

SFA 302 / 1002

Reminder: All about SFS (except programming and interpreting output) on Thursday (Jan 19)

---

SFS: GLM procedure  
"General Linear Model"

- Alternative in proc reg
- doesn't have all the regression features of proc reg  
  eg: variable selection techniques
  - does have lots of other features; better for debugging

with categorical predictor variables

Output: ANOVA table as in previous

Then Type I SS (Sequential Sum of Squares)  
- we won't use in this course

Then Type III SS - contribution to regression sum of squares over and above everything else in model

From solution option: estimates of  $\beta$ 's and t-tests for tests  $H_0: \beta_j = 0$

Comparing output from `proc test`, `reg`, `glm`:

From output:

t<sub>test</sub>

Means:  $\bar{x}_{other} = 29.49$   
 $\bar{x}_{spoke} = 14.62$   
diff = 14.87

Test stat for testing  
equal means (assuming  
equal variances)

$$t_{obs} = 5.67$$

Distr under  $H_0$ :  $t_{44}$

p-value  $< .0001$

reg With  $x_i = I_{spoke,i}$

Estimates of parameters

$$b_0 = 29.49$$

$$b_1 = -14.87$$

$b_0$  is est  $E(Y_i | x_i = 0)$

$b_0 + b_1$  is est  $E(Y_i | x_i = 1)$

Test  $H_0: \beta_1 = 0$

argue about the both  
groups have same mean

Test stat:  $t_{obs} = -5.67$

Distr under  $H_0$ :  $t_{44}$

p-value  $< .0001$

glm

Estimates of parameters

$$b_0 = 14.62$$

$$b_1 = 14.87$$

$b_0$  is est  $E(Y_i | \text{spoke}_i = 0)$   
(i<sup>th</sup> obs is judge)

$b_0 + b_1$  is est

$E(Y_i | \text{i<sup>th</sup> obs} = 1)$   
(other judge)

Test  $H_0: \beta_1 = 0$

Test stat: 5.67

Distr under  $H_0$ :  $t_{44}$

p-value:  $< .0001$

Model being fit by  $prz$   $glm$ :

$$Y_i = \beta_0 + \beta_1 I_{\text{other}, i} + \epsilon_i$$

Type III  $SS$

F-test  $H_0: \beta_1 = 0$  vs  $H_a: \beta_1 \neq 0$   
Test Stat:  $F_{obs} = 32.15$  ( $= 5.67^2$ )

Distr under  $H_0$ :  $F_{1, 44}$

p-value:  $p < .0001$

Why does  $prz$   $GLM$  give  $NBTE$ : "The  $X'X$  matrix has been found to be singular ..."

Review of Regression in Matrix Terms

Model:  $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$   
 $i=1, \dots, N = \# \text{ of observations}$

Matrix form

where  $\tilde{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$ ,  $\tilde{Y} = X\tilde{\beta} + \tilde{\epsilon}$   
 $\tilde{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$ ,  $\tilde{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$

$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \dots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$   
 $n \times (p+1)$

Least squares estimator of  $\beta$ 's:

$$\hat{\beta} = (X'X)^{-1} X'y$$

$X'X$  has dimension  $(p+1) \times (p+1)$

Need  $X'X$  to be full rank to be invertible

$$\text{rank } X'X = \text{rank } X$$

Need  $X$  to be rank  $(p+1)$ ; need columns of  $X$  to be linearly independent

Class statement in pre-quiz creates dummy variable for all categories of the variable. So if I created I speak and I other =  $\begin{cases} 1 & \text{if other judge} \\ 0 & \text{if speaker judge} \end{cases}$

$$\Rightarrow X = \begin{pmatrix} \text{Total} & \text{I speak} \\ 0 & 0 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ \vdots & \vdots \\ 1 & 37 \end{pmatrix} \quad \left( \begin{array}{l} \text{linearly dependent} \\ \text{columns} \end{array} \right)$$

Prof also dropped the last column

For a categorical variable with 6 categories, need 5-1 indicator variables in the model

Answer to 1st question: Does Spork's judge differ from other judges?

We have very strong evidence ( $p < .0001$ ) that the mean  
to women differs between Spork's judge's venires and the  
other judge's venires. On average, Spork's has 14.6%  
women, and the other judge has 29.5% women.

2nd question: Is there evidence of any differences  
among other 6 judges?

$$H_0: \mu_A = \mu_B = \mu_C = \mu_D = \mu_E = \mu_F$$

where  $\mu_A$  is mean  
to women on  
venires for judge A

$H_a:$  at least one of these means is  
not equal to another



Linear model;

$$Y_i = \beta_0 + \beta_1 I_{A,i} + \beta_2 I_{B,i} + \beta_3 I_{C,i} + \beta_4 I_{D,i} + \beta_5 I_{E,i} + \epsilon_i$$

where  $I_{A,i} = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ observation is in group A} \\ 0 & \text{otherwise} \end{cases}$   
etc.

In general, have  $G$  groups, want to test whether they all have the same mean, need  $G-1$  indicator variables

$$E(Y_i) = \begin{cases} \beta_0 + \beta_1 & \text{if } i^{\text{th}} \text{ obs'n in 1st group} \\ \beta_0 + \beta_2 & \text{if } i^{\text{th}} \text{ obs'n in 2nd group} \\ \vdots & \vdots \\ \beta_0 + \beta_{G-1} & \text{if } i^{\text{th}} \text{ obs'n in } G^{\text{th}} \text{ group} \end{cases}$$

To test whether the means of all  $G$  groups are equal

$$H_0: \beta_1 = \beta_2 = \beta_3 = \dots = \beta_{G-1} = 0$$

Can show least squares estimates of  $\beta_0, \dots, \beta_5$  are

$$\overline{b_0} = \overline{y_6} \quad (= \overline{y_5}) \quad \text{average of observations}$$

$$b_1 = \overline{y_1} - \overline{y_6} \quad (= \overline{y_1} - \overline{y_5}) \quad (G=6 \text{ for our example})$$

etc.

$$\hat{y}_i = \text{predicted value of } Y_i$$

$$= \begin{cases} b_0 + b_1 & \text{1st observation in 1st group} \\ b_0 + b_2 & \text{" " " 2nd " "} \\ \vdots & \text{" " " " " "} \\ b_0 & \text{Gth group} \end{cases}$$

$$= \begin{cases} \bar{y}_1 \\ \bar{y}_2 \\ \vdots \\ \bar{y}_g \end{cases}$$

For our example

Test  $H_0: \beta_1 = \dots = \beta_r = 0$   
vs  $H_a$ : at least one not zero

Analysis of variance F-test:  $F_{obs} = 1.22$

$F_{obs}$  is an observation from  
F distribution with 5, 31 distribution

$\rightarrow$  # of being tested  
 $\leftarrow$  df error

No evidence of differences in means  
among other judges

Sums and squares and Analysis of Variance Table

Total SS (corrected SS)

$$SST = \sum_{i=1}^N (y_i - \bar{y})^2$$

$$N = n_1 + n_2 + \dots + n_g$$

$$df = N - 1$$

$$\text{Decomposition} = \underbrace{\sum_{i=1}^N (y_i - \hat{y}_i)^2}_{SS_{Res}} + \underbrace{\sum_{i=1}^N (y_i - \bar{y})^2}_{SS}$$

In 1-way analysis of variance: (predictor are indicator variables that classifying the observations one way)

$\hat{y}_i =$  mean of observations for group,  $g_i$  belongs to  $g$  in which the  $i$ th observation belongs to

$$SS_{Res} (Model SS) = \sum_{i=1}^N (y_i - \hat{y}_i)^2, \hat{y}_i \text{ is one of } g$$

$i$  is index for  
group  $g = 1, 2, \dots, G$

$$= \sum_{g=1}^G n_g (\bar{y}_g - \bar{y})^2 \quad \bar{y}_1, \bar{y}_2, \dots, \bar{y}_G$$

RSS (Residual  
or Error SS)

$$df = n - (G-1) - 1 \\ = n - G$$

$$= \sum_{i=1}^N (y_i - \hat{y}_i)^2 \\ = \sum_{g=1}^G \sum (y_i - \bar{y}_g)^2$$

$\sum_{(g)}$  indicates summation over  
observations in group  $g$

# ANOVA Table for 1-Way classification

<u>Source</u>	<u>df</u>	<u>SS</u>	<u>MS</u>	<u>F-stat</u>
<u>Group</u>	$g-1$	$SS_{\text{Reg}}$	$SS_{\text{Reg}} / (g-1)$	$\frac{SS_{\text{Reg}} / (g-1)}{MSE}$
<u>Error</u>	$N-g$	$PSS$	$PSS / (N-g) = MSE$	
<u>Total</u>	$N-1$	$SST$		

$MS = \text{"Mean Square"}$

In 1-way analysis of variance often called "between group"  $SS$  often called "within group"  $SS$

Idea of F-test: If between group variation is large compared to within group variation, there is evidence of differences in means among the groups.