

Elliot Smith  
STAT616 - HW #2

Problem 1

a. Effects Model:  $y_{ij} = \mu + \alpha_i + \epsilon_{ij}$

Using Dummy Coding:

$$y \begin{matrix} M \\ F \end{matrix} \begin{matrix} \mu & \alpha_2 \\ \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \end{matrix} \quad \alpha_1 = 0$$

So:  $y = \mu + \alpha_2 x_2$

$\rightarrow \hat{\mu} = \bar{y}_{..}$

$\rightarrow \hat{\alpha}_i = \bar{y}_{i.} - \hat{\mu} = \bar{y}_{i.} - \bar{y}_{..}$

b. Effects Model - Dummy Coding

$$y \begin{matrix} M \\ F \end{matrix} \begin{matrix} \mu & \alpha_2 \\ \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \end{matrix}$$

$\hat{\alpha}_1 = \bar{y}_{1.} - \hat{\mu} = \bar{y}_{1.} - \frac{1}{2}(\bar{y}_{1.} + \bar{y}_{2.})$

$\hat{\alpha}_2 = \bar{y}_{2.} - \hat{\mu} = \bar{y}_{2.} - \frac{1}{2}(\bar{y}_{1.} + \bar{y}_{2.})$

$\hat{\mu} = \frac{1}{2}(\bar{y}_{1.} + \bar{y}_{2.})$

$\hat{\sigma}_{LS}^2 = \frac{\sum_i \sum_j (y_{ij} - \bar{y}_{i.})^2}{N - k}$

$k = \# \text{ of treatments}$

$N = n_1 + n_2$

$SE(\text{intercept}) \rightarrow \text{Var}(\text{intercept}) = \text{Var}(\hat{\mu}) = \text{Var}(\bar{y}_{..}) = \frac{\sigma^2}{n}$

So:  $\hat{SE}(\text{intercept}) = \frac{\hat{\sigma}}{\sqrt{n}}$

$SE(\hat{\alpha}_1) \rightarrow \text{Var}(\hat{\alpha}_1) = \bar{y}_{1.} - \bar{y}_{..} = \text{Var}(\bar{y}_{1.}) + \text{Var}(\bar{y}_{..}) = \frac{\sigma^2}{n_1} + \frac{\sigma^2}{n} = \frac{2\sigma^2}{n}$

So:  $SE(\hat{\alpha}_1) = \frac{\sqrt{2}}{\sqrt{n}} \hat{\sigma}$

$SE(\hat{\alpha}_2) \rightarrow \text{Var}(\hat{\alpha}_2) = \bar{y}_{2.} - \bar{y}_{..} = \text{Var}(\bar{y}_{2.}) + \text{Var}(\bar{y}_{..}) = \frac{\sigma^2}{n_2} + \frac{\sigma^2}{n} = \frac{2\sigma^2}{n}$

So:  $SE(\hat{\alpha}_2) = \frac{\sqrt{2}}{\sqrt{n}} \hat{\sigma}$

Since a model is defined as  $E[y|x_2]$  (since  $\mu$  takes on  $x_1$  when  $x_1=0$ ) we should examine  $n_2$  fitted values.

### c. Effects Model - Zero Sum Coding

$$Y = \mu + \alpha_1 X_1 - \alpha_2 X_2$$

$$\alpha_1 + \alpha_2 = 0$$

$$\alpha_1 = -\alpha_2$$

$$\text{So: } y = \mu + \alpha_1 x_1 - \alpha_2 x_2$$

$$\rightarrow \hat{\mu} = \frac{1}{2}(\bar{y}_{1..} + \bar{y}_{2..})$$

$$\rightarrow \hat{\alpha}_1 = \bar{y}_{1..} - \hat{\mu} = \bar{y}_{1..} - \frac{1}{2}(\bar{y}_{1..} + \bar{y}_{2..})$$

$$\rightarrow \hat{\alpha}_2 = \bar{y}_{2..} - \hat{\mu} = \bar{y}_{2..} - \frac{1}{2}(\bar{y}_{1..} + \bar{y}_{2..})$$

$$\hat{\sigma}_\epsilon^2 = \frac{\sum_{i \in j} (y_{ij} - \bar{y}_{i..})^2}{N-k}$$

$$k = \# \text{ of treatments}$$

$$N = n_1 + n_2$$

$$\hat{\mu} = \bar{y}_{..} \quad SE(\hat{\mu}) = \sqrt{\frac{\sigma^2}{N}} = \sqrt{\frac{\sigma^2}{kn}}$$

$$Var(\hat{\mu}) = Var(\bar{y}_{..}) = \frac{\sigma^2}{N}$$

$$\hat{\alpha}_1 = \bar{y}_{1..} - \bar{y}_{..} \quad SE(\hat{\alpha}_1) = \sqrt{\frac{2}{3} \frac{\sigma^2}{n}}$$

$$Var(\hat{\alpha}_1) = Var(\bar{y}_{1..} - \bar{y}_{..}) = Var(\bar{y}_{1..}) + Var(\bar{y}_{..}) - 2Cov(\bar{y}_{1..}, \bar{y}_{..})$$

$$= \frac{\sigma^2}{n} + \frac{\sigma^2}{N} - \frac{2}{3} \frac{\sigma^2}{n}$$

Since a model is defined as  $E[Y|X_1, X_2]$  we should have exactly  $n_1$  values if  $n_1 < n_2$ .

- d. The estimators are different because of the coding. For example, with dummy coding, we establish a baseline  $\beta$ , meaning that we set this estimator to zero and interpret all results based on this "new zero". On the contrary, with zero-sum coding, we set a different baseline, with the sum of all the estimators equal to zero and thus establishing a different "new zero". In summary, different coding methods decide on a particular baseline, which determines the value of the estimators.

e. To find the MLEs, perform the following steps:

- ① Find the joint probability function of the data
- ② Simplify with logs
- ③ Maximize function with respect to  $\mu_i, \sigma^2$  ( $\frac{\partial l}{\partial \mu_i} = 0, \frac{\partial l}{\partial \sigma^2} = 0$ )

The MLEs come out to:

$$\hat{\mu}_i = \bar{y}_i$$

$$\hat{\sigma}_{ML}^2 = \frac{\sum_i \sum_j (y_{ij} - \bar{y}_i)^2}{N}$$

$$y_{ij} = \mu_i + \varepsilon_{ij}, \varepsilon_{ij} \sim N(0, \sigma^2)$$

$$\textcircled{1} \text{ pdf: } f(y; \mu_1, \mu_2, \sigma^2) = \prod_{ij} f(y_{ij}; \mu_i, \sigma^2)$$

$$= \prod_{ij} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_{ij} - \mu_i)^2}$$

$$\textcircled{2} f(y) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^N e^{-\frac{1}{2\sigma^2} \sum_{ij} (y_{ij} - \mu_i)^2}$$

$$\log[f(y)] = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{ij} (y_{ij} - \mu_i)^2$$

$$\textcircled{3} \arg \max_{\mu_i, \sigma^2} l(\mu_i, \sigma^2; y)$$

$$\frac{\partial l}{\partial \mu_i} = 0, \frac{\partial l}{\partial \sigma^2} = 0$$

$$\hat{\sigma}_{LS}^2 = \frac{\sum_i \sum_j (y_{ij} - \bar{y}_i)^2}{N - k} \rightarrow \text{unbiased} \quad E[\hat{\sigma}_{LS}^2] = \sigma^2$$

\*  $k = \#$  of  $\mu$ 's

$$\hat{\sigma}_{ML}^2 = \frac{\sum_i \sum_j (y_{ij} - \bar{y}_i)^2}{N} \rightarrow \text{biased} \quad E[\hat{\sigma}_{ML}^2] \neq \sigma^2$$

As  $N \gg k$ , these will converge!

## Problem 2

a.

<u>Source</u>	<u>SS</u>	<u>df</u>	<u>MS</u>	<u>F-ratio</u>
Between (SSB)	258	4	64.5	10.134
Within (SSW)	76	12	6.333	N/A
Total (SST)	334	16		

We may conclude that  $\frac{258}{334} = \sim 77\%$  of the variance is explained by the group-to-group variance.

Since the F-statistic is significantly greater than 1, we may conclude that at least one of the population means differs from the others.

Based on only the F-ratio, we cannot say if the variation is statistically significant. We will first want to design a hypothesis test and calculate a p-value to test our hypothesis of all equal means.

- b. I tested my null hypothesis of equal means and got a p-value of 0.0008 and thus rejected my null hypothesis. My methodology was to randomly permute the labels of the data and then generate an F-statistic for each run through (I performed 5000). To calculate my p-value I found the number of generated F-statistics that were greater or equal to my observed F-statistic (there were 4). I then divided this number by the total number of generated F-statistics (5000) to obtain a p-value of 0.0008 and rejected my null hypothesis.

\*See code appendix

- c. Please see my code appendix for my construction of the model and the accompanying plots.

Based on my QQ Plot and Histogram of the residuals, I conclude that the errors are mostly normal. The QQ Plot shows that the residuals mostly follow the line, telling us that they are normally distributed. The histogram as well shows us that the residuals are mostly normal as the middle value is tallest and the fringe values tend to be less represented as we approach the extremes.

Based on my test for equal variances, I conclude this test passes as well. As we can see, the residual values are equally and fairly distributed around the zero-axis line, showing that the variances are equally distributed and mostly equal.

I conclude that all tests pass.

- d. The ANOVA table shows us the sources and contribution of variance and whether or not, through hypothesis testing using the F-statistic, all of the means are equal; it does NOT tell us which groups are different. To answer this question, we will fit an ANOVA model using regression. So in summary, the ANOVA task will tell us if all means are equal and then our ANOVA model via regression may tell us which means are not equal.

- e. To compare groups B and C and B and E, we will use a t-test. Please see my code appendix for R calculations.

#### B and C

t-statistic: 0.7071

df = 5

p-value = 0.5123

I will NOT reject  $H_0: \mu_B = \mu_C$  in favor of  $H_1: \mu_B > \mu_C$

#### B and E

t-statistic: -2.1669

df = 5

p-value = 0.0983

I will NOT reject  $H_0: \mu_B = \mu_E$  in favor of  $H_1: \mu_B > \mu_E$

STAT 616

### Problem 1

After some trial and error, the lowest possible p-value that I could obtain was  $\approx 0.02$  (see code appendix for values chosen). I am not entirely sure what this suggests, however, a fair consideration might be that, with a certain confidence level, we may not reject  $H_0$  with a p-value of 0.02 (for example if  $\alpha = 0.01$ ). So we may say that having only 2 observations for each of the 3 treatments is not enough and we must add more for better results.

1. ...

# Code Appendix

*Elliot Smith*

*2/13/2018*

```
options(scipen = 999)
```

## Problem 2

### Part b

```
# Problem 2

## Part b

values <- c(1,3,5,9,5,5,5,6,6,3,3,3,0,6,14,10,18)
labels <- c("A","A","A","B","B","B","B","C","C","C","D","D","D","D","E","E","E")

data <- data.frame(labels, values)

aov_out <- aov(values ~ labels, data = data)
summary(aov_out)

##              Df Sum Sq Mean Sq F value    Pr(>F)
## labels         4     258    64.50   10.18 0.000782 ***
## Residuals     12       76     6.33
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

f_stat <- as.vector(summary.lm(aov_out)$fstatistic["value"])

f_stats <- numeric()
n <- 5000

for (i in 1:n) {

  labels_temp <- sample(labels)
  data_temp <- data.frame(labels_temp, values)
  aov_out_temp <- aov(values ~ labels_temp, data = data_temp)
  f_stats[i] <- as.vector(summary.lm(aov_out_temp)$fstatistic["value"])

}

count_greater <- sum(f_stats >= f_stat)

count_greater/n

## [1] 0.0004
```

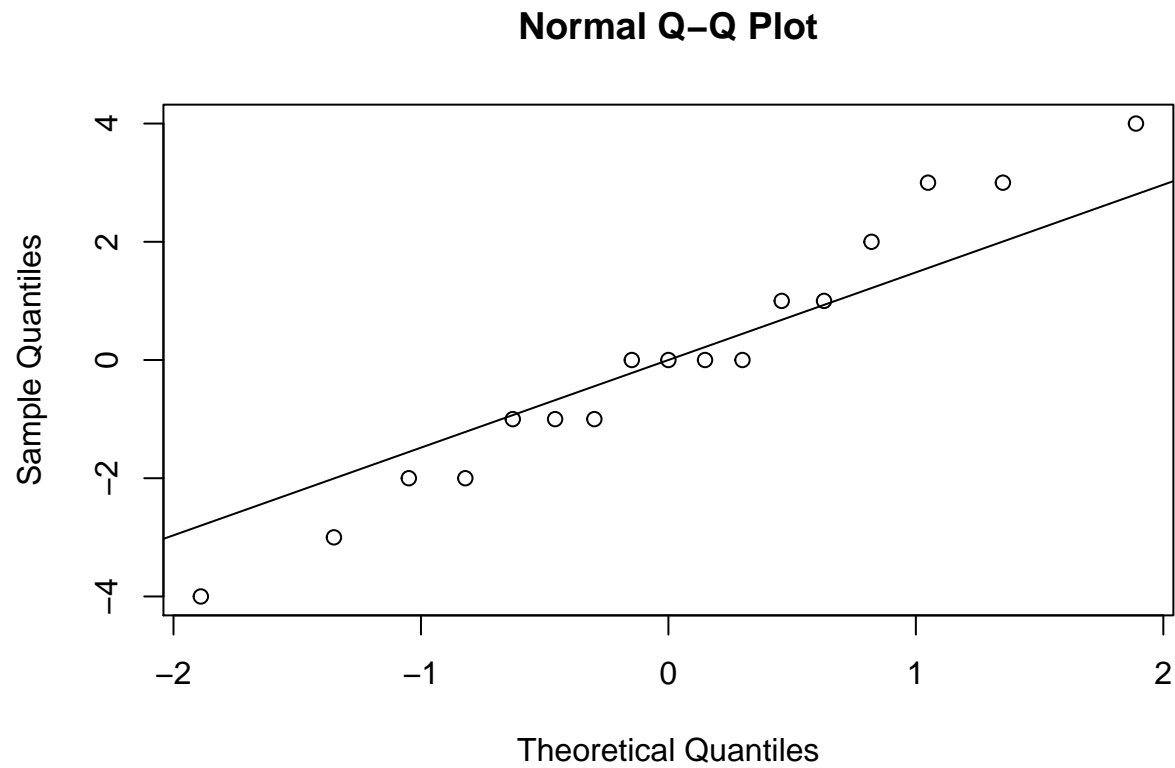


## Part c

```
## Part c

### Check for error normality

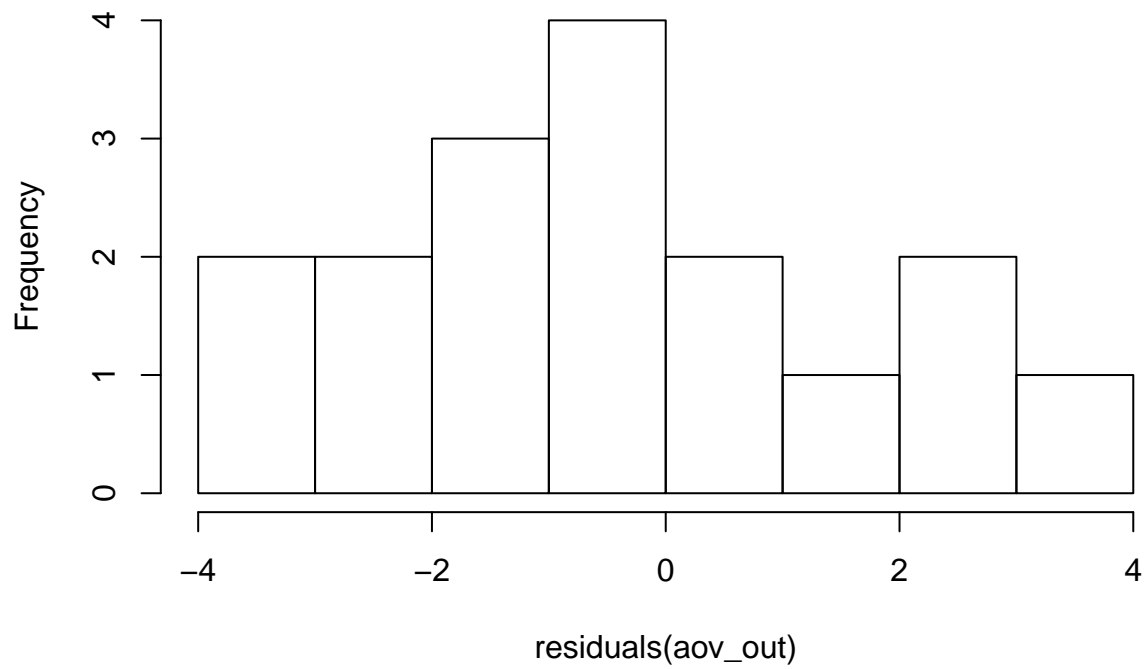
qqnorm(residuals(aov_out))
qqline(residuals(aov_out))
```



```
hist(residuals(aov_out))
```

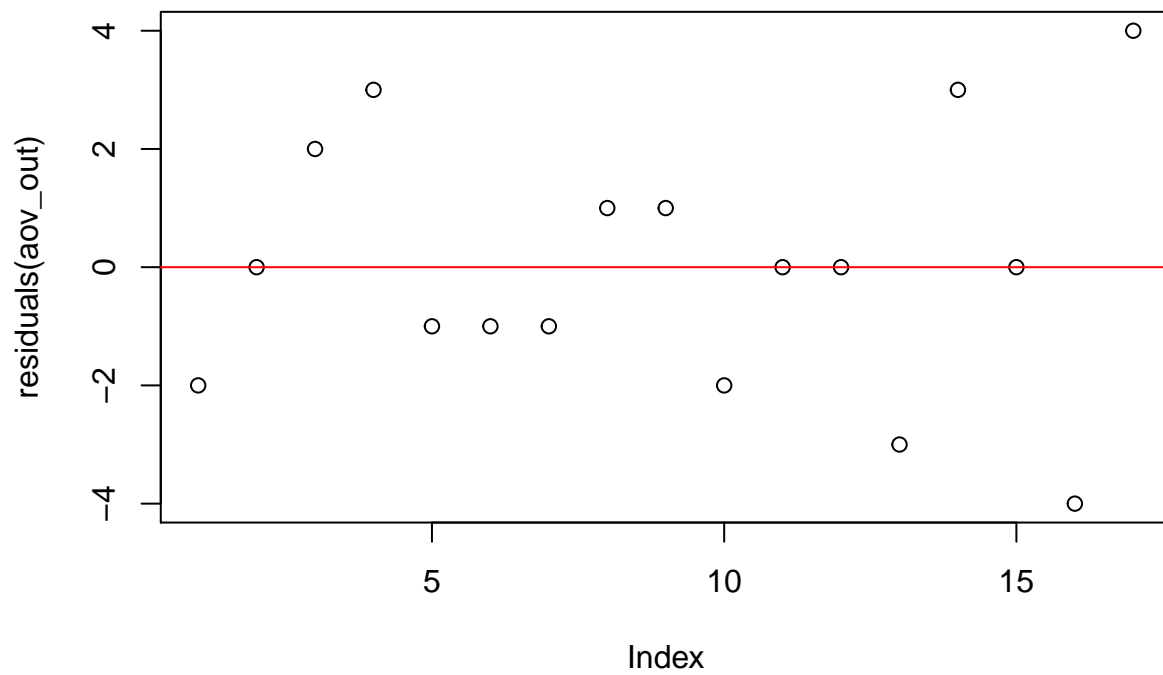


**Histogram of residuals(aov\_out)**



```
### Check residuals for equal variance
```

```
plot(residuals(aov_out))  
abline(h = 0, col = "red")
```



## Part e

```
## Part e

### Compare B and C

t.test(x = data[4:7,2], y = data[8:10,2])

##
## Welch Two Sample t-test
##
## data: data[4:7, 2] and data[8:10, 2]
## t = 0.70711, df = 4.8, p-value = 0.5123
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -2.681387 4.681387
## sample estimates:
## mean of x mean of y
##      6      5

### Compare B and E

t.test(x = data[4:7,2], y = data[14:17,2])

##
## Welch Two Sample t-test
##
## data: data[4:7, 2] and data[14:17, 2]
## t = -2.1669, df = 3.8802, p-value = 0.09825
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -13.78215 1.78215
## sample estimates:
## mean of x mean of y
##      6      12
```

## STAT616 - Problem 1

```
# STAT616 - Problem 1

values <- c(1,2,3,4,60000001,60000000)
labels <- c("A","A","B","B","C","C")

data <- data.frame(labels, values)

aov_out <- aov(values ~ labels, data = data)
summary(aov_out)
```

	Df	Sum Sq	Mean Sq	F value
## labels	2	47999999680000010	23999999840000005	47999999696203265
## Residuals	3	1	0	
##		Pr(>F)		
## labels		<0.0000000000000002	***	

```

## Residuals
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

f_stat <- as.vector(summary.lm(aov_out)$fstatistic["value"])

f_stats <- numeric()
n <- 5000

for (i in 1:n) {

  labels_temp <- sample(labels)
  data_temp <- data.frame(labels_temp, values)
  aov_out_temp <- aov(values ~ labels_temp, data = data_temp)
  f_stats[i] <- as.vector(summary.lm(aov_out_temp)$fstatistic["value"])

}

count_greater <- sum(f_stats >= f_stat)

count_greater/n

## [1] 0.024

```