# Chapter 3  Logistic regression

## Linear regression

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + e_i$$

$$\text{Assume} = e_i \overset{iid}{\sim} N(0, \sigma^2)$$

$$\mu_i = E(y_i) = \underbrace{\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + e_i}_{\text{linear predictor}}$$

## Logistic regression    binary data

① Ungrouped data
$$y_i = 1, 0 \qquad i = 1, \cdots, n$$

$$\underline{y} \quad \underline{x_1} \quad \underline{x_2} \quad \cdots \quad \underline{x_p} \qquad y_i \sim \text{Bernoulli}(p_i)$$

$$\begin{array}{c} 1 \\ \text{or } 0 \end{array}$$

② grouped data

$$\underline{n} \quad \underline{r} \quad \underline{x_1}, \cdots, \underline{x_p} \qquad r_i \sim \text{Binomial}(n_i, p_i)$$

| n | r |
|---|---|
| 100 | 60 |
| 500 | 95 |
| ⋮ | ⋮ |

e.g. MATH 3424

$p_i$ — depend on $x$

$x_1$ — assignment score

$x_2$ — quiz score

$x_3$ — final exam score

$$y - \begin{cases} 1 \\ 0 \end{cases}$$



$$100 \left\{ \begin{array}{c} \underline{y} \\ 1 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{array} \right. \begin{array}{c} \left.\vphantom{\begin{array}{c}1\\\vdots\end{array}}\right\} 60 \\ \left.\vphantom{\begin{array}{c}0\\\vdots\end{array}}\right\} 40 \end{array}$$

## Linear regression model

① Normality ✗     C.L.T.

$$r_i \sim \text{Binomial}(n_i, p_i)$$

$$n_i \longrightarrow \infty, \quad r_i \overset{\cdot}{\sim} N(n_i p_i, \; n_i p_i q_i)$$

② constant variance ? NO

$$\boxed{y_i} = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \rho_i \qquad y_i \sim \text{Bernoulli}(P_i)$$

$$\parallel$$

$$1, 0 \qquad E(y_i) = P_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + E(\rho_i)$$

$$\uparrow$$

$$\text{Assume}$$

$$\Rightarrow \quad \rho_i = \beta_i \quad y_i - P_i$$

$$= \begin{cases} 1 - P_i & \text{when } y_i = 1 \\ - P_i & \text{when } y_i = 0 \end{cases}$$

$$\text{Var}(\rho_i) = (1 - P_i)^2 * \Pr(y_i = 1) + (- P_i)^2 * \Pr(y_i = 0)$$

$$= (1 - P_i)^2 * P_i + (- P_i)^2 * (1 - P_i)$$

$$= P_i Q_i \qquad P_i - \text{diff. for diff } i$$

$$\neq \text{constant}$$

③ linearity $X$

$$\hat{P}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \cdots + \hat{\beta}_p x_{ip}$$

$$\downarrow \qquad \qquad \qquad \downarrow$$

$$0 < P_i < 1 \text{ here.} \qquad \text{give any real no.}$$

$\Rightarrow$ ~~make transformation of $\hat{P}_i$~~ ?

$$P_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip}$$

$\Rightarrow$ make transformation of $\mu_i$ ?

$$g(\mu_i) = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip}$$

$$\uparrow$$

link function

②

eg. ① $Y_i \sim N \Rightarrow$ linear regression

$$\mu_i = E(y_i) = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip}$$

$$\Downarrow$$

$$g(\mu_i) = \mu_i \quad - \quad \text{link} \quad - \quad \text{identity}$$

② $Y_i \sim \text{Bernoulli}(P_i) \Rightarrow$ logistic regression

$$g(\mu_i) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}$$

$$\uparrow$$

$$0 < P_i < 1 \qquad *$$

~ logit function

$$g(\mu_i) = \ln\left(\frac{\mu_i}{1-\mu_i}\right)$$

$$\Rightarrow g(\mu_i) \, \log_e\left(\frac{P_i}{1-P_i}\right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}$$

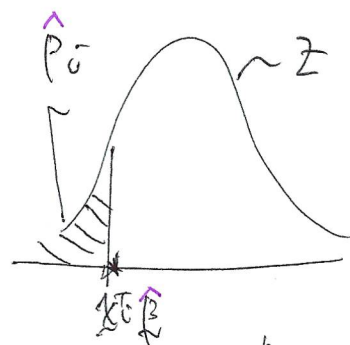$$\Rightarrow \quad \frac{P_i}{1-P_i} = \exp(\beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip})$$

$$\Rightarrow \quad P_i = \frac{\overset{+ve}{\boxed{\exp}}(\beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip})}{\boxed{1+}\exp(\beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip})}$$

$$0 < P_i < 1$$

~ Probit function

$$g(\mu) = \boxed{\Phi}^{-1}(\mu)$$

$$\uparrow$$

c.d.f. of $Z$



$\hat{P}_i \quad \sim Z$

$X_i'\beta$

$$\Rightarrow g(P_i) = \Phi^{-1}(P_i) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}$$

$$\Rightarrow P_i = \Phi(\beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip})$$

③

③ $Y_i \sim \text{Poisson}(\mu_i)$  $\xrightarrow{\quad >0 \quad}$  Poisson Regression

$$g(\mu_i) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}$$

Choose $g(\mu_i) = \log_e(\mu_i)$  — Link function = log function

$$\Rightarrow \mu_i = \exp(\underbrace{\beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}})$$

Generalized linear regression
$\uparrow$

DISTRIBUTION & LINK

- Estimation / Hypothesis Testing

Estimation

grouped data  ① $r_i \sim \text{Binomial}(n_i, p_i)$

$$n_i \rightarrow \infty$$

② Link function $g(\mu_i) = g(p_i) = \ln\left(\frac{p_i}{1-p_i}\right)$

③ ~~Got~~ Constant variance

⟶ If we ~~fit~~ fit

$$\log\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right) \text{ on } X_1, \cdots, X_p$$

Is $\text{Var}\left(\log\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right)\right)$ constant ?

④

Example 1 – one independent variable

Data

population prop. $\longleftarrow$ Sample proportion

$x \quad n \quad r \quad \log_e(x)$

$P_i$ depends on $x$

$\ln\left(\frac{P_i}{1-P_i}\right) = \beta_0 + \beta_1 * \text{logload}_i$

$\Rightarrow$ grouped data

| Obs | load | speci | fail | logload |
|-----|------|-------|------|---------|
| 1 | 5 | 600 | 13 | 1.60944 |
| 2 | 35 | 500 | 95 | 3.55535 |
| 3 | 70 | 600 | 189 | 4.24850 |
| 4 | 80 | 300 | 95 | 4.38203 |
| 5 | 90 | 300 | 130 | 4.49981 |

$P_1$
$P_2$
$P_3$
$P_4$
$P_5$

$r_i \sim \text{Binomial}(n_i, P_i)$

$n_i$ large $\forall i$

Fit $\ln\left(\frac{\hat{P}_i}{1-\hat{P}_i}\right)$ on $\text{logload}_i$

$$\hat{P}_i = \frac{r_i}{n_i} \qquad i = 1, 2, \cdots, 5$$

Model $= \ln\left(\frac{P_i}{1-P_i}\right) = \boxed{\beta_0} + \boxed{\beta_1} X_{i1} \qquad i = 1, \cdots, 5$

| Obs | $n$ | $r$ | $\hat{P}$ | $y$ | logload |
|-----|-----|-----|-----------|-----|---------|
| 1 | 600 | 13 | 13/600 | $\ln\left(\frac{13/600}{1-13/600}\right) = -3.81008$ | |
| 2 | 500 | 95 | 95/500 | $\ln\left(\frac{95/500}{1-95/500}\right) = -1.45001$ | |
| 3 | 600 | 189 | | $-0.77685$ | |
| 4 | 300 | 95 | | $-0.76913$ | |
| 5 | 300 | 130 | | $-0.26826$ | |

$\Rightarrow$ Fit a model of $y$ on $X$

$\ln\left(\frac{\frac{r}{n}}{1-\frac{r}{n}}\right) \qquad$ log load

Is $\text{Var}\left(\ln\left(\frac{\hat{P}}{1-\hat{P}}\right)\right)$ constant?

3423    Delta Method

$$Var\left(\log\left(\frac{\hat{P_i}}{1-\hat{P_i}}\right)\right) \overset{?}{=} \text{constant}$$

$\underbrace{\phantom{xxxxxxxxxx}}$
$$f(\hat{P_i})$$

$$f(\hat{P_i}) = f(P_i) + f'(P_i) * (\hat{P_i} - P_i) + \cdots$$

$\underbrace{\phantom{xxxxxxx}}_{\text{constant}}$

$$Var(f(\hat{P_i})) = (f'(P_i))^2 \, Var(\hat{P_i}) \qquad \neq$$

$$f(P_i) = \ln\left(\frac{P_i}{1-P_i}\right) \qquad \hat{P_i} = \frac{r_i}{n_i}$$

$$f'(P_i) = \frac{1}{P_i(1-P_i)}$$

Assume $r_i \sim$ Binomial
$$(n_i, P_i)$$

$$\Rightarrow Var(r_i) = n_i P_i (1-P_i)$$

$$Var(\hat{P_i}) = \frac{1}{n_i^2} Var(r_i)$$

$$Var\left(\ln\left(\frac{\hat{P_i}}{1-\hat{P_i}}\right)\right) = \left(\frac{1}{P_i(1-P_i)}\right)^2$$

$$= \frac{1}{n_i^2} n_i P_i(1-P_i)$$

$$= \frac{P_i(1-P_i)}{n_i}$$

$$* \frac{P_i(1-P_i)}{n_i}$$

$$= \boxed{\frac{1}{n_i \, P_i (1-P_i)}} \sim \sigma_i^2 \quad \neq \text{constant}$$

depends on $x$

$$i = 1, \cdots, S$$

Weighted least squares

Method of least squares from Chapter 1 $\quad Var(\underline{e}) = \sigma^2 \underline{I}$

$$\text{Min} \sum_{i=1}^{n} \hat{e_i}^2 \quad \text{OR Min Res S.S.}$$

$$\Rightarrow \text{Min} \sum_{i=1}^{n} (y_i - \hat{y_i})^2 \Rightarrow \text{Min} (\underline{Y} - \underline{X}\underline{\beta})^T (\underline{Y} - \underline{X}\underline{\beta})$$

$$\Rightarrow \hat{\underline{\beta}} = (\underline{X}^T\underline{X})^{-1} \underline{X}^T \underline{Y} \quad \& \quad Var(\hat{\underline{\beta}}) = (\underline{X}^T\underline{X})^{-1}\sigma^2$$

$\text{(6)}$

# Weighted least squares

$$\underset{\sim}{Y} = X\underset{\sim}{\beta} + \underset{\sim}{e}$$

If $Var(\underset{\sim}{e}) = V \implies Var(\underset{\sim}{Y}) = V$

Find $\tilde{\beta}_V$ s.t.

$$\text{Min } SS_{res,V} = (\underset{\sim}{Y} - X\tilde{\beta}_V)^T \underset{\sim}{V}^{-1} (\underset{\sim}{Y} - X\tilde{\beta}_V)$$

$$= (\underset{\sim}{Y}^T - \tilde{\beta}_V^T X^T) \underset{\sim}{V}^{-1} (\underset{\sim}{Y} - X\tilde{\beta}_V)$$

$$= (\underset{\sim}{Y}^T \underset{\sim}{V}^{-1} - \tilde{\beta}_V^T X^T \underset{\sim}{V}^{-1})(\underset{\sim}{Y} - X\tilde{\beta}_V)$$

$$= \underset{\sim}{Y}^T \underset{\sim}{V}^{-1} \underset{\sim}{Y} - - - - - \quad \text{is minimized}$$

$$\implies \tilde{\beta}_V = \boxed{(X^T \underset{\sim}{V}^{-1} X)^{-1} (X^T \underset{\sim}{V}^{-1} \underset{\sim}{Y})}$$

If $\underset{\sim}{V} = \sigma^2 \underset{\sim}{I} \implies \hat{\beta} = (X^T X)^{-1} X^T \underset{\sim}{Y}$

From Theorem 3.2 in Chapter 1

$$Var(\underset{\sim}{C} \underset{\sim}{Y}) = \underset{\sim}{C} Var(\underset{\sim}{Y}) \underset{\sim}{C}^T$$

$$\implies Var(\tilde{\beta}_V) = (\cancel{X^T V})(X^T \underset{\sim}{V}^{-1} X)^{-1} X^T \underset{\sim}{V}^{-1} \underset{\sim}{V}$$

$$\underset{\sim}{V}^{-1} X (X^T \underset{\sim}{V}^{-1} X)^{-1}$$

$$= (X^T \boxed{\underset{\sim}{V}^{-1}} X)^{-1} \qquad {\color{green} Var(\hat{\beta}) = (X^T X)^{-1} \sigma^2}$$

$${\color{green} \text{Consider } \underset{\sim}{V} = \sigma^2 \underset{\sim}{I}}$$

Let $\underset{\sim}{V} = \begin{pmatrix} \sigma_1^2 & & \underset{\sim}{0} \\ & \ddots & \\ \underset{\sim}{0} & & \sigma_s^2 \end{pmatrix}$ ⟵ obs. are indep.

$$\implies SS_{res,\text{weighted}} = (\underset{\sim}{Y} - X\tilde{\beta}_V)^T \underset{\sim}{V}^{-1} (\underset{\sim}{Y} - X\tilde{\beta}_V)$$

$$\hat{\sigma}_i^2 = \hat{Var}\left(\log\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right)\right) \qquad = \sum_{i=1}^{S} \frac{(y_i - \hat{y}_i)^2}{\sigma_i^2}$$

$$= \frac{1}{n_i \hat{p}_i (1-\hat{p}_i)} \qquad = \sum_{i=1}^{S} \boxed{\frac{1}{\sigma_i^2}} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{S} W_i (y_i - \hat{y}_i)^2$$

$$W_i = \frac{1}{\hat{\sigma}_i^2} = n_i \hat{p}_i (1-\hat{p}_i) = n_i \frac{r_i}{n_i}\left(1 - \frac{r_i}{n_i}\right) = \frac{r_i(n_i - r_i)}{n_i} \quad ⑦$$

Example 1 – one independent variable

Data

$x$    $n$    $r$    $\log(x)$

$$\ln\left(\frac{\frac{r_i}{n_i}}{1-\frac{r_i}{n_i}}\right)$$

$$n_i\left(\frac{r_i}{n_i}\right)\left(1-\frac{r_i}{n_i}\right) = \frac{r_i(n_i-r_i)}{n_i}$$

| Obs | load | speci | fail | logload | y | w |
|-----|------|-------|------|---------|-----------|---------|
| 1 | 5 | 600 | 13 | 1.60944 | -3.81608 | 12.718 |
| 2 | 35 | 500 | 95 | 3.55535 | -1.45001 | 76.950 |
| 3 | 70 | 600 | 189 | 4.24850 | -0.77685 | 129.465 |
| 4 | 80 | 300 | 95 | 4.38203 | -0.76913 | 64.917 |
| 5 | 90 | 300 | 130 | 4.49981 | -0.26826 | 73.667 |

x =logload

**Deviance and Pearson Goodness-of-Fit Statistics**

| Criterion | Value | DF | Value/DF | Pr > ChiSq |
|-----------|-------|-----|----------|------------|
| Deviance | 5.3883 | 3 | 1.7961 | 0.1455 |
| Pearson | 5.3792 | 3 | 1.7931 | 0.1460 |

Parameter estimates with confidence interval

**Analysis of Maximum Likelihood Estimates**

| Parameter | DF | Estimate | Standard Error | Wald Chi-Square | Pr > ChiSq |
|-----------|-----|----------|----------------|-----------------|------------|
| Intercept | 1 | -5.5784 | 0.3682 | 229.4877 | <.0001 |
| logload | 1 | 1.1400 | 0.0893 | 163.0932 | <.0001 |

Covariance matrix

**Estimated Covariance Matrix**

| Parameter | Intercept | logload |
|-----------|-----------|---------|
| Intercept | 0.1356 | -0.03253 |
| logload | -0.03253 | 0.007968 |

1

$$\underset{\sim}{X} = \begin{pmatrix} 1 & 1.60944 \\ 1 & 3.55535 \\ 1 & 4.24850 \\ 1 & 4.38203 \\ 1 & 4.49981 \end{pmatrix} \qquad \underset{\sim}{X} = \begin{pmatrix} 1 & x_{11} \\ \vdots & \vdots \\ 1 & x_{51} \end{pmatrix}$$

$$\underset{\sim}{V}^{-1} = \begin{pmatrix} 12.718 & & & & 0 \\ & 76.950 & & & \\ & & 129.465 & & \\ & & & 64.917 & \\ 0 & & & & 73.667 \end{pmatrix} \qquad \underset{\sim}{V}^{-1} = \begin{pmatrix} w_1 & & 0 \\ & \ddots & \\ 0 & & w_s \end{pmatrix}$$

$$\underset{\sim}{Y} = \begin{pmatrix} -3.81008 \\ -1.45001 \\ -0.77685 \\ -0.76913 \\ -0.26826 \end{pmatrix} \qquad \underset{\sim}{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_s \end{pmatrix}$$

$$\underset{\sim}{\tilde{\beta}} = (\underset{\sim}{X}^T \underset{\sim}{V}^{-1} \underset{\sim}{X})^{-1} (\underset{\sim}{X}^T \underset{\sim}{V}^{-1} \underset{\sim}{Y})$$

$$= \begin{pmatrix} \sum\limits_{i=1}^{s} w_i & \sum\limits_{i=1}^{s} w_i x_{i1} \\ \sum\limits_{i=1}^{s} w_i x_{i1} & \sum\limits_{i=1}^{s} w_i x_{i1}^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum\limits_{i=1}^{s} w_i y_i \\ \sum\limits_{i=1}^{s} w_i x_{i1} y_i \end{pmatrix}$$

$$= \begin{pmatrix} 357.717 & 1460.041 \\ 1460.041 & 6080.621 \end{pmatrix}^{-1} \begin{pmatrix} -330.31 \\ -1209.7 \end{pmatrix}$$

$$= \begin{pmatrix} 0.146036 & -0.0336246 \\ -0.0336246 & 0.00823819 \end{pmatrix} \begin{pmatrix} -336.31 \\ -1209.7 \end{pmatrix}$$

$$= \begin{pmatrix} -5.5784 \\ 1.1405 \end{pmatrix}$$