsque  $\pi = 0.5$ .
- est ( $> 0$ )  $> 1$  lorsque  $\pi > 0.5$ .
- ( $\rightarrow -\infty$ )  $\rightarrow 0$  lorsque  $\pi \rightarrow 0$ .
- ( $\rightarrow +\infty$ )  $\rightarrow +\infty$  lorsque  $\pi \rightarrow 1$ .

$$\hat{\pi} = \frac{40}{55} = 0.73$$

$$\hat{\exp(\theta)} = \frac{\hat{\pi}}{1 - \hat{\pi}} = \frac{0.73}{0.27} = 2.67$$

$$\hat{\theta} = \ln(2.67) = 0.98.$$

## Le logarithme du rapport de côtés :

- On peut calculer la côte de succès dans différentes conditions.

## Définition

*Le rapport de cotés  $\psi$  permet alors d'évaluer l'influence du facteur considéré :*

$$\Psi = \frac{\exp(\theta_2)}{\exp(\theta_1)} = \exp(\theta_2 - \theta_1).$$

- Lorsque  $\Psi$  est  $> 1$  ( $< 1$ ) le succès a une côte supérieure (inférieure) pour le deuxième niveau du facteur.
  - Le logarithme du rapport de côtes,  $\theta_2 - \theta_1$ , est  $> 0$  ( $< 0$ ) lorsque le succès a une probabilité supérieure (inférieure) pour le deuxième niveau du facteur.

Frédéric Bertrand et Myriam Mauim

<b>Introduction</b>	Régression logistique : variable explicative qualitative
<b>Exemple</b>	Régression logistique : variable explicative continue
<b>Rapport des côtés</b>	Rapport des cotés
<b>Intervalle de confiance</b>	Intervalle de confiance

# Intervalle de confiance

- Si pour chaque individu, la probabilité de succès est  $\pi$ , alors le nombre  $Y$  de succès parmi  $n$  individus indépendants suit une loi binomiale  $B(n, \pi)$ . Ainsi :

$$\mathbb{E}\left[\hat{\pi} = \frac{Y}{n}\right] = \frac{1}{n}\mathbb{E}[Y] = \pi \quad ; \quad \text{Var}\left[\hat{\pi}\right] = \frac{1}{n^2}\text{Var}[Y] = \frac{\pi(1-\pi)}{n}.$$

- Un intervalle de confiance (dans le cadre d'application de l'approximation de la loi binomiale par une loi normale) à 95 % pour  $\pi$  est donné par :

$$\hat{\pi} \pm 1.96 \times \sqrt{\frac{\hat{\pi}(1 - \hat{\pi})}{n}}.$$

## Exemple

- La côte du succès (= « développer une tumeur ») observée est égale à :*

$$\left\{ \begin{array}{l} Côte(\text{succès}/\text{Treatment}) = \exp(\hat{\theta}_2) = \frac{21}{2} = 10.5 \\ Côte(\text{succès}/\text{Contrôle}) = \exp(\hat{\theta}_1) = \frac{19}{13} = 1.46. \end{array} \right.$$

$$D'où \quad \hat{\Psi} = \frac{21 \cdot 13}{2 \cdot 19} = 7.18 > 1$$

et  $\ln(\hat{\Psi}) = \hat{\theta}_2 - \hat{\theta}_1 = 1.97 > 0$ .

ès de la tâche est supérieure (multipliée par

La cote de succès de la tumeur est supérieure (multipliée par 7) lorsque les souris sont exposées à la fumée de cigarettes.

Frédéric Bertrand et Myriam Maumy | a régression logistique

<b>Introduction</b>	Régression logistique : variable explicative qualitative
<b>Exemple</b>	Régression logistique : variable explicative continue
	Rapport des côtes
	Intervalle de confiance

- Dans notre exemple on souhaiterait comparer les probabilités  $\pi_1$  et  $\pi_2$  de développer une tumeur sous et sans exposition à la fumée de cigarettes et déterminer si elles sont significativement différentes. Cela reviendrait à déterminer s'il existe un lien entre le développement de la tumeur et le facteur risque considéré.

$$(\hat{\pi}_1 - \hat{\pi}_2) \pm 1.96 \times \sqrt{\frac{\hat{\pi}_1(1 - \hat{\pi}_1)}{n_1} + \frac{\hat{\pi}_2(1 - \hat{\pi}_2)}{n_2}}.$$

## Example

$$0 \notin (0.114, 0.524)$$

On en déduit que la différence  $\pi_1 - \pi_2$  est significativement écartée de 0 au seuil  $\alpha = 5\%$ .

Ainsi on sait non seulement que la fumée de cigarettes a un effet significatif sur le nombre de cancers développés mais surtout on a quantifié cet effet.

Frédéric Bertrand et Myriam Maumy	Régression logistique : variable explicative qualitative	La régression logistique
	Introduction	Définitions
	Introduction	Les modèles
	Régression logistique - variable explicative continue	Déviance
	Régression logistique : variables explicatives mixtes	

Frédéric Bertrand et Myriam Matmy  
Régression logistique : variable explicative qualitative  
Regression logistique : variable explicative continue  
Régression logistique : variables explicatives mixtes

La régression logistique  
Définitions  
Les modèles  
Déviance

Remarque

Dans des situations plus complexes, à savoir par exemple dans des cas où il y a plus que deux variables qualitatives ou plus que deux niveaux du facteur qui est joué par la variable qualitative (on rappelle que l'on parle de facteur lorsque l'on a à faire à des variables qualitatives (cf l'ANOVA)), l'approche précédente est trop lourde.

⇒ On travaille alors avec les côtés de succès que nous allons définir.

Définition

*Le logarithme de la côte de succès sous le premier niveau du facteur vaut  $\mu$ .*

Définition

*Le logarithme du rapport des cotés du succès sous les kème et 1er niveau du facteur vaut  $\theta_k - \theta_1 = \alpha_k$ .*

### Remarque

Par conséquent une valeur de  $\alpha_k > 0$  ( $< 0$ ) indique que la côte du succès observée est plus grande (petite) sous le  $k^{\text{ème}}$  niveau du facteur que sous le 1<sup>er</sup> niveau du facteur.

$$\Rightarrow \pi_k = \frac{\exp(\mu + \alpha_k)}{1 + \exp(\mu + \alpha_k)}.$$

avec

$$\text{logit}(\pi_k) = \ln \left( \frac{\pi_k}{1 - \pi_k} \right) = \theta_k = \mu + \alpha_k; (\alpha_1 = 0)$$

Si  $X$  est une variable explicative à  $K$  niveaux, le modèle logistique suppose que :

*Le logarithme de la côte de succès sous le premier niveau du facteur vaut  $\mu$ .*

## Estimation des $\alpha_k$

- On estime les  $\alpha_k$  à l'aide d'une méthode statistique appelée méthode du maximum de vraisemblance.
- Dans ce cas, on sait qu'asymptotiquement (lorsque la taille de l'échantillon tend vers l'infini) les estimateurs des  $\alpha_k$  suivent une loi normale de moyenne  $\alpha_k$  et de variance  $\text{Var}[\hat{\alpha}_k]$ .
- De plus, ces estimateurs sont sans biais.

Par conséquent un intervalle de confiance à 95 % approximatif pour les  $\alpha_k$  est donné par :

$$\hat{\alpha}_k \pm 1.96 \times \sigma(\hat{\alpha}_k).$$

## Les différents modèles possibles pour l'exemple sont :

- Modèle 1 avec « effet du traitement » :
- Modèle 2 sans « effet du traitement » ( $\alpha_2 = 0$  ci-dessus) :

$$\text{logit}(\pi_k) = \theta_k = \mu + \alpha_k \quad \text{où } k = 1 \text{ ou } 2.$$

On compare alors la probabilité de succès estimée dans le groupe  $k$ , notée  $\tilde{\pi}_k$  et la proportion de succès observée notée  $\hat{\pi}_k$ .

### Définition

La déviance  $D$  est alors définie ainsi :

$$\begin{aligned} D &= -2 \sum_k \left\{ y_k \ln \left( \frac{\tilde{\pi}_k}{\hat{\pi}_k} \right) + (n_k - y_k) \ln \left( \frac{1 - \tilde{\pi}_k}{1 - \hat{\pi}_k} \right) \right\} \\ &= -2(I(\tilde{\pi}_k) - I(\hat{\pi}_k)). \end{aligned}$$

Cette quantité est à rapprocher de la somme des carrés à minimiser dans la régression linéaire simple ou multiple. Elle évalue globalement la qualité de l'ajustement obtenu.

Le deuxième modèle ne fait pas intervenir de variable explicative. Il peut servir à tester la nullité de toutes les pentes : l'équivalent du test de Fisher global dans le cadre de la régression logistique.

On calcule la statistique  $G^2 = D_2 - D_1 = -2(l_2 - l_1)$  comparant la déviance des deux modèles.

### Définition

Sous l'hypothèse nulle  $H_0$  que les restrictions impliquées par le modèle 2 au modèle 1 sont correctes,

$$G \underset{H_0}{\sim} \chi^2_{ddl_2 - ddl_1}.$$

### Exemple

**Sous l'hypothèse nulle**

$$H_0 : \alpha_2 = 0$$

on a

$$G_2 = 7.635, \quad dd|_1 = 0, dd|_2 = 1, \quad \text{et} \quad p = 0.006.$$

Ce qui permet de décider que  $\alpha_2$  est significativement différent de 0 au niveau  $\alpha = 5\%$ . On obtient également les informations suivantes :  $\hat{\mu} = 0.38$  et  $\hat{\alpha}_2 = 1.97$ . Ceci permet de calculer les probabilités de succès : 0.59 et 0.91. Le rapport des cotes du groupe exposé contre le groupe de contrôle est estimé par  $\exp(\hat{\alpha}_2) = 7.24$  soit une côte de succès plus de 7 fois plus grande pour le groupe des traités.

### Exemple

Voici un second exemple que l'on va traiter avec Minitab.  
Relation entre les habitudes tabagiques d'étudiants en Arizona et les habitudes de leurs parents (Agresti, 1990, p. 124).

Nombre de parents fumeurs	Enfant fumeur	Enfant non fumeur	Total
Deux	400	1380	1780
Un seul	416	1823	2239
Aucun	188	1168	1358

On peut construire un intervalle de confiance (approximatif)  $(1 - \alpha) \cdot 100\%$  pour le logarithme du rapport de cotes (abrégué en LRC) du groupe  $k$  contre le groupe de référence  $\alpha_k$  avec

$$\hat{\alpha}_k \pm 1.96 \times \sigma(\hat{\alpha}_k).$$

### Exemple

Dans notre exemple, on obtient :  $\alpha_2 \in (0.36; 3, 58)$  confirmant le rejet de l'hypothèse nulle  $H_0$  (avec  $\alpha = 5\%$ ) et l'augmentation significative de développer un cancer du poumon après exposition à la fumée de cigarettes. L'intervalle de confiance approximatif pour le rapport de côte est alors égal à (1.43, 36.0).

On définit le succès comme étant le fait de fumer pour l'enfant, le modèle logistique précédent devient :

$$\text{logit}(\pi_k) = \theta_k = \mu + \alpha_k; (\alpha_1 = 0).$$

La catégorie de référence est par défaut "Aucun". On utilise

Minitab pour mener à bien l'analyse. On peut tester l'hypothèse null

$$H_0 : \alpha_2 = \alpha_3 = 0$$

en comparant la déviance de ce modèle avec celle du précédent.  $G^2_{obs} = 38.37$  d'où une  $p$ -valeur de 0.000.

**Conclusion du test :** Association significative au niveau  $\alpha = 5\%$  entre habitudes tabagiques des parents et des enfants.

## Exemple

*Effet de la cyperméthrine à différentes doses (en  $\mu\text{g}$ ) sur la survie de parasites. Pour chaque niveau de dose, 20 parasites sont exposés. La survie éventuelle de l'animal est évaluée après 72 heures. Les animaux peuvent être distingués par leur sexe (Collett, 1991, CRC, p. 75).*

Dose Mâle	N morts	Dose Femelle	N morts
1	1	1	0
2	4	2	2
4	9	4	6
8	13	8	10
16	18	16	12
32	20	32	16

## Variabile explicative continue

*Effet de la cyperméthrine à différentes doses (en  $\mu\text{g}$ ) sur la survie de parasites. Pour chaque niveau de dose, 20 parasites sont exposés. La survie éventuelle de l'animal est évaluée après 72 heures. Les animaux peuvent être distingués par leur sexe (Collett, 1991, CRC, p. 75).*

Ignorons le sexe de l'animal en premier lieu.

**Question :** Existe-t-il un lien entre la mort d'une larve et la dose reçue ? Si oui quelle est la nature de cette relation ?

Frédéric Bertrand et Myriam Maumy

Introduction  
Régression logistique : variable explicative qualitative  
Régression logistique : variable explicative continue  
Régression logistique : variables explicatives mixtes

La régression logistique

Frédéric Bertrand et Myriam Maumy

Introduction  
Régression logistique : variable explicative qualitative  
Régression logistique : variable explicative continue  
Régression logistique : variables explicatives mixtes

La régression logistique

- On cherche donc à déterminer comment la probabilité de succès  $\pi$  change avec une ou plusieurs variables explicatives continues à partir des observations de  $y_i$  succès en  $n_i$  expériences indépendantes sous des valeurs de  $X$  observées égales à  $x_i$ , ( $i = 1, \dots, J$ ).
- On souhaite utiliser une modélisation de la côte de succès sachant que  $X = x$ , c'est-à-dire :

$$(Y|X = x_i) \sim \mathcal{B}(n_i, \pi_i)$$

$$\text{logit}(\pi_i) = \theta_i = \theta_i(x_i).$$

On s'aperçoit qu'une transformation logarithmique serait la bienvenue.

$$\tilde{\theta}_i = \ln \left( \frac{y_i + 0.5}{n_i - y_i + 0.5} \right).$$

Pour avoir une première idée de la relation entre la côte de succès et  $X$ , on examine le **logarithme de la côte empirique** contre  $x_i$  :

## Régression logistique : variables explicatives mixtes

Le modèle suggéré est donc :

$$(Y|X = x_i) \sim \mathcal{B}(\eta_i, \pi_i)$$

avec

$$\text{logit}(\pi_i) = \theta_i = \alpha_0 + \beta_1 x_i$$

où

$$x_i = \log(\text{dose}_i).$$

- Dans l'exemple précédent, on a ignoré l'influence potentielle du sexe sur la probabilité de succès. L'analyse précédente indique que la dose influe de manière significative sur la probabilité qu'une larve meurt.
- Considérons le cas simple où on a à la fois une variable continue  $X$  et une variable qualitative  $Z$ . Les données sont donc du type  $(y_{ki}, \eta_{ki}, x_{ki}, z_{ki})$ . Le modèle suggéré est donc :

$$(Y|X = x_{ki}, Z = z_{ki}) \sim \mathcal{B}(\eta_{ki}, \pi_{ki})$$

avec

$$\text{logit}(\pi_{ki}) = \theta_{ki}.$$

Nous avons donc 5 modèles à notre disposition :

- $X + Z + X^*Z, (\alpha_0 + \alpha_k) + (\beta_1 + \tau_k)x_{ki}$ .
- $X + Z, (\alpha_0 + \alpha_k) + \beta_1 x_{ki}$ .
- $X, \alpha_0 + \beta_1 x_{ki}$ .
- $Z, \alpha_0 + \alpha_k$ .
- $1, \alpha_0$ .

Reste à détecter les modèles convenables à l'aide du test du  $G^2$ . Pour cela, on utilise Minitab et le fichier de données disponible sur le site.