

Compléments en régression logistique

① Interprétation des coefficients en régression logistique

En régression logistique, les paramètres β_1, \dots, β_p que l'on estime s'interprètent comme des logarithmes d'odds ratio (rapports de cotes)

①.1 Odds = cote = chance = risque

Soit Y une variable binaire $\begin{cases} 1 = \text{succès, présence...} \\ 0 = \text{échec, absence...} \end{cases}$

$$Y \sim \text{Bernoulli}(p) : \begin{cases} p = \mathbb{P}(Y=1) \\ 1-p = \mathbb{P}(Y=0) \end{cases}$$

On définit alors:

cote = $\frac{p}{1-p}$: Si un événement (ici $Y=1$) a une probabilité p , il a de grandes chances risques d'intervenir:

$$\begin{cases} p=0 \Leftrightarrow \text{cote} = 0 \\ p \approx 1 \Leftrightarrow \text{cote} \approx +\infty \end{cases}$$

①.2 Odds ratio = rapport de cotes

Mesure la liaison entre deux variables qualitatives binaires

X et Y . On note alors:

$$p_1 = \mathbb{P}(Y=1 \mid X=1)$$

$$p_0 = \mathbb{P}(Y=1 \mid X=0)$$

On définit alors l'odds ratio (OR) de Y associé à X

par:

$$\text{OR} = \frac{\frac{p_1}{1-p_1}}{\frac{p_0}{1-p_0}}$$

Par exemple: $Y = \text{présence/absence d'une maladie}$

X = variable d'exposition à la maladie par
exemple le sexe : $\left. \begin{array}{l} 1 = \text{homme} \\ 0 = \text{femme} \end{array} \right\}$

$\Rightarrow \frac{p_1}{1-p_1}$: risque de présence de la maladie chez les hommes

$\frac{p_0}{1-p_0}$ = risque de présence de la maladie chez les femmes

\Rightarrow Si :

$OR \approx 1$, la probabilité d'être malade est identique chez les hommes et chez les femmes donc le risque de maladie n'est pas associé au sexe

$\left\{ \begin{array}{l} OR > 1, \text{ le risque d'être malade est plus grand} \\ \text{chez les hommes que chez les femmes} \\ OR < 1, \text{ c'est le contraire} \end{array} \right.$

\rightarrow Dans les deux cas, il existe une association entre le risque de maladie et le sexe.

Autre exemple : Y = présence / absence d'une maladie
 X = âge du patient (quantitative)

\Rightarrow On définit la variable binaire $\left. \begin{array}{l} X = z \\ X = x+1 \end{array} \right\}$ où x = âge donné.

et on note

$$p_1 = P(Y=1 / X=x+1)$$

$$p_0 = P(Y=1 / X=x)$$

Dans ce cas :

$$OR = \frac{p_1}{1-p_1} \quad \text{mesure le risque que } Y=1 \text{ (maladie)}$$

$\frac{p_1}{1-p_1}$ par exemple) lorsque la variable X (l'âge par exemple) augmente d'une unité.

⇒ Si :

$OR \approx 0$, le risque de maladie n'est pas associé à l'âge.

$OR > 1$, le risque de maladie \uparrow lorsque l'âge augmente (d'un an).

$OR < 1$, c'est l'inverse.

1.3) les coefficients en régression logistique = log(odds ratio)

On peut montrer qu'en régression logistique les coefficients β_1, \dots, β_p associés aux p variables explicatives X^1, \dots, X^p s'interprètent comme des logarithmes d'odds ratios. On a vu que les variables explicatives X^j ont soit quantitatives, soit binaires (indicateurs des modalités des variables qualitatives).

→ Si X^j binaire :
$$\begin{cases} p_1 = P(Y=1 / X^j=1) \text{ et} \\ p_0 = P(Y=1 / X^j=0) \end{cases}$$

$e^{\beta_j} = \frac{p_1 / (1-p_1)}{p_0 / (1-p_0)}$ est l'odds ratio du risque de survenue de l'événement ($Y=1$) chez les personnes exposées au facteur X^j par rapport aux personnes non exposées, ajusté sur les autres variables explicatives du modèle (toute chose égale par ailleurs - TCEPA)

→ Si X^j quantitative :
$$p_1 = P(Y=1 / X^j = x+1)$$

$$| p_0 = \pi (Y=1 / X=x)$$

e^{β_j} est l'odds ratio du risque de survenue de l'événement $Y=1$ pour une augmentation d'une unité du facteur X^j , toute chose égale par ailleurs.

② Log-vraisemblance en régression logistique

Echantillon $(X_1, Y_1), \dots, (X_i, Y_i), \dots, (X_n, Y_n)$

et d'après le modèle :

$$p_i = \mathbb{P}(Y=y_i / X=x_i) = \frac{\exp(\beta_0 + x_i^T \beta)}{1 + \exp(\beta_0 + x_i^T \beta)}$$

$$\Rightarrow \ell(\beta_0, \beta) = \log \prod_{i=1}^n \mathbb{P}(Y_i=y_i / X_i=x_i)$$

$$= \log \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i}$$

$$= \sum_{i=1}^n [y_i \log p_i + (1-y_i) \log(1-p_i)]$$

$$= \sum_{i=1}^n \left[y_i \log \left(\frac{p_i}{1-p_i} \right) + \log(1-p_i) \right]$$

$$\text{Or } 1-p_i = 1 - \frac{\exp(\beta_0 + x_i^T \beta)}{1 + \exp(\beta_0 + x_i^T \beta)} = \frac{1}{1 + \exp(\beta_0 + x_i^T \beta)}$$

$$\log \left(\frac{p_i}{1-p_i} \right) = \text{logit}(p) = \beta_0 + x_i^T \beta$$

$$\Rightarrow \underline{\ell(\beta_0, \beta) = \sum_{i=1}^n y_i (\beta_0 + x_i^T \beta) - \log(1 + \exp(\beta_0 + x_i^T \beta))}$$
