

ORIGIN ≡ 0

## GLM Test of Linear Fit

W. Stein

The General Linear Models (GLM) Test using comparison between "full" and "reduced" statistical models allows one to formally argue whether a linear fit of the data is sufficient to describe the distribution of X,Y points in a dataset, or whether a more complex non-linear model (or transformed linear model) is required instead. This example is "creepy" because on first sight the "full" model doesn't *look* as constrained as the "reduced" model. However, one can tell that the "full" model is in fact the more constrained by calculating the Sum of Squares Error for Full Model SSEF, and seeing that it is less than Sum of Squares Error for Reduced model SSER. One can also compare the total number of parameters estimated for each model. A Full Model always has more than the corresponding Reduced Model. For the Full Model here, an expected level in dependent variable Y for each "bin" of independent variable X involves estimation of a separate parameter (mean  $\bar{Y}_j$  for each bin  $X_j$  - estimating  $\mu_j$  for each j). Thus, for the Full Model, a total of  $\max(i)=c$  parameters are estimated. The Reduced Model, by contrast, is fit using a single regression line, and this involves estimation of only 2 parameters ( $b_0$  &  $b_1$  estimating  $\beta_0$  &  $\beta_1$  respectively). So, this set up works when  $c > 2$ .

This example is also instructive in giving us an expanded meaning for the Null Hypothesis  $H_0$ . In many GLM Full vs Reduced model tests,  $H_0$  involves setting one or more parameters of a linear full model to zero. The two models neatly interconnect with "reduced" models comprising a subset of the "full" model. However, in this case Full versus Reduced Models have different forms. This points the generality and power of the GLM approach. Worked example drawn from Kuter et al. (KNNL) *Applied Linear Statistical Models* 5th Edition.

### Assumptions:

- Standard Linear Regression depends on specifying in advance which variable is to be considered 'dependent' and which 'independent'. This decision matters as changing roles for Y & X usually produces a different result.
- $Y_1, Y_2, Y_3, \dots, Y_n$  (dependent variable) is a random sample  $\sim N(\mu, \sigma^2)$ .
- $X_1, X_2, X_3, \dots, X_n$  (independent variable) with each value of  $X_i$  matched to  $Y_i$

Within this setup, two models for the relationship between X and Y variables are explicitly compared:

### Full Model:

$Y_{ij} = \mu_j + \varepsilon_{ij}$       where:  $\mu_j$  are parameter means at specified Bin levels of X. Bins are indicated by index variable j, with j having c levels with  $c > 2$   
 $\varepsilon_{ij}$  are "within" errors compared to each mean  $\mu_j \sim N(0, \sigma^2)$

### Reduced Model:

$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$       where:  $\beta_0$  is the y intercept of the regression line (translation),  
 $\beta_1$  is the slope of the regression line (scaling coefficient),  
 $\varepsilon_i$  is the error factor in prediction of  $Y_i$  and  
 a random variable  $\sim N(0, \sigma^2)$ .

### Example: Reading Bank Example KNNL Table 3.4:

```
K := READPRN("c:/2008LinearModelsData/Bank.txt")
```

^ The data was first sorted in Excel by ascending values of the first column (X) before importing.

K =

	0	1
0	75	28
1	75	42
2	100	112
3	100	136
4	125	160
5	125	150
6	150	152
7	175	156
8	175	124
9	200	124
10	200	104

**Assigning Variables and Calculating Simple Statistics:**

$$X := K^{(0)}$$

$$Y := K^{(1)}$$

$$n := \text{length}(X) \quad n = 11 \quad < \text{total observations}$$

**Setting up Bins for the Full Model:**

$$c := 6 \quad < \text{defined number of bins}$$

$$j := 0..c - 1 \quad < \text{range variable j for c bins}$$

$$X_{B_0} := \begin{pmatrix} X_0 \\ X_1 \end{pmatrix} \quad Y_{B_0} := \begin{pmatrix} Y_0 \\ Y_1 \end{pmatrix} \quad X_{B_0} = \begin{pmatrix} 75 \\ 75 \end{pmatrix} \quad Y_{B_0} = \begin{pmatrix} 28 \\ 42 \end{pmatrix}$$

$$X_{B_1} := \begin{pmatrix} X_2 \\ X_3 \end{pmatrix} \quad Y_{B_1} := \begin{pmatrix} Y_2 \\ Y_3 \end{pmatrix} \quad X_{B_1} = \begin{pmatrix} 100 \\ 100 \end{pmatrix} \quad Y_{B_1} = \begin{pmatrix} 112 \\ 136 \end{pmatrix}$$

$$X_{B_2} := \begin{pmatrix} X_4 \\ X_5 \end{pmatrix} \quad Y_{B_2} := \begin{pmatrix} Y_4 \\ Y_5 \end{pmatrix} \quad X_{B_2} = \begin{pmatrix} 125 \\ 125 \end{pmatrix} \quad Y_{B_2} = \begin{pmatrix} 160 \\ 150 \end{pmatrix}$$

$$X = \begin{array}{|c|c|} \hline & 0 \\ \hline 0 & 75 \\ \hline 1 & 75 \\ \hline 2 & 100 \\ \hline 3 & 100 \\ \hline 4 & 125 \\ \hline 5 & 125 \\ \hline 6 & 150 \\ \hline 7 & 175 \\ \hline 8 & 175 \\ \hline 9 & 200 \\ \hline 10 & 200 \\ \hline \end{array}$$

$$Y = \begin{array}{|c|c|} \hline & 0 \\ \hline 0 & 28 \\ \hline 1 & 42 \\ \hline 2 & 112 \\ \hline 3 & 136 \\ \hline 4 & 160 \\ \hline 5 & 150 \\ \hline 6 & 152 \\ \hline 7 & 156 \\ \hline 8 & 124 \\ \hline 9 & 124 \\ \hline 10 & 104 \\ \hline \end{array}$$

$$X_{B_3} := X_6 \quad Y_{B_3} := Y_6 \quad X_{B_3} = 150 \quad Y_{B_3} = 152$$

$$X_{B_4} := \begin{pmatrix} X_7 \\ X_8 \end{pmatrix} \quad Y_{B_4} := \begin{pmatrix} Y_7 \\ Y_8 \end{pmatrix} \quad Y_{B_4} = \begin{pmatrix} 156 \\ 124 \end{pmatrix} \quad Y_{B_4} = \begin{pmatrix} 156 \\ 124 \end{pmatrix}$$

$$X_{B_5} := \begin{pmatrix} X_9 \\ X_{10} \end{pmatrix} \quad Y_{B_5} := \begin{pmatrix} Y_9 \\ Y_{10} \end{pmatrix} \quad X_{B_5} = \begin{pmatrix} 200 \\ 200 \end{pmatrix} \quad Y_{B_5} = \begin{pmatrix} 124 \\ 104 \end{pmatrix}$$

< Note: this was done by hand. Statistical software typically have automated ways to efficiently handle the problem of binning.

**Calculating Sum of Squares Error for the Full Model:**

$$Y_{\text{bar}_j} := \text{mean}(Y_{B_j}) \quad Y_{\text{bar}} = \begin{pmatrix} 35 \\ 124 \\ 155 \\ 152 \\ 140 \\ 114 \end{pmatrix}$$

$$e_j := Y_{B_j} - Y_{\text{bar}_j} \quad \sum_j (Y_{B_j} - Y_{\text{bar}_j})^2 = \begin{pmatrix} 574 \\ 574 \end{pmatrix}$$

$$SSE_F := \sum \left[ \sum_j (e_j)^2 \right] \quad SSE_F = 1148$$

^ summing the square differences between  $Y_{B_j}$  and  $Y_{\text{bar}}$   
 Note that Mathcad had difficulty and reported a partial sum instead. So I added them together with the extra  $\Sigma$  here.  
 < Results confirmed p. 122.

## Performing Linear Regression for the Reduced Model:

### Initial calculations:

$$i := 0..n - 1 \quad < \text{range variable } i$$

$$X_{\text{bar}} := \text{mean}(X) \quad X_{\text{bar}} = 136.3636 \quad Y_{\text{bar}} := \text{mean}(Y) \quad Y_{\text{bar}} = 117.0909 \quad < \text{means for } X \text{ \& } Y$$

$$L_{xx} := \sum_i (X_i - X_{\text{bar}})^2 \quad L_{xx} = 21704.5455 \quad < \text{Corrected Sum of Squares for } X_i$$

$$L_{yy} := \sum_i (Y_i - Y_{\text{bar}})^2 \quad L_{yy} = 19882.9091 \quad < \text{Corrected Sum of Squares for } Y_i$$

$$L_{xy} := \sum_i (X_i - X_{\text{bar}}) \cdot (Y_i - Y_{\text{bar}}) \quad L_{xy} = 10563.6364 \quad < \text{Corrected Sum of Squares for Cross Product}$$

### Regression coefficients:

$$b_1 := \frac{L_{xy}}{L_{xx}} \quad b_1 = 0.4867 \quad < \text{sample estimate of slope } \beta_1$$

$$b_0 := Y_{\text{bar}} - b_1 \cdot X_{\text{bar}} \quad b_0 = 50.7225 \quad < \text{sample estimate of intercept } \beta_0$$

### Point Estimate of mean response (i.e., Regression line):

$$Y_{h_i} := b_0 + b_1 \cdot X_i \quad < \text{vector of points along the regression line.}$$

### Residuals:

$$e_i := Y_i - Y_{h_i} \quad < \text{vector of deviations of each value } Y_i \text{ from Regression line} = Y_{h_i}$$

## Calculating Sum of Squares Error for the Reduced Model:

$$SSE_R := \sum_i (e_i)^2 \quad SSE_R = 14741.5707 \quad < \text{Results confirmed p. 123}$$

## Hypotheses:

$H_0$ :  $E(Y) = \beta_0 + \beta_1 X$ , that is a Linear association between  $Y$  &  $X$

$H_1$ :  $E(Y)$  requires more specification than linear association such as  $\mu_j$  with  $j$  up to  $c > 2$

Note that Null Hypotheses are, in general, a formal statement of parsimony (i.e., "simplicity" of explanation). The null hypothesis says that the simpler of two alternatives is to be preferred unless the data require us to reject it. In a one population  $t$ -test of mean, for example, we ask whether an observed mean value  $X_{\text{bar}}$  is statistically indistinguishable from some specified value  $\mu_0$ . We normally interpret the Null Hypothesis  $H_0$  to say "the differences we observe between  $X_{\text{bar}}$  and  $\mu_0$  are the expected result of random behavior". However, random must always be defined in light of some model of what we expect for random, such as  $\sim N(\mu, \sigma)$ . We might more accurately claim instead that  $H_0$  says "unless compelled to do so, prefer the simpler hypothesis about difference between  $X_{\text{bar}}$  and  $\mu_0$ ", namely that there is nothing more to explain about the relationship between  $X_{\text{bar}}$  and  $\mu_0$  than  $\sim N(\mu, \sigma)$ . The Null Hypothesis for GLM above works exactly this way.

**Degrees of Freedom:**

$$df_F := n - c \quad df_F = 5$$

&lt; "full" model with c estimated parameters for bin means

$$df_R := n - 2 \quad df_R = 9$$

&lt; "reduced" model with 2 parameters for regression coefficients

**GLM Test Statistic:**

$$F := \frac{\frac{SSE_R - SSE_F}{df_R - df_F}}{\frac{SSE_F}{df_F}} \quad F = 14.8014 \quad < \text{results confirmed p. 124}$$

**Critical Value of the Test:**

$$\alpha := 0.01 \quad < \text{Probability of Type I error must be explicitly set}$$

$$CV := qF(1 - \alpha, c - 2, n - c) \quad CV = 11.3919 \quad < \text{note degrees of freedom here}$$

**Decision Rule:**

^ CV confirmed p. 124

IF  $F > CV$ , THEN REJECT  $H_0$  OTHERWISE ACCEPT  $H_0$ 

$$F = 14.8014$$

$$CV = 11.3919$$

**Probability Value:**

$$P := 1 - pF(F, c - 2, n - c) \quad P = 0.0056$$

&lt; results confirmed p. 124.

**Note on Bins:**

Although this test is formally designed to cover X,Y data with exact **replicates** in X, KNNL allow that sets of nearby X's with *approximately similar* Y's can be binned also. Results with these so-called **pseudoreplicates** should then be considered only approximate. The only problem I can see with pseudoreplicate binning is that perhaps one should be careful not to "cherry pick" sets of X's so as to force  $\bar{Y}_{\text{bar}}$ 's into a specific *a priori* pattern. Doing so necessarily biases conclusions against parsimony as described above, resulting in rejection of  $H_0$  perhaps more than one should. On the other hand, the whole question is whether a linear model based on  $H_0$  suffices to explain the data. If  $H_0$  suffices despite a concerted attempt to bias against it, perhaps one can consider the case well made.

**Prototype in R:**

```
# READ A STRUCTURED DATA TABLE
K=read.table("c:/2008LinearModelsData/bankR.txt")
K
X=K$deposit
Y=K$newacc
B=factor(K$bin)

# FINDING N & C
n=length(X)
n
c=nlevels(B)
c
```

```
K
deposit newacc bin
1      75      28  1
2      75      42  1
3     100     112  2
4     100     136  2
5     125     160  3
6     125     150  3
7     150     152  4
8     175     156  5
9     175     124  5
10    200     124  6
11    200     104  6
```

n = 11      c = 6

```
# ANOVA OF FITTED MODELS
anova(lm(Y~B)) #FULL MODEL
anova(lm(Y~X)) #REDUCED MODEL
```

```
# CALCULATING SUM OF SQUARES ERROR
# FOR FULL MODEL
MF=summary(lm(Y~B),digits=10)
MSEF=MF$sigma^2
dfF=MF$df[2]
SSEF=dfF*MSEF
SSEF
```

```
# CALCULATING SUM OF SQUARES ERROR
# FOR REDUCED MODEL
MR=summary(lm(Y~X),digits=10)
MSER=MR$sigma^2
dfR=MR$df[2]
SSER=dfR*MSER
SSER
```

```
# CALCULATING GLM F STATISTIC
F=((SSER-SSEF)/(dfR-dfF))/(SSEF/dfF)
F
# FINDING CRITICAL VALUE
alpha=0.01
CV=qf(1-alpha,c-2,n-c)
CV
# PROBABILITY VALUE
P=1-pf(F,c-2,n-c)
P
```

```
# THE EFFICIENT WAY TO DO THIS
# SPECIFY FULL VS REDUCED MODELS:
FM=lm(Y~B)
RM=lm(Y~X)
# CALCULATE ANOVA TABLE OF COMPARISON
anova(RM,FM)
```

```
> anova(lm(Y~B)) #FULL MODEL
Analysis of Variance Table
```

```
Response: Y
      Df Sum Sq Mean Sq F value Pr(>F)
B      5 18734.9  3747.0  16.320 0.004085 **
Residuals 5  1148.0   229.6
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
> anova(lm(Y~X)) #REDUCED MODEL
Analysis of Variance Table
```

```
Response: Y
      Df Sum Sq Mean Sq F value Pr(>F)
X      1  5141.3  5141.3   3.1389 0.1102
Residuals 9 14741.6  1638.0
```

```
> SSEF
```

```
[1] 1148
```

```
> SSER
```

```
[1] 14741.57
```

```
> F
```

```
[1] 14.80136
```

```
> CV
```

```
[1] 11.39193
```

```
> P
```

```
[1] 0.005593812
```

```
> anova(RM,FM)
```

```
Analysis of Variance Table
```

```
Model 1: Y ~ X
```

```
Model 2: Y ~ B
```

```
  Res.Df  RSS Df Sum of Sq    F    Pr(>F)
1      9 14742
2      5  1148  4     13594 14.801 0.005594 **
```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

^ the parsimonious model is rejected...