

*Le modèle linéaire généralisé (logit)*  
*Master 2 Recherche SES-IES Analyse de données*

Ana Karina Fermin

Université Paris Nanterre

<http://fermin.perso.math.cnrs.fr/>

- 1 Modèle de régression logistique
- 2 Cotes et rapports de cotes
- 3 Références

# Objectif.

On souhaite “expliquer” une variable réponse  $Y$  par une variable explicative  $X$  (ou plusieurs variables explicatives  $X_1, X_2, \dots, X_p$ ) lorsque  $Y$  est 0 (échec) ou 1 (succès).

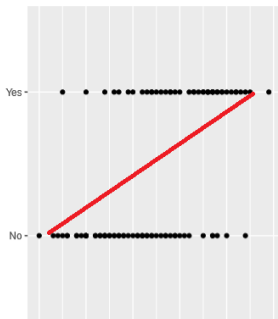
## Exemples:

- Médecine :  $Y$  vaut 1 si le patient atteint la maladie, 0 sinon. La variable  $X$  est l'âge.
- Banque :  $Y$  vaut 1 si le client fait défaut sur sa dette. La variable  $X$  est par exemple l'âge, la profession, le montant moyen mensuel d'utilisation de la carte de crédit, ...
- Sociologie :  $Y$  vaut 1 si le fils est cadre, 0 sinon. La variable  $X$  est par exemple le niveau d'éducation du père.,
- Socio-Eco :  $Y$  vaut 1 si vote pour président A, 0 sinon. La variable  $X$  est par exemple revenu, sexe, appartenance à une minorité, ...

Socio-Eco :  $Y$  vaut 1 si vote pour président A, 0 sinon. La variable  $X$  est par exemple le revenu



Socio-Eco :  $Y$  vaut 1 si vote pour président A, 0 sinon. La variable  $X$  est par exemple le revenu



Est ce que un ajuster un modèle de régression linéaire simple est raisonnable ?

# Modélisation (cas simple)

- On cherche à estimer

$$P(Y = 1|X)$$

- Nous supposons

$$P(Y = 1|X) = f(\beta_0 + \beta_1 X),$$

où  $f$  est une fonction inversible donnée avec  $\beta_0, \beta_1$  inconnus.

# Régression logistique (cas simple)

## Fonction $f$ cas de reg. logistique

On parle de régression **logit** ou **logistique** lorsque pour tout  $t \in \mathbb{R}$ ,

$$f(t) = \frac{\exp(t)}{1 + \exp(t)}.$$

$$P(Y = 1|X) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$

$$\log \left( \frac{P(Y = 1|X)}{1 - P(Y = 1|X)} \right) = \beta_0 + \beta_1 x$$

# Modélisation (cas multiple avec $p$ variables)

- $X = (X_1, X_2, \dots, X_p)$
- Notons  $\pi(X) = P(Y = 1|X) = P(Y = 1|X_1, X_2, \dots, X_p)$
- Nous supposons  $\pi(X) = f(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p)$ , où  $f$  est une fonction inversible donnée avec  $\beta_0, \beta_1, \dots, \beta_p$  inconnus.



# Régression logistique

## Fonction f cas de reg. logistique

On parle de régression **logit** ou **logistique** lorsque pour tout  $t \in \mathbb{R}$ ,

$$f(t) = \frac{\exp(t)}{1 + \exp(t)}.$$

$$\pi(x) = \frac{\exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p)}{1 + \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p)}$$

$$\log\left(\frac{\pi(x)}{1 - \pi(x)}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

## Estimation

- En pratique, les coefficients  $\beta_0, \beta_1, \dots, \beta_p$  doivent être déterminés à l'aide des données.
- On utilise la méthode du Maximum de Vraisemblance (MV).
- En général la méthode de MV fournit des estimateurs avec des bonnes propriétés statistiques.

# Loi Bernoulli

- **But:** prédire  $Y \in \{0, 1\}$  sachant  $X$
- Bernoulli  $\mathcal{B}(p)$ : loi sur  $\{0, 1\}$  tel que

$$Y \sim \mathcal{B}(p) \Leftrightarrow \begin{cases} P(Y = 1) = p \\ P(Y = 0) = 1 - p \end{cases}$$

## Loi Bernoulli Conditionnel

- Si  $Y \in \{0, 1\}$  alors

$$Y|X \sim \mathcal{B}(P(Y = 1|X))$$

- $P(Y = 1|X = x) = f(\beta_0 + \beta_1x_1 + \beta_2x_2 + \dots + \beta_px_p)$
- Commençons par définir la fonction log-vraisemblance associée au modèle logit, probit, ...

### log-Vraisemblance

$$LV(\beta) = \sum_{i=1}^n Y_i \log(f(X_i^t \beta)) + (1 - Y_i) \log(1 - f(X_i^t \beta))$$

avec  $\beta = (\beta_0, \beta_1, \dots, \beta_p)$ .

Les logiciels de statistiques calculent la fonction  $LV(\beta)$  et cherchent les coefficients  $\beta_0, \beta_1, \dots, \beta_p$  que maximisent cette fonction à l'aide d'un algorithme itérative.

Dans ce cours on va juste utiliser et interpréter les résultats donnés par le logiciel R (vous n'avez pas besoin de connaître les résultats théoriques de la log-vraisemblance associée au modèle) !!!

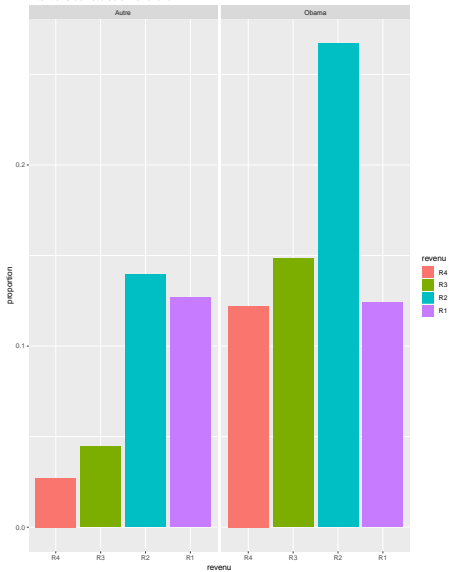
# Exemple : Élection présidentielle américaine de 2008.

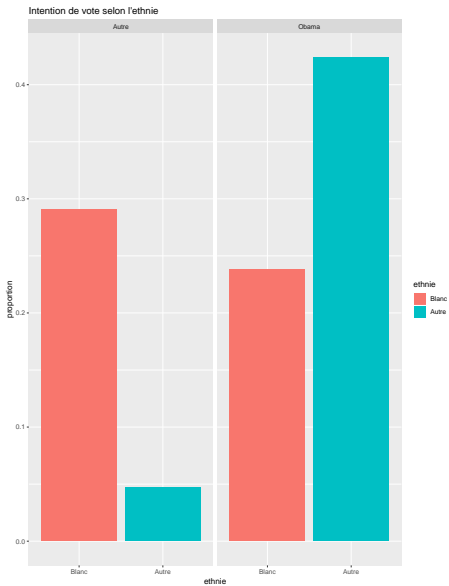
- On dispose d'un échantillon de taille 2322
- On s'intéresse pas exemple aux 4 variables
  - vote : 1 = vote pour Obama, 0 = ne vote pas pour Obama (Autre)
  - sexe : F= femmes , H = hommes
  - ethnie : Blanc et Autre (minorités ethniques)
  - revenu : R1 = très hauts revenus, R2 = hauts revenus, R3 = bas revenus, R4 = très bas revenus
- 5 électeurs au hasard

	sexe	ethnie	revenu	vote
446	H	Blanc	R3	Autre
59	H	Autre	R4	Obama
1258	F	Autre	R1	Autre
934	H	Blanc	R2	Autre
78	H	Autre	R2	Obama

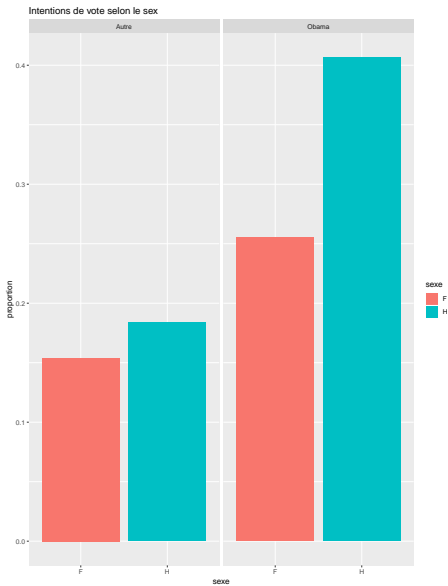
sexe	ethnie	revenu	vote
F:590	Blanc:762	R4:215	Autre:487
H:851	Autre:679	R3:278	Obama:954
R2:586			
R1:362			

Intentions de vote selon le revenu









- Données de l'élection présidentielle américaine de 2008.
- Est-ce que les électeurs de minorités ethniques (vs les Blancs), les gens moins favorisés économiquement (vs les plus favorisés) et les femmes (plus que les hommes) ont plus grande probabilité d'appuyer le candidat démocrate Barack Obama ?.
- But : Estimer la probabilité de voter pour Obama à partir des variables explicatives 'sexe', 'ethnie' et 'revenu'.

```
Call: glm(formula = vote ~ sexe + ethnie + revenu,
family = binomial("logit"), data = ElectionUSA2008_A)
```

Coefficients:

(Intercept)	sexeH	ethnieAutre	revenuR3	revenuR2	revenuR1
0.3161	0.2036	2.2842	-0.2809	-0.5742	-1.0881

Degrees of Freedom: 1440 Total (i.e. Null); 1435 Residual

Null Deviance: 1844

Residual Deviance: 1455 AIC: 1467

- $\text{sexe} = \text{"F"}, \text{ethnie} = \text{"Autre"}, \text{revenu} = \text{"R4"}$
- $\text{sexe} = \text{"H"}, \text{ethnie} = \text{"Autre"}, \text{revenu} = \text{"R4"}$

1                      2

0.9308788    0.9428834

- Une femme très défavorisée et appartenant à une minorité aurait voté pour Obama avec une proba de 93,1%. La probabilité pour un homme du même profil est de 94,3 %.

- 1 Modèle de régression logistique
- 2 Cotes et rapports de cotes
- 3 Références

# Cotes (odds) et rapports de cotes (odds ratios)

Dans le cas où la variable réponse  $Y$  est à valeurs dans  $\{0, 1\}$  et  $x = (x_1, x_2, \dots, x_p)$ , on définit :

La cote : 
$$C(x) = \frac{\pi(x)}{1 - \pi(x)}.$$

Le rapport de cotes : 
$$OR := R(x', x) = \frac{C(x')}{C(x)}.$$

## Cotes (odds) et rapports de cotes (odds ratios)

**Cas Simple** : Supposons qu'on dispose d'une unique variable explicative  $X$  de type qualitative à deux modalités  $\{0,1\}$ .

Si l'on suppose que  $\pi(x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$  on a alors

$$\log\left(\frac{\pi(x)}{1 - \pi(x)}\right) = \beta_0 + \beta_1 x$$

avec  $\beta_0$  et  $\beta_1$  inconnus.

- **La cote** :  $C(x) = \frac{\pi(x)}{1 - \pi(x)}$ .
- Coefficients estimés

$$\hat{\beta}_0 = \log(C(0)) \quad \text{et} \quad \hat{\beta}_1 = \log(C(1)/C(0)) = \log(OR)$$





## Rapports de cotes (odds ratios)

- Expliquer le vote en fonction du sexe

```
Call: glm(formula = vote ~ sexe,  
family = binomial("logit"), data = ElectionUSA2008_A)
```

Coefficients:

(Intercept)	sexeH
0.5054	0.2882

Degrees of Freedom: 1440 Total (i.e. Null); 1439 Residual

Null Deviance: 1844

Residual Deviance: 1837 AIC: 1841

- Odds-ratio

```
exp(0.2882)  
[1] 1.334024
```

# Cotes (odds) et rapports de cotes (odds ratios)

- Le rapport de cotes :  $OR := R(x', x) = \frac{C(x')}{C(x)}$ .
- Si toutes les coordonnées de  $x$  et  $x'$  sont identique hormis la  $j$ ème, pour laquelle  $x'_{[j]} = x_{[j]} + 1$ , alors on a

$$R(x', x) = \exp(\beta_j)$$

- 1 Modèle de régression logistique
- 2 Cotes et rapports de cotes
- 3 Références

## Références :

- Statistiques en sciences sociales avec R, Jean-Herman Guay
- An introduction to Generalized Linear Models, A.J. Dobson (2002)
- Statistiques avec  $\mathbb{R}$ , Pierre-André Cornillon et al. (2010), Presses universitaires de Rennes.
- Applied econometrics with R, Christian Kleiber et Achim Zeileis (2011), Springer.