

Polytomous Responses

Three response scales:

- ▶ Nominal response:
 - ▶ red, green, blue;
- ▶ Ordinal response:
 - ▶ young, middle aged, old
 - ▶ dislike very much, dislike, no opinion, like, like very much
- ▶ Interval response (scores attached):
 - ▶ < 200 , $200 - 300$, $300 - 400$, > 400 (scores: 100, 250, 350, 500)

When $J = 2$, degenerate to binary response model; when $J > 2$, different models for different response types.

Multinomial Distribution

If the response variable is categorical, with more than two categories, then the response is polytomous. The J possible response values are called the response categories.

Consider a r.v. Y with the potential outcome in J categories. Let π_1, \dots, π_J denote the respective probabilities, with

$$\pi_1 + \dots + \pi_J = 1.$$

Y follows a categorical distribution.

If we group n independent observations, and use y_i to denote the number of observations in category i . Then, $\mathbf{y} = (y_1, \dots, y_J)$ follows a multinomial distribution with pmf

$$f(\mathbf{y}) = \frac{n!}{y_1! \dots y_J!} \pi_1^{y_1} \dots \pi_J^{y_J}.$$

- It is in the **exponential family** with $J - 1$ canonical parameters

$$\theta_j = \log \frac{\pi_j}{\pi_1}, \quad j = 2, \dots, J$$

- ▶ The Binomial distribution is a special case with $J = 2$.
- ▶ $Y_j \sim \text{Binomial}(n, \pi_j)$,
- ▶ $\text{cov}(Y_j, Y_k) = -n\pi_j\pi_k, (j \neq k)$
- ▶ The joint dist of $(Y_1, Y_2, n - Y_1 - Y_2)$ is multinomial with three parameters π_1, π_2 and $1 - \pi_1 - \pi_2$.
- ▶ The dist of (Y_1, \dots, Y_J) conditional on $Y_j = y_j$ is again multinomial on the reduced set of categories with parameters $n - y_j$ and $\pi_i/(1 - \pi_j)$.

Nominal Logistic Regression

Consider the nominal polytomous response.

Data: $\mathbf{y}_i = (y_{i1}, \dots, y_{iJ})$, $\sum_j y_{ij} = m_i$, \mathbf{x}_i , $i = 1, \dots, n$

Goal: study how $\boldsymbol{\pi}_i = (\pi_{i1}, \dots, \pi_{iJ})$ is related to \mathbf{x}_i

- ▶ Random component: $\mathbf{y}_i \sim \text{Multi}(m_i, \boldsymbol{\pi}_i)$
- ▶ One category is chosen as the reference category. Suppose this is the first category.
- ▶ Fit $J - 1$ models, one for each remaining category.

- ▶ Canonical link: $\theta_{ij} = \log \frac{\pi_{ij}}{\pi_{i1}} = \eta_{ij} = \mathbf{x}_i^T \boldsymbol{\beta}_j$

$$\pi_{ij} = \frac{\exp(\eta_{ij})}{1 + \sum_{j=2}^J \exp(\eta_{ij})}, \quad j = 2, \dots, J$$

$$\pi_{i1} = \frac{1}{1 + \sum_{j=2}^J \exp(\eta_{ij})}$$

- ▶ Parameter interpretation: θ_{ij} is the log odds for response category j vs 1. Thus, β_{jc} is the log odds ratio with one unit change in x_c .

- Categories are exchangeable. Reference category can be changed for different comparison.

$$\begin{aligned}\log \frac{\pi_{ij}}{\pi_{ir}} &= \log \frac{\pi_{ij}}{\pi_{i1}} - \log \frac{\pi_{ir}}{\pi_{i1}} = \eta_{ij} - \eta_{ir} \\ &= \mathbf{x}_i^T (\boldsymbol{\beta}_j - \boldsymbol{\beta}_r) \\ \log \frac{\pi_{i1}}{\pi_{ir}} &= -\mathbf{x}_i^T \boldsymbol{\beta}_r\end{aligned}$$

- When $J = 2$, this reduces to the standard logistic model.

Goodness-of-Fit

- Generalized Pearson χ^2 statistic

$$G = \sum_{i=1}^n \sum_{j=1}^J \frac{(y_{ij} - m_i \hat{\pi}_{ij})^2}{m_i \hat{\pi}_{ij}} = \sum_{i=1}^n \sum_{j=1}^J R_{p_{ij}}^2$$

- Deviance statistic

$$D = 2 \sum_{i=1}^n \sum_{j=1}^J y_{ij} \log \frac{y_{ij}}{m_i \hat{\pi}_{ij}} = \sum_{i=1}^n \sum_{j=1}^J R_{d_{ij}}^2$$

- When model is correct and $m_i \pi_{ij}$ are large for all $i = 1, \dots, n; j = 1, \dots, J$, both statistics are approximately $\chi^2((n-p)(J-1))$

Example: Car Preferences

In a study of motor vehicle safety, 150 men and 150 women were interviewed to rate how important air conditioning and power steering were to them when they were buying a car.

Sex	Age	Response			Total
		Unimportant	Import	Very Import	
Women	18-23	26 (58%)	12 (27%)	7 (16%)	45
	24-40	9 (20%)	21 (47%)	15 (33%)	45
	> 40	5 (8%)	14 (23%)	41 (68%)	60
Men	18-23	40 (62%)	17 (26%)	8 (12%)	65
	24-40	17 (39%)	15 (34%)	12 (27%)	44
	> 40	8 (20%)	15 (37%)	18 (44%)	41
Total		105	94	101	300

- ▶ sex and age group are covariates
- ▶ 6 covariate patterns in total
- ▶ 3 response categories (conceptually ordinal, but treated as nominal temporarily)

To fit a nominal logistic regression, we first choose the category “unimportant” as the reference category.

Define the following three dummy variables,

- ▶ x_1 : the indicator of men;
- ▶ x_2 : the indicator of age 24-40 years;
- ▶ x_3 : the indicator of age > 40 years.

The model is

$$\log\left(\frac{\pi_2}{\pi_1}\right) = \beta_{02} + \beta_{12}x_1 + \beta_{22}x_2 + \beta_{32}x_3$$

$$\log\left(\frac{\pi_3}{\pi_1}\right) = \beta_{03} + \beta_{13}x_1 + \beta_{23}x_2 + \beta_{33}x_3$$

Question: how to compare with a reduced model where $\beta_{3j} = 2\beta_{2j}$ ($j = 2, 3$)?