



北京航空航天大学

— 经济管理学院 —

BEIHANG UNIVERSITY  
SCHOOL OF ECONOMICS AND MANAGEMENT

# Generalized Linear Models

Lecture 3: Binary response



- 1 Logistic regression
- 2 Inference for logistic regression
- 3 Diagnostics for logistic regression
- 4 Model selection

Suppose response variable  $Y_i$  takes values 0 or 1 with probability  $P(Y_i = 1) = p_i$ . (i.e., a Bernoulli distribution)

We relate  $p_i$  to the predictors:

$$\eta_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_q x_{i,q}$$

$$\eta_i = g(p_i)$$

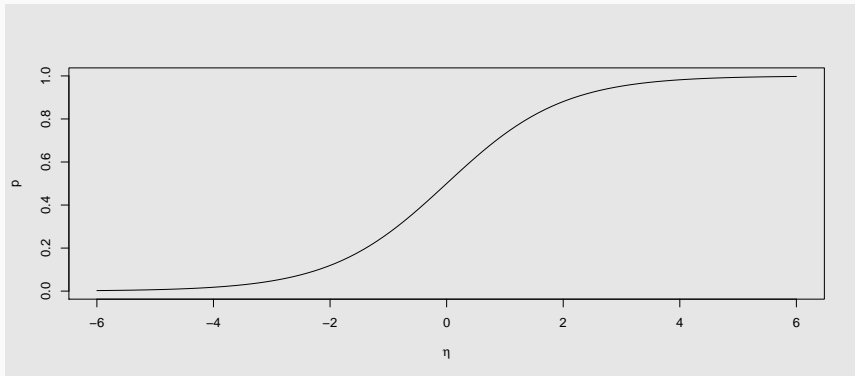
- The link function  $g$  should be monotone.
- $0 \leq g^{-1}(\eta) \leq 1$  for any  $\eta$ .

## logit link function

- $\eta = g(p) = \log\left(\frac{p}{1-p}\right)$  is the *logit* function. It maps  $(0, 1) \rightarrow \mathbb{R}$ .
- The inverse logit is  $g^{-1}(\eta) = \frac{e^\eta}{1+e^\eta}$  which maps  $\mathbb{R} \rightarrow (0, 1)$ .
- $p_i = \frac{e^{\eta_i}}{1+e^{\eta_i}} = P(Y_i = 1)$ .

## Inverse logit function

```
curve(ilogit(x), -6, 6, xlab=expression(eta), ylab='p')
```



- If  $p_i \approx 0.5$ , logistic and linear regression similar.

$$\eta_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_q x_{i,q}$$

$$\begin{aligned}\log L(\beta) &= \sum_{i=1}^n \log(p_i) 1_{Y_i=1} + \log(1 - p_i) 1_{Y_i=0} \\ &= \sum_{i=1}^n y_i [\eta_i - \log(1 + e^{\eta_i})] + (1 - y_i) \log \left( 1 - \frac{e^{\eta_i}}{1 + e^{\eta_i}} \right) \\ &= \sum_{i=1}^n y_i [\eta_i - \log(1 + e^{\eta_i})] + (1 - y_i) \log \left( \frac{1}{1 + e^{\eta_i}} \right) \\ &= \sum_{i=1}^n y_i [\eta_i - \log(1 + e^{\eta_i})] - (1 - y_i) \log(1 + e^{\eta_i}) \\ &= \sum_{i=1}^n [y_i \eta_i - \log(1 + e^{\eta_i})]\end{aligned}$$

$$\eta_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_q x_{i,q}$$

$$\begin{aligned}\log L(\beta) &= \sum_{i=1}^n \log(p_i) 1_{Y_i=1} + \log(1 - p_i) 1_{Y_i=0} \\ &= \sum_{i=1}^n y_i [\eta_i - \log(1 + e^{\eta_i})] + (1 - y_i) \log \left( 1 - \frac{e^{\eta_i}}{1 + e^{\eta_i}} \right) \\ &= \sum_{i=1}^n y_i [\eta_i - \log(1 + e^{\eta_i})] + (1 - y_i) \log \left( \frac{1}{1 + e^{\eta_i}} \right) \\ &= \sum_{i=1}^n y_i [\eta_i - \log(1 + e^{\eta_i})] - (1 - y_i) \log(1 + e^{\eta_i}) \\ &= \sum_{i=1}^n [y_i \eta_i - \log(1 + e^{\eta_i})]\end{aligned}$$

- Uses MLE to obtain  $\hat{\beta}$
- See Chapter 8 for estimation details.

```
fit <- glm(y ~ x1 + x2, family=binomial,  
          data=df)
```

- Bernoulli is equivalent to Binomial with only two levels.
- First (alphabetical) level is set to 0, other to 1. Use `relevel` if you want to change it.
- `glm` uses MLE with a logit link function (when `family=binomial`).



$$\text{Odds} = p/(1 - p).$$

- Odds are unbounded.
- $\log(\text{odds}) = \log(p/(1 - p)) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$
- A unit increase in  $x_1$  with  $x_2$  held fixed increases log-odds of success by  $\beta_1$ .
- A unit increase in  $x_1$  with  $x_2$  held fixed increases odds of success by a factor of  $e^{\beta_1}$ .

- 1 Logistic regression
- 2 Inference for logistic regression
- 3 Diagnostics for logistic regression
- 4 Model selection

- The likelihood ratio statistic:  $2 \log \frac{L_L}{L_S}$ , where  $L_L$  and  $L_S$  are likelihood of a larger model with  $l$  parameters and a smaller model with  $s$  parameters.

- The likelihood ratio statistic:  $2 \log \frac{L_L}{L_S}$ , where  $L_L$  and  $L_S$  are likelihood of a larger model with  $l$  parameters and a smaller model with  $s$  parameters.
- Now suppose we choose a saturated larger model, which has as many parameters as cases and has fitted values  $\tilde{p}_i = y_i$ . Then the statistic becomes:

$$\begin{aligned} D &= 2 \log L_L - 2 \log L_S \\ &= -2 \log L_S \\ &= -2 \sum_{i=1}^n \tilde{p}_i \log(\hat{p}_i) + \log(1 - \hat{p}_i) \end{aligned}$$

## *D* is called the Deviance.

- Difference between deviances equivalent to a likelihood ratio test.
- $D_S - D_L \sim \chi^2_{l-s}$ , where  $l$  and  $s$  are the number of parameters, assuming
  - 1 smaller model is correct
  - 2 models are nested
  - 3 distributional assumptions true
- Null deviance is for model with only an intercept.

In R:

```
fit <- glm(y ~ x1 + x2, family='binomial',  
          data=df)  
anova(fit, test="Chisq")  
drop1(fit, test="Chisq")  
anova(fit1, fit2, test="Chisq")
```

- NOT equivalent to t-tests on coefficients
- Deviance tests preferred

- Constructed using normal approximations for the parameter estimates.
- A  $100(1 - \alpha)\%$  confidence interval for  $\beta_i$  would be :  
$$\hat{\beta}_i \pm z^{\alpha/2} se(\hat{\beta}_i).$$
- Implemented in R using `confint`.

- 1 Logistic regression
- 2 Inference for logistic regression
- 3 Diagnostics for logistic regression
- 4 Model selection



**Response residuals:** Observation - estimate

$$e_i = y_i - \hat{y}_i$$

**Response residuals:** Observation - estimate

$$e_i = y_i - \hat{y}_i$$

**Pearson residuals:** Standardized

$$r_i = e_i / \hat{\sigma}$$

**Response residuals:** Observation - estimate

$$e_i = y_i - \hat{y}_i$$

**Pearson residuals:** Standardized

$$r_i = e_i / \hat{\sigma}$$

- Mean 0, variance 1.

**Response residuals:** Observation - estimate

$$e_i = y_i - \hat{y}_i$$

**Pearson residuals:** Standardized

$$r_i = e_i / \hat{\sigma}$$

- Mean 0, variance 1.

**Deviance residuals:** Signed root contribution to  $-2 \log L$ .

$$-2 \log L = c + \frac{1}{\hat{\sigma}^2} \sum e_i^2 = c + \sum d_i^2$$

$$d_i = e_i / \hat{\sigma}$$

**Response residuals:** Observation - estimate

$$e_i = y_i - \hat{p}_i$$

# Logistic regression residuals

**Response residuals:** Observation - estimate

$$e_i = y_i - \hat{p}_i$$

**Pearson residuals:** Standardized

$$r_i = \frac{y_i - \hat{p}_i}{\sqrt{\hat{p}_i(1 - \hat{p}_i)}}$$

**Response residuals:** Observation - estimate

$$e_i = y_i - \hat{p}_i$$

**Pearson residuals:** Standardized

$$r_i = \frac{y_i - \hat{p}_i}{\sqrt{\hat{p}_i(1 - \hat{p}_i)}}$$

- Mean 0, variance 1.

**Response residuals:** Observation - estimate

$$e_i = y_i - \hat{p}_i$$

**Pearson residuals:** Standardized

$$r_i = \frac{y_i - \hat{p}_i}{\sqrt{\hat{p}_i(1 - \hat{p}_i)}}$$

- Mean 0, variance 1.

**Deviance residuals:** Signed root contribution to  $-2 \log L$ .

$$-2 \log L = -2 \sum [\log(\hat{p}_i)y_i + \log(1 - \hat{p}_i)(1 - y_i)]$$

$$d_i = \text{sign}(y_i - \hat{p}_i) \sqrt{-2 [y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i)]}$$



In R:

```
fit <- glm(y ~ x1 + x2, family='binomial',  
          data=df)  
e <- residuals(fit, type='response')  
r <- residuals(fit, type='pearson')  
d <- residuals(fit, type='deviance')
```

- deviance is the default
- Residual plots can be hard to interpret
- Don't expect residuals to be normally distributed

- 1 Logistic regression
- 2 Inference for logistic regression
- 3 Diagnostics for logistic regression
- 4 Model selection

### Akaike's Information Criterion

$$\text{AIC} = -2 \log L + 2q$$

or

$$\text{AIC} = \text{deviance} + 2q$$

- Select model with smallest AIC
- Beware that model selection is not magical