# M6: Categorical regression

## Notation

$(Y_i, x_i)$    independent pairs

response     covariates $k$    $k+1 = p$
as before (M2-M5

$\in \{1, 2, \ldots, c+1\}$    $(c+1)$ ordered    og    unordered categories
(ordinal)      (nominal)

reference category      mental impairment      alligator food type

dummy variable

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_c \end{bmatrix} \quad 0/1 \quad \text{a } c\text{-dim response vector}$$

multivariable response ← NEW

$y = (0, 0, \ldots, 0) \rightarrow$ cat $c+1$
$y = (1, 0, 0, \ldots, 0) \rightarrow$ cat 1
$y = (0, 1, 0, \ldots, 0) \rightarrow$ cat 2

## Multinomial distribution

$$\pi_r = P(Y = r) \qquad r = 1, \ldots, c$$

$$\pi_{c+1} = P(Y = c+1) = 1 - \pi_1 - \pi_2 - \cdots - \pi_c$$

for ungrouped data

One observation
$y = (y_1, y_2, \ldots, y_c)$
$0/1$

$$f(y \mid \pi) = \pi_1^{y_1} \cdot \pi_2^{y_2} \cdots \pi_c^{y_c} (1 - \pi_1 - \cdots - \pi_c)^{1 - y_1 - \cdots - y_c}$$

$M(1, \pi)$

1

m independent trials

$y_r$ = # obs from category $r$

$$f(y|\pi) = \frac{m!}{y_1! \; y_2! \cdots y_c! \; (m-y_1-y_2\cdots -y_c)!} \pi_1^{y_1} \cdots \pi_c^{y_c} (1-\pi_1-\cdots -\pi_c)^{m-y_1-y_2-\cdots -y_c}$$

$M(m, \pi)$

can be shown $\quad E(Y) = m \cdot \pi = \begin{bmatrix} m \cdot \pi_1 \\ \vdots \\ m \cdot \pi_c \end{bmatrix}$

Q: $E(Y_{c+1}) = m \cdot \pi_{c+1} = m \cdot (1 - \pi_1 - \pi_2 - \cdots - \pi_c)$

$\quad Cov(Y_1, Y_{c+1}) = -\pi_1 \cdot \pi_{c+1} \cdot m$

AIM: Model $\pi_{ir} = P(Y_i = r)$ as a function of covariates

# Regression with nominal response  [6.2]

a) We generalize the binary logit model

$$\log\left(\frac{P(Y_i=1)\overset{\pi_{i1}}{}}{P(Y_i=0)\underset{\pi_{i0}}{}}\right) = \eta_i$$

$$\pi_{i1} = \frac{\exp(\eta_i)}{1+\exp(\eta_i)}$$

to $c$ models of $\pi_{ir}$ vs $\pi_{i,c+1}$ for $r=1,..,c$

(*)  $$\log\left(\frac{\pi_{ir}}{\pi_{i,c+1}}\right) = \eta_{ir} = x_i^T \beta_r$$  ← $p{\times}1$ vector of regr. param.

need **one** vector for each response cat $1,..,c$

Observe: effects $\beta$ vary by pairing respective category with reference category.

Q: How many regression parameters do we need to estimate? $c \cdot p$

This means that for some pair of response categories $(a,b)$

$$\log\left(\frac{\pi_{ia}}{\pi_{ib}}\right) = \underbrace{\log\left(\frac{\pi_{ia}}{\pi_{i,c+1}}\right)}_{x_i^T \beta_a} - \underbrace{\log\left(\frac{\pi_{ib}}{\pi_{i,c+1}}\right)}_{x_i^T \beta_b} = x_i^T(\beta_a - \beta_b)$$

3

b) Alternatively:

$$P(Y_i = r) = \pi_{ir} = \frac{\exp\left(x_i^T \beta_r\right)}{1 + \sum_{s=1}^{c} \exp(x_i^T \beta_s)}$$

$r = 1, \ldots, C$

<span style="color:red">observe all</span>

<span style="color:red">$\beta$'s contribute to $\pi_{ir}$</span>

$$P(Y_i = c+1) = 1 - \pi_{i1} - \cdots - \pi_{ic} = \frac{1}{1 + \sum_{s=1}^{c} \exp(x_i^T \beta_s)}$$

## Multivariate GLM

1. Random component $\quad Y_i \sim$ multinomial $\quad E(Y_i) = \underset{c \times 1}{\pi_i}$

2. Systematic component: $\quad \underset{(c \times 1)}{\eta_i} = \begin{bmatrix} \eta_{i1} \\ \vdots \\ \eta_{ic} \end{bmatrix} = \begin{bmatrix} x_i^T \beta_1 \\ \vdots \\ x_i^T \beta_c \end{bmatrix} \quad p \cdot C$ unknown param

3. Link function $\quad g(\mu_i) = \eta_i \quad$ when
   $\quad C \times 1 \qquad \qquad \pi_i$

   $$g_r(\pi_i) = \ln\left(\frac{\pi_{ir}}{\pi_{i,c+1}}\right)$$

   element $r$ of $C \times 1$ vector $\qquad \quad 1 - \pi_{i1} - \cdots - \pi_{ic}$

4

Ex: Alligators:

Q1: Why $df = 12$ :     8 groups and 5 responses $= 8 \cdot 4 = 32$
"free parameters"

$$4 \cdot (\underset{1}{\text{intercept}} + \underset{+\ 3}{\text{lake}} + \underset{+\ 1}{\text{size}}) = 20 \text{ param in } \beta\text{-vec}$$

$\Rightarrow$ residual $df = 32 - 20 = \underline{\underline{12}}$

Q2: Size & invertebrate : $\overset{\exp\left(\hat{\beta}_{\text{size}:1}\right) = 4.3}{\underbrace{\phantom{xxxxxxxxx}}_{\text{inv}}}$

$$\dfrac{\ln\left(\dfrac{P(\text{inv} \mid \text{size}=1, \text{lake a})}{P(\text{fish} \mid \text{size}=1, \text{lake a})}\right)}{\ln\left(\dfrac{P(\text{inv} \mid \text{size}=0, \text{lake a})}{P(\text{fish} \mid \text{size}=0, \text{lake a})}\right)} = \left(X_i^T - X_i^{T*}\right) \cdot \beta_{\text{inv}} = \beta_{\text{inv, size}}$$

$p \times 1$     $\hat{\beta}_{\text{size}:1}$

$\exp\left(\hat{\beta}_{\text{inv, size}}\right) = 4.3 = $ increase in OR of inv vs fish
if size change from 0 to 1

odds ratio

Regression with ordinal data: the cumulative model [6.3.1]

An unobserved latent variable $U_i$ drives the observation $Y_i$

$$Y_i = r \iff \theta_{r-1} < U_i \leq \theta_r$$

$$-\infty = \theta_0 < \theta_1 < \dots < \theta_c < \theta_{c+1} = \infty$$

Where $U_i = -x_i^T \beta + \varepsilon_i$

↗ no intercept

↖ error with cdf $F$
RV

logistic distribution ↓

Common $\beta$ for all categories

$$P(Y_i \leq r) = P(U_i \leq \theta_r) = P(-x_i^T \beta + \varepsilon_i \leq \theta_r)$$

$$= P(\varepsilon_i \leq \theta_r + x_i^T \beta) = F(\theta_r + x_i^T \beta)$$

no intercept

↗ no latent variables ($U_i$)

Different $F$ give different models and we will only consider
$F$ as the cdf in the logistic distribution

$$P(Y_i \leq r) = \frac{\exp(\theta_r + x_i^T \beta)}{1 + \exp(\theta_r + x_i^T \beta)}$$

⇕

$$\ln\left(\frac{P(Y_i \leq r)}{P(Y_i > r)}\right) = \theta_r + x_i^T \beta$$