

## Binomial Regression II

## Reminder: questions from Challenger disaster

- Forecast probability of an O-ring being damaged when the launch temperature is  $29^{\circ}F$ .
- How good is our forecast? Can we provide a confidence interval?
- Is temperature useful to predict the O-ring failing?

## Is temperature useful to predict the O-ring failing?

$Y_i$ , the number of damaged O-rings on the  $i$ -th launch, has distribution

$$Y_i \sim \text{bin}(6, p_i)$$

where

$$\log p_i / (1 - p_i) = \eta_i = \beta_0 + \beta_1 t_i.$$

Test for association between the number of damaged O-rings and temperature:

$$H_0 : \beta_1 = 0 \quad H_a : \beta_1 \neq 0$$

- Wald Test
- Likelihood Ratio Test

## Wald Test

Reminder: asymptotic normality of MLE

$$\hat{\theta}_i \sim \text{asy. } N(\theta_i^*, [\mathcal{I}(\hat{\theta})^{-1}]_{i,i}).$$

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} \quad \text{se}(\hat{\beta}_i) = \sqrt{\mathcal{I}(\hat{\beta})^{-1}_{ii}}$$
$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i \\ \vdots \\ \theta_m \end{bmatrix} \leftarrow \hat{\theta}_i \quad \text{MLE for } i\text{th parameter}$$

Test for association between the number of damaged O-rings and temperature

$$H_0 : \beta_1 = 0 \quad H_a : \beta_1 \neq 0$$

Wald test statistic:

$$z^* = \frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)} \sim \text{asy. } N(0, 1) \text{ under } H_0.$$

$$\text{so } \frac{\hat{\beta}_i - \beta_i^*}{\text{se}(\hat{\beta}_i)} \sim N(0, 1)$$

under  $H_0$ ,  $\beta_i^* = 0$

Challenger disaster

- See R script and result in “Wald Test” of Challenger.pdf
- $|z^*| = 4.07 > 1.96$  (critical value  $N(0, 1)$  at  $\alpha = 0.05$ )  $\Rightarrow$  reject  $H_0$ .
- p-value = 0.0000476.

# Likelihood Ratio Test (LRT)

Test for association between the number of damaged O-rings and temperature

$$H_0 : \beta_1 = 0 \quad H_a : \beta_1 \neq 0$$

Full model (F):

- $\eta_i = \beta_0 + \beta_1 t_i$
- Maximum log likelihood:  $\log \mathcal{L}(\hat{\beta}^F)$
- $\hat{\beta}^F$ : MLE of the parameters in the full model.

Reduced model (R): model under  $H_0$

- $\eta_i = \beta_0$
- Maximum log likelihood:  $\log \mathcal{L}(\hat{\beta}^R)$
- $\hat{\beta}^R$ : MLE of the parameters in the reduced model.

Compare two models.

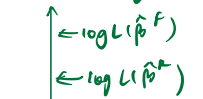
- Likelihood ratio test statistic:

$$LR^* = -2 \left[ \log \mathcal{L}(\hat{\beta}^R) - \log \mathcal{L}(\hat{\beta}^F) \right] \sim \text{asy. } \chi_1^2 \text{ under } H_0.$$

- $LR^* > \text{critical value from } \chi_1^2 \text{ at } \alpha \Rightarrow \text{reject } H_0.$

intuition =

model with more parameters will have a high likelihood.



if full model include useful parameters, then  $\log \mathcal{L}(\hat{\beta}^F)$  will increase a lot,

∴ difference in the number of parameters between full and reduced model is large,  $LR^*$  is large

full:  $\beta_0, \beta_1$   
reduced:  $\beta_0$

## Wald test vs Likelihood Ratio Test (LRT)

Wald test and LRT are asymptotically equivalent.

$$[z^*]^2 = \left[ \frac{\hat{\beta}_1}{se(\hat{\beta}_1)} \right]^2 \sim asy. \chi_1^2$$

- Precisely speaking, two tests are asymptotically equivalent in the sense that under  $H_0$  they reach the same decision with probability approaching 1 as  $n$  goes to infinite.

However, the chi-squared approximation to the log likelihood ratio is generally better than the normal approximation to the MLE.

*small sample size  $\Rightarrow$  prefer log likelihood ratio test*

## Likelihood Ratio Test (LRT): Challenger disaster

- See R script and result in “Likelihood Ratio test” and “Wald Test vs Likelihood Ratio test” of Challenger.pdf
- $LR^* = 21.98 > 3.84$  (critical value from  $\chi_1^2$  at  $\alpha = 0.05$ )  $\Rightarrow$  reject  $H_0$ .
- p-value = 0.0000027.

can obtain one model by constrain parameters

model 1 & 2 nested  
model 2 & 3 nested

model 1:

$$y_i = \beta_0 + \beta_1 \text{age}_i + \beta_2 \text{weight}_i$$

model 2:

$$y_i = \beta_0 + \beta_1 \text{age}_i$$

model 3:

$$y_i = \beta_0 + \beta_2 \text{weight}_i$$

## Likelihood Ratio Test (LRT) for model selection

In general, likelihood ratio test is used to select between two nested models (one model can be obtained by constraining parameters of another model).

Full model (F):

- Maximum log likelihood:  $\log \mathcal{L}(\hat{\theta}^F)$

Reduced model (R):

- Maximum log likelihood:  $\log \mathcal{L}(\hat{\theta}^R)$

Let  $k$  indicate the difference in the number of parameters between two models.

Compare two nested models.

- Under the reduced model

$$LR^* = -2 \left[ \log \mathcal{L}(\hat{\theta}^R) - \log \mathcal{L}(\hat{\theta}^F) \right] \sim \text{asy. } \chi_k^2.$$

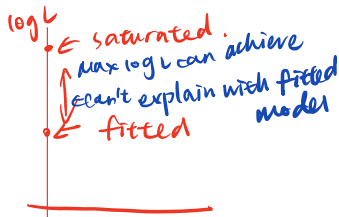
- $LR^* > \text{critical value from } \chi_k^2 \text{ at } \alpha \Rightarrow \text{select the full model.}$



we know:

temperature is useful

but is the model good enough



if  $D$  too big,  
means we miss lots of info

## (Scaled) Deviance

The scaled deviance is used to judge model adequacy.

For the binomial regression model the deviance is the same as the *scaled deviance*, which is defined as the log likelihood ratio for the fitted model compared to the saturated model.

Full model (F):

# of parameter = # of observation

- The saturated model has the same number of parameters and the observations.
- Maximum log likelihood:  $\log \mathcal{L}(\hat{\theta}^F)$

Reduced model (R):

- The fitted model.
- Maximum log likelihood:  $\log \mathcal{L}(\hat{\theta}^R)$

The scaled deviance:

$$D = -2 \left[ \log \mathcal{L}(\hat{\theta}^R) - \log \mathcal{L}(\hat{\theta}^F) \right].$$

## (Scaled) Deviance



**Warning:** the number of parameters in the saturated model is  $n$ , which is not fixed, so the theory of maximum likelihood does not apply, and  $D$  may not converge to a chi-squared distribution.

but for binomial regression, still can use  
chi-squared approximation to the log-likelihood  
ratio

## (Scaled) Deviance for binomial regression model

$$Y_i \sim \text{Bi}(m_i, p_i) \quad p_i = g^{-1}(x_i^T \beta)$$

The saturated model allows for one parameter for each observation. For

binomial regression the saturated model has  $p_1, p_2, \dots, p_n$  as parameters.

Clearly, for this model we estimate  $p_i$  by  $y_i/m_i$ . Let  $\hat{p}_i = g^{-1}(x_i^T \hat{\beta})$  be our (not saturated) model estimate of  $p_i$ , then the scaled deviance is

fitted model  
fitted parameter.

fitted

$$D = -2 \left[ \sum_{i=1}^n (y_i \log \hat{p}_i + (m_i - y_i) \log (1 - \hat{p}_i)) - \sum_{i=1}^n (y_i \log \frac{y_i}{m_i} + (m_i - y_i) \log (1 - \frac{y_i}{m_i})) \right]$$

$$\begin{aligned} D &= -2 \sum_{i=1}^n \left( y_i (\log \hat{p}_i - \log \frac{y_i}{m_i}) + (m_i - y_i) (\log (1 - \hat{p}_i) - \log (1 - \frac{y_i}{m_i})) \right) \\ &= -2 \sum_{i=1}^n \left( y_i \log \frac{\hat{y}_i}{y_i} + (m_i - y_i) \log \frac{m_i - \hat{y}_i}{m_i - y_i} \right) \end{aligned}$$

where  $\hat{y}_i = m_i \hat{p}_i$  is the  $i$ -th fitted value.

## (Scaled) Deviance for binomial regression model for testing model adequacy

**It just happens:** if  $m_i p_i$  and  $m_i(1 - p_i)$  are large enough ( $\geq 5$  is a common rule of thumb), then for a binomial regression model, if the model is adequate then  $D \approx \chi^2_{n-k}$ , where  $k$  is the number of parameters in the fitted model (including  $\beta_0$ ).

In this case the (scaled) deviance can be used as a test for model adequacy. If  $D$  is too large (as compared to a  $\chi^2_{n-k}$ ), then the model is missing something.

For a binomial model with small  $m_i$  we can't use the (scaled) deviance directly to test model adequacy, but we can still use it for model selection.

why?

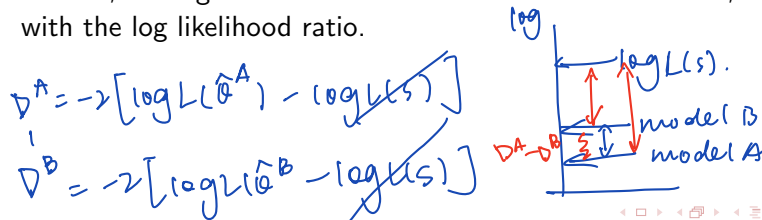
## Use the scaled deviance for model selection (LRT)

If model A is nested within model B, and model A has (scaled) deviance  $D^A$  and model B has (scaled) deviance  $D^B$ , then

$$D^A - D^B = -2 \left[ \log \mathcal{L}(\hat{\theta}^A) - \log \mathcal{L}(\hat{\theta}^B) \right],$$

where  $\hat{\theta}^A$  and  $\hat{\theta}^B$  are MLEs for the models A and B, respectively.

That is, the log likelihood for the saturated model cancels, and we are left with the log likelihood ratio.



## Use the scaled deviance for model selection (AIC)

The Akaike Information Criterion is used for model selection:

$$\text{AIC} = 2k - 2 \log \mathcal{L}(\hat{\theta})$$

where  $k$  is the number of parameters in the model. Given a choice, we prefer that model with the smaller AIC.

If model B has  $s$  more parameters than model A (not necessarily nested within B), then

$$\begin{aligned} \text{AIC}^B - \text{AIC}^A &= 2s - 2 \log \mathcal{L}(\hat{\theta}^B) + 2 \log \mathcal{L}(\hat{\theta}^A) \\ &= 2s - D^A + D^B. \end{aligned}$$

## (Scaled) deviance: Challenger disaster

- See R script and result in “Deviance” of Challenger.pdf

## Reminder: questions from Challenger disaster

- Forecast probability of an O-ring being damaged when the launch temperature is  $29^{\circ}F$ .
- How good is our forecast? Can we provide a confidence interval?
- Is temperature useful to predict the O-ring failing?



## Learning goals

Understand binomial regression

- know when you should use binormal regression.
- be able to write binomial regression model and its likelihood.
- be able to obtain estimators of parameters or function of parameters using R script.
- be able to quantify uncertainty of the estimators (e.g., computing CI).
- be able to test hypothesis.
- be able to do model selection.

Understand asymptotic properties of MLEs (maximum likelihood estimators)

- use asymptotic normality of MLEs to quantify uncertainty of the estimators (e.g., Wald CI).

Understand Wald test and likelihood ratio test (LRT)

- use them to test hypothesis in binomial regression

Understand (scaled) deviance

- use it to test model adequacy or perform LRT and model comparison.