

STA 303 / 1002

Note Title

2/16/2012

TEST: Monday Feb. 27 EX 310 / 320

All the details are on Blackboard

See Blackboard for office hours next week  
(announcement)

COME EARLY !!

Assignment 1 marks will be on Blackboard sometime  
to monitor

Last: Binomial logistic regression model

Deviance GOF test compares fitted model to saturated model

Test statistic:  $\hat{y}_i$  is binomial count est'd from fitted model

$$G^2 = 2 \sum_{i=1}^n \left\{ y_i \log(y_i) - y_i \log(\hat{y}_i) + (m_i - y_i) \log(m_i - y_i) - (m_i - y_i) \log(m_i - \hat{y}_i) \right\}$$

Why is there no Deviance GOF test in the Binary logistic regression case?

In Binary case,  $m_i = 1$  (all  $i$ ) and  $y_i = \begin{cases} 0 \\ 1 \end{cases}$

$$G^2 \text{ becomes } 2 \sum_{i=1}^n \left\{ -n_i \log(\hat{y}_i) - (n_i - n_i) \log(n_i - n_i) \right\}$$

The terms that come from logs are gone so  
 $G^2$  no longer includes both  $\ln \hat{y}_i$  &  $\ln$

For model assessment in Binomial Logistic Regression:

- ① Plot observed logits (log<sub>e</sub> of response proportions) versus quantitative explanatory variable to see if linear relationship is appropriate.
- ② Deviance GOF test

### ③ Examine residuals

"Residuals"

$$y_i - \hat{\pi}_{M,i}$$

observed proportion      fitted proportion

- useful to look for outliers
- need to standardize them

2 choices:

① Pearson residual

$$R_{res,i} = \frac{y_i - m_i \hat{\pi}_{m,i}}{\sqrt{m_i \hat{\pi}_{m,i} (1 - \hat{\pi}_{m,i})}}$$

$\sum_{i \in \text{set}}$   
of all  
binomial  
dists.

SRs: RESCH1

- ② Deviance Residual - defined so that  
the sum of the squares of the deviance  
 residuals is the deviance ( $G^2$ )

$$D_{res,i} = \frac{\text{sign}(y_i - m_i \hat{\pi}_{m,i}) * \sqrt{2 \left\{ y_i \log\left(\frac{y_i}{m_i \hat{\pi}_{m,i}}\right) + (m_i - y_i) \log\left(\frac{m_i - y_i}{m_i - m_i \hat{\pi}_{m,i}}\right) \right\}}}{m_i - m_i \hat{\pi}_{m,i}}$$

SAs: PESTON

Which standardized residual should you use?

- I usually look at both
- Both are asymptotically normally distributed (standard)

- if  $|r_{res,i}|$  or  $|D_{res,i}| > 2$ , possibly an outlier, particularly if a gap between the residual and the next largest residual
- if  $|r_{res,i}|$  or  $|D_{res,i}| > 3$ , classify as outlier
- For small sample sizes, the deviance residuals are closer to normally distributed than the Pearson residuals
- Pearson residuals are unstable when  $n$  close to 0 or 1

Looking at both Pearson and Deviance residuals for  
Kinnitt example, all are  $< 2$  in absolute value,  
so no outliers