

STA 303 / 1002

Note Title

2/13/2012

TEST Monday Feb. 27

Where?

Last name

Location

A-K

EX 310

L-Z

EX 320

Bring: calculator + T-cards.

Coverage: Analysis of variance + logistic regression
To read of lecture on Thurs Feb 16

Formula sheet: similar to last year's test

↑ I suspect

Office hours during Reading Week (regular hours cancelled)

Thurs Feb 23

10:00 - 12:00

SS 5016A

Fri Feb 24

10:00 - 12:00

The Binomial Logistic Regression Model

Data: y_i observed Binomial count, $i=1, \dots, n$
 $y_i \sim \text{Binomial}(n_i, \pi_i)$

$$\text{Model: } \log\left(\frac{\pi_i}{1-\pi_i}\right) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

Example: Kruunut Iskads - off Finland

Objective: Learn how to create nature preserves that help preserve endangered species.
Are large or small preserves better?

Have counts of bird species for 18 islands.

Data: - area of island in km^2
- Number of species on each island in 1949
- " " " " " " in 1959

$N = 18$

$\pi_i =$ probability of "extinction" (extinction = species no longer on particular island in 1959)

Assume that this is the same for each species of bird on a particular island.

$m_i =$ number of species on island i in 1949

$y_i =$ number of species no longer there on island i in 1959

Assume species survival is independent

Then $y_i \sim \text{Binomial}(m_i, \pi_i)$

With these data (unlike Donner party example) we estimate π_i from data

"Response proportion" - est. of π_i from y_i, m_i

$$\hat{\pi}_{s,i} = \frac{y_i}{m_i}$$

S = "saturated"

"Observed" or "Empirical" $\text{Logit} : \text{Log} \left(\frac{\hat{\pi}_{s,i}}{1 - \hat{\pi}_{s,i}} \right)$

~~Proposed~~
Model : $\text{Log} \left(\frac{\pi_i}{1 - \pi_i} \right) = \beta_0 + \beta_1 \text{area}$

Can plot observed logs versus area to see if a linear model seems appropriate

For our example, from this plot, decided we should look at $\log(\text{area})$

The relationship between $\log(\text{area})$ and empirical

logs seems linear

Model we will fit:
$$\log\left(\frac{\pi_i^2}{1-\pi_i}\right) = \beta_0 + \beta_1 \log(\text{area})$$

SAs: Model statement:
model $y_i / m_i = x_i$;

Need to specify both count of # of successes and total number of trials

Output:

Number of observations: 18 (n) (# of trials)

Sum of frequencies: 632 = $\sum_{i=1}^{18} m_i$

Fitted model: $\logit(\hat{\pi}) = -1.196 - .297 \log(\text{area})$

Wald test: strong evidence that coef of $\log(\text{area})$ is not zero ($p < .0001$)

Same test as in binary logistic regression

$$H_0: \beta_1 = 0 \quad \text{vs} \quad H_a: \beta_1 \neq 0$$

$$\text{Test stat } z_{\text{obs}} = \frac{-0.2971}{0.0549} \quad \text{or} \quad \left(\frac{-0.2971}{0.0549} \right)^2$$

$$\text{Chi square (1)} = 29.3$$



$$95\% \text{ CI for } \beta_1: -0.2971 \pm 1.96 (0.0549) \\ = (-1.119, -0.074)$$

Interpreting β_1 (since area was log-transformed)

$$\text{logit}(\pi) = \beta_0 + \beta_1 \text{log}(x)$$

$$\Rightarrow \frac{\pi}{1-\pi} = e^{\beta_0} e^{\beta_1 \text{log}(x)}$$

Changing x by a factor of k , changes odds by a factor of k^{β_1} , changes

Commit example: Double island area $\rightarrow .2971$
odds change by a factor of $2 = .81$

So odds of extinction are 81% of odds on island half size.

Island Tests are same as for

binary logistic regression:

Fitted model: $\text{logit}(\hat{\pi}_i) = -1.196 - .297 \log(\text{area})$

For $\overline{\text{Wlkalcrummi}}$ ($i=1$) area = 185.5 km²

$$\text{logit}(\hat{\pi}_1) = -1.196 - .297 \log(185.5) \\ = -2.75$$

Estimated prob. of extinction for a species
on $\overline{\text{Wlkalcrummi}}$

$$\hat{\pi}_1 = \frac{e^{-2.75}}{1 + e^{-2.75}} = 0.060$$

Call this $\hat{\pi}_{M,1}$, $M_{i,j}$ for model^j

For comparison, the response proportion for
Altothmini is $\hat{\pi}_{s,1} = \frac{5}{75} = .067$

EffectPlot.png (add plots (only) = effect
- plots $\hat{\pi}_M$ & $\hat{\pi}_S$ on prob (logistic statement)
 \leftarrow points
 \leftarrow with smooth
 \leftarrow smooth curve

What's the big difference between binary and binomial
logits regression?

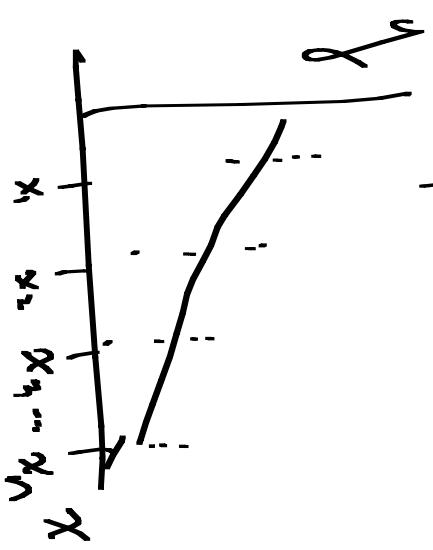
Can do more tests for model adequacy in binomial LR.

One such test: DEVIANCE GOODNESS-OF-FIT TEST (GOF)

Is a linear model appropriate / adequate?

The idea behind this test

(explained in terms of simple regression)



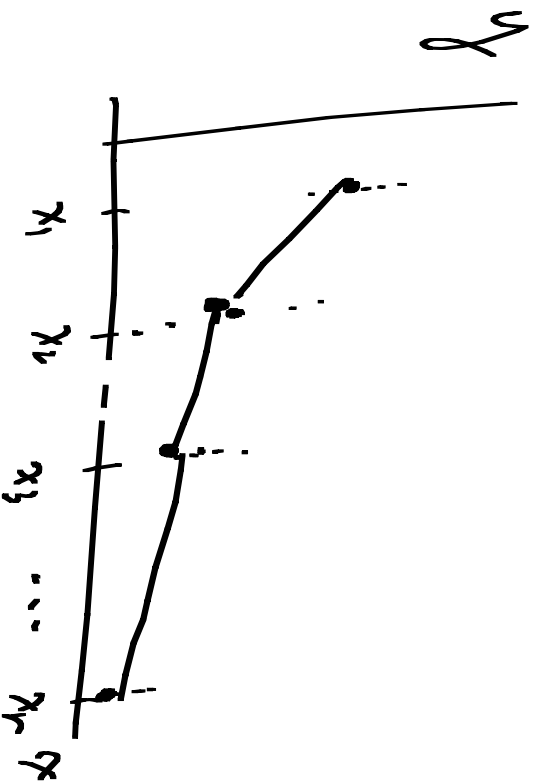
Suppose we have replicate observations at each value of x

$$y = \beta_0 + \beta_1 x + \epsilon$$

(n different values of x)

Alternative model:

Analysis of Variance: $y = \beta_0 + \beta_1 I_{x=x_1} + \beta_2 I_{x=x_2}$
 $+ \dots + \beta_{n-1} I_{x=x_{n-1}} + e$



-more β 's
but better model if
relationship isn't
linear.

There is a linear regression GOF test that
compares these 2 models

In logistic regression for binomial counts, think of

data for each x_i as m_i Bernoulli observations;
then have m_i observations at x_i

The Deviance G-D-F test compares:

① Model of Interest: $\text{logit}(\pi_i) = \beta_0 + \beta_1 x_i$ (β_2 parameter)

② Saturated Model: $\text{logit}(\pi_i) = \beta_0 + \beta_1 I_1 + \beta_2 I_2 + \dots + \beta_{n-1} I_{n-1}$ (n parameters)

In saturated model, each explanatory variable is an indicator variable. x_i is treated as a categorical variable with n categories

Use Likelihood Ratio Test to compare Model of Interest (reduced model) to Saturated Model (full model)

This called the "Drop-in-Deviance" test (sometimes)

Test statistic is
$$-2 \log \left(\frac{L_R}{L_F} \right) = -2 \log \left(\frac{L_M}{L_S} \right)$$

H_0 : fitted model fits data as well as saturated model

H_a : Saturated model is better

Under H_0 , the test statistic is an observation

from a chi-square distribution with $N - (p+1)$ df (This is an approximation if approximations work well)

Calculation of test statistics:

For saturated model (full model)

$$\hat{\pi}_{s,i} = \frac{y_i}{m_i}$$

For fitted model (reduced model)

$$\hat{\pi}_{m,i} = \hat{y}_i = m_i \hat{\pi}_{m,i}$$

estimated using MLE

For binomial logistic regression,
likelihood function

$$L = \prod_{i=1}^n \binom{m_i}{y_i} \pi_i^{y_i} (1 - \pi_i)^{m_i - y_i}$$

where $\pi_i =$

$$\frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}$$

$$\log L = \sum_{i=1}^n \left[y_i \log(\pi_i) + (m_i - y_i) \log(1 - \pi_i) + \log \binom{m_i}{y_i} \right]$$

Test statistic for Deviance GOF test is called the "deviance"

$$\underline{\text{Deviance}} = -2 \left(\log L_n - \log L_s \right) = 2 \left(\log L_s - \log L_n \right)$$

$$= 2 \sum_{i=1}^n \left\{ y_i \log \left(\frac{y_i}{m_i} \right) + (m_i - y_i) \log \left(\frac{m_i - y_i}{m_i} \right) + \log \binom{m_i}{y_i} - y_i \log \left(\frac{y_i}{m_i} \right) - (m_i - y_i) \log \left(\frac{m_i - y_i}{m_i} \right) - \log \binom{m_i}{y_i} \right\}$$

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

Small deviance \Rightarrow Model fits data well

To get deviance GOF test in SAS, we scale = none
option on model statement

Kruskal example: Test statistic = 12.06

H_0 : fitted model is equal
in fit to saturated model
 $df = 16 (= 18 - 2)$

H_a : saturated model fits better
 $p = .7397$

The data are consistent with H_0 ; we'll use the simpler model with linear function of log(sales)

For the deviance GOF test:

- Large p-value means
 - fitted model is adequate
 - OR - test is not powerful enough to detect inadequacies
 - Small p-value means
 - the fitted model is not correct (e.g. need more explanatory variables, or need polynomial model, or ...)
- OR - response distribution is not adequately modeled

by the General distribution

OR - there are some severe outliers