

## 3. Oneway ANOVA

Suppose  $Y_{i1}, \dots, Y_{in_i} \stackrel{\text{iid}}{\sim} N(\mu_i, \sigma^2)$  for  $i = 1, \dots, k$  groups. Let  $\bar{Y}_i$  be the group  $i$  average and let  $\bar{Y}$  be the overall average. Consider a test of  $H_0: \mu_1 = \dots = \mu_k$  vs.  $H_a: \text{"not } H_0"$ .

a) Explain how the  $F$  statistic generalizes the pooled 2-sample  $t$  statistic.

$$F\text{-stat} = \frac{SSM/(k-1)}{SSE/(n-k)} = \frac{\sum_{i=1}^k n_i (\bar{Y}_i - \bar{Y})^2 / (k-1)}{\sum_{i=1}^k (n_i - 1) s_i^2 / (n-k)}$$

For example, show how rejecting for large values of this is the same as rejecting for large  $T^2$  when  $k = 2$  and  $N = n_1 + n_2$  (what would unusually small values of  $T^2$  suggest?).

a)  $Y_{i1}, \dots, Y_{in_i} \stackrel{\text{iid}}{\sim} N(\mu_i, \sigma^2)$   $i=1 \dots k$  groups  $\bar{Y}_i$  = group  $i$  average  $\bar{Y}$  = overall average

$H_0: \mu_1 = \dots = \mu_k$   $H_a$ : not  $H_0$

$F$ -Stat generalizes the pooled  $Z$ -sample  $t$ -stat:

$$F\text{-stat} = \frac{SSM/(k-1)}{SSE/(n-k)} = \frac{\sum_{i=1}^k n_i (\bar{Y}_i - \bar{Y})^2 / (k-1)}{\sum_{i=1}^k (n_i - 1) s_i^2 / (n-k)}$$

$\rightarrow$  variation b/w group and overall  $\rightarrow$  variation in each group

$Z$ -sample  $t$ -test:

$$T = \frac{|\bar{Y}_1 - \bar{Y}_2|}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

When  $k=2$  and  $N = n_1 + n_2$

$$F = \frac{SSM/(2-1)}{SSE/(n-2)} = \frac{SSM}{SSE/(n-2)} = \frac{\sum n_i (\bar{Y}_i - \bar{Y})^2}{\sum s_i^2 (n_i - 1) / (n-2)} = \frac{n_1 (\bar{Y}_1 - \bar{Y})^2 + n_2 (\bar{Y}_2 - \bar{Y})^2}{\frac{s_1^2 (n_1 - 1) + s_2^2 (n_2 - 1)}{(n_1 + n_2 - 2)}}$$

$$\bar{Y} = \text{overall mean} = \frac{n_1 \bar{Y}_1 + n_2 \bar{Y}_2}{N}$$

numerator

$$\begin{aligned} & n_1 \left[ \bar{Y}_1 - \left( \frac{n_1 \bar{Y}_1 + n_2 \bar{Y}_2}{N} \right) \right]^2 + n_2 \left[ \bar{Y}_2 - \left( \frac{n_1 \bar{Y}_1 + n_2 \bar{Y}_2}{N} \right) \right]^2 = n_1 \left[ \frac{N \bar{Y}_1 - n_1 \bar{Y}_1 - n_2 \bar{Y}_2}{N} \right]^2 + n_2 \left[ \frac{N \bar{Y}_2 - n_1 \bar{Y}_1 - n_2 \bar{Y}_2}{N} \right]^2 \\ & = n_1 \left[ \frac{(N - n_1) \bar{Y}_1 - n_2 \bar{Y}_2}{N} \right]^2 + n_2 \left[ \frac{(N - n_2) \bar{Y}_2 - n_1 \bar{Y}_1}{N} \right]^2 = n_1 \left[ \frac{n_2 \bar{Y}_1 - n_2 \bar{Y}_2}{N} \right]^2 + n_2 \left[ \frac{n_1 \bar{Y}_2 - n_1 \bar{Y}_1}{N} \right]^2 \\ & = \frac{n_1 n_2^2}{N^2} (\bar{Y}_1 - \bar{Y}_2)^2 + \frac{n_2 n_1^2}{N^2} (\bar{Y}_2 - \bar{Y}_1)^2 = \left( \frac{n_1 n_2^2}{N^2} + \frac{n_2 n_1^2}{N^2} \right) (\bar{Y}_1 - \bar{Y}_2)^2 = \frac{n_1 n_2 (n_1 + n_2)}{N^2} (\bar{Y}_1 - \bar{Y}_2)^2 = \frac{n_1 n_2}{N} (\bar{Y}_1 - \bar{Y}_2)^2 \\ & = \frac{\frac{1}{N}}{\frac{1}{n_1} + \frac{1}{n_2}} (\bar{Y}_1 - \bar{Y}_2)^2 = \frac{1}{\frac{1}{n_1} + \frac{1}{n_2}} (\bar{Y}_1 - \bar{Y}_2)^2 \end{aligned}$$

$$\text{numerator} = \frac{(\bar{Y}_1 - \bar{Y}_2)^2}{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$F = \frac{(\bar{Y}_1 - \bar{Y}_2)^2}{\left( \frac{1}{n_1} + \frac{1}{n_2} \right) \left( \frac{s_1^2 (n_1 - 1) + s_2^2 (n_2 - 1)}{(n_1 + n_2 - 2)} \right)}$$

Thus  $F = T^2$  for  $k=2$

You reject for large  $F$ . Indicating that variation b/w groups is larger than

variation within groups. Since  $F = T^2$  for  $k=2$

then reject for large  $T^2$  too. Small  $T^2$  suggest

difference b/w groups is small compared to within groups

$$T^2 = \left( \frac{\bar{Y}_1 - \bar{Y}_2}{\text{Sp} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \right)^2 = \frac{(\bar{Y}_1 - \bar{Y}_2)^2}{\text{Sp}^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)} = \frac{(\bar{Y}_1 - \bar{Y}_2)^2}{\frac{(n_1 - 1) s_1^2 + (n_2 - 1) s_2^2}{(n_1 + n_2 - 2)} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

$\hookrightarrow$  since pooled  $s_1 = s_2 = \text{Sp}$

pooled variance formula

If  $n_1 = n_2 = n$  it is much more simple!

$k$  groups,  $n_i$  individuals,  $Y_i$  = group  $i$  ave,  $Y_{ij}$  = individual  $j$  in group  $i$

$$SSM = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2 = \sum_{i=1}^k n_i (Y_i - \bar{Y})^2$$

when  $k=2$ :

$$\bar{Y} = \frac{n_1 Y_1 + n_2 Y_2}{n_1 + n_2} = \frac{n(Y_1 + Y_2)}{2n} = \frac{1}{2}(Y_1 + Y_2)$$

$$\text{Then: } Y_1 - \bar{Y} = \frac{1}{2}(Y_1 - Y_2)$$

$$Y_2 - \bar{Y} = \frac{1}{2}(Y_2 - Y_1)$$

$$\text{so, } \sum_{i=1}^2 n(Y_i - \bar{Y})^2 = \frac{n}{2^2}(Y_1 - Y_2)^2 + \frac{n}{2^2}(Y_2 - Y_1)^2 = (Y_1 - Y_2)^2 \left(\frac{n}{2} + \frac{n}{2}\right)$$

$$\text{Thus } SSM = \frac{n}{2}(Y_1 - Y_2)^2$$

$$T\text{-stat} = \frac{Y_1 - Y_2}{\text{Sp} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} = \frac{Y_1 - Y_2}{\text{Sp} \sqrt{\frac{2}{n}}}$$

$$T^2 = \frac{n(Y_1 - Y_2)^2}{2 \text{Sp}^2} \quad \text{where } \text{Sp}^2 = \text{MSE} = \text{SSE} / (N - k)$$

$$F = \frac{SSM / (2 - 1)}{\text{Sp}^2} = \frac{\frac{n}{2}(Y_1 - Y_2)^2 / (2 - 1)}{\text{Sp}^2} = \frac{n(Y_1 - Y_2)^2}{2 \text{Sp}^2}$$

Therefore when  $k=2$   $F$  stat and  $T^2$  are the same.

- b) To simplify, suppose  $n_i = n$ , for  $i = 1, \dots, k$ , so  $N = nk$  and  $\bar{Y} = \frac{1}{k} \sum_{i=1}^k Y_i$ . Show, using facts we have already proved, that  $\frac{\text{SSE}}{\sigma^2} = \sum_{i=1}^k \sum_{j=1}^n \frac{(Y_{ij} - Y_i)^2}{\sigma^2} \sim \chi^2_{(N-k)}$  and independent of SSM, and that  $\frac{\text{SST}}{\sigma^2} = \sum_{i=1}^k \sum_{j=1}^{n_i} \frac{(Y_{ij} - \bar{Y})^2}{\sigma^2} \sim \chi^2_{(N-1)}$  if  $H_0$  is true. Also show  $\text{SST} = \text{SSM} + \text{SSE}$ , and infer that  $\frac{\text{SSM}}{\sigma^2} \sim \chi^2_{(k-1)}$  if  $H_0$  is true. Conclude that the null sampling distribution of the  $F$ -statistic is  $F_{(k-1, N-k)}$ . Note that  $\text{MSE} = \text{SSE} / (N - k)$  is an unbiased estimate for  $\sigma^2$ . If  $H_0$  is true, then  $\text{MSM} = \text{SSM} / (k - 1)$  is also unbiased for  $\sigma^2$ , but otherwise is biased high.

Step 1: Show that  $\frac{SSE}{\sigma^2} = \sum_{i=1}^k \sum_{j=1}^{n_i} \frac{(Y_{ij} - \bar{Y}_i)^2}{\sigma^2} \sim \chi^2_{(N-k)}$

Note that we have already shown that:

$$Y_{11}, \dots, Y_{1n_1} \stackrel{iid}{\sim} N(\mu_1, \sigma^2)$$

$$S^2 = \frac{\sum_{j=1}^n (Y_{ij} - \bar{Y}_i)^2}{n-1}$$

$Y_{ij} - \bar{Y}_i$  independent of  $\bar{Y}_i$

$$S^2 = \frac{\sum_{j=1}^n (Y_{ij} - \bar{Y}_i)^2}{n-1} \text{ independent of } \bar{Y}_i$$

Also seen that:  $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{(n-1)}$  so  $S^2 \sim \text{Gamma}(\frac{n-1}{2}, \frac{n-1}{2\sigma^2})$

So since  $Y_{ij} \sim N(\mu_i, \sigma^2)$  then

$$Y_{ij} - \bar{Y}_i \sim N(0, \sigma^2) \text{ (independently for each } j)$$

And we know that  $\frac{(n_i-1)S_i^2}{\sigma^2} \stackrel{iid}{\sim} \chi^2_{(n_i-1)}$

$$\text{Thus } \frac{SSE}{\sigma^2} = \sum_{i=1}^k \frac{(n_i-1)S_i^2}{\sigma^2} = \sum_{i=1}^k \frac{\sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2}{\sigma^2} \sim \chi^2_{(N-k)} \text{ where } N = \sum_{i=1}^k n_i$$

summing  $k$  independent groups so add degrees of freedom  $n_1 + n_2 + \dots + n_k - k = N - k$

Show that  $MSE = SSE / (N-k)$  is unbiased for  $\sigma^2$

$$MSE = \frac{SSE}{N-k}$$

$$E(MSE) = E\left(\frac{SSE}{N-k}\right) = \frac{(N-k)\sigma^2}{(N-k)} = \sigma^2$$

Also note that  $MSE \sim \text{Gamma}\left(\frac{N-k}{2}, \frac{N-k}{2\sigma^2}\right)$  and  $E(MSE) = \frac{N-k}{2} \left(\frac{2\sigma^2}{N-k}\right) = \sigma^2$

So unbiased estimate.

SSM and SSE are independent.

- we know that  $Y_{ij} \sim N(\mu_i, \sigma^2)$  thus  $(Y_{ij} - \bar{Y}_i)$  or  $S_i^2$  is independent of  $\bar{Y}_i$

-  $\bar{Y}_i$  are calculated from  $Y_{ij}$ 's but  $(Y_{ij} - \bar{Y}_i)$  is independent of  $\bar{Y}_i$ .

So, SSM depends on group means  $\bar{Y}_i$ , but SSE depends on deviations from the group means

We already know that  $(Y_{ij} - \bar{Y}_i)$  is independent of  $\bar{Y}_i$  so SSE and SSM are indep.

Show that  $\frac{SST}{\sigma^2} = \sum_{i=1}^k \frac{\sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2}{\sigma^2} \sim \chi^2_{(N-1)}$  under  $H_0$

Note: Under  $H_0$ :  $\mu_1 = \mu_2 = \mu_3 = \dots = \mu_k = \mu$  so  $Y_{ij} \stackrel{iid}{\sim} N(\mu, \sigma^2)$

$$SST = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$$

$$\text{Each } (Y_{ij} - \mu)^2 \sim N(0, \sigma^2)$$

Summing standard normals is chi-square (proved at earlier date)

$$\text{So } \frac{SST}{\sigma^2} \sim \chi^2_{(N-1)}$$

It is  $N-1$  because  $\bar{Y}$  is estimated

Since  $SST = SSM + SSE$  and  $SSE$  and  $SSM$  are independent.

$$\frac{SST}{\sigma^2} = \frac{SSM}{\sigma^2} + \frac{SSE}{\sigma^2}$$

$$\chi^2_{(N-1)} = ? + \chi^2_{(N-k)} \quad - \text{DF add}$$

$$\text{So } \frac{SSM}{\sigma^2} = \chi^2_{(k-1)} \text{ only if } H_0 \text{ is true}$$

Conclude null sampling of f-stat is  $F(k-1, N-k)$

$$F\text{-stat} = \frac{SSM/(k-1)}{SSE/(N-k)} = \frac{MSM}{MSE}$$

$$\text{under null } \frac{SSE}{\sigma^2} \sim \chi^2_{N-k} \quad \frac{SSM}{\sigma^2} \sim \chi^2_{k-1}$$

The ratio of two independent  $\chi^2$  variables follows an F distribution

thus  $F(k-1, N-k)$

Note:  $X_1 \sim \text{Gamma}(a, \lambda)$  independent  $X_2 \sim \text{Gamma}(b, \lambda)$

Here  $X_1 = SSM \sim \text{Gamma}(\frac{k-1}{2}, \frac{1}{2\sigma^2})$  and  $X_2 = SSE \sim \text{Gamma}(\frac{N-k}{2}, \frac{1}{2\sigma^2})$

$$\frac{X_1}{X_2} \sim F^*(a, b, c)$$

$$\frac{X_1/a}{X_2/b} \sim F(2a, 2b) \quad \text{So } \frac{SSM/(k-1)}{SSE/(N-k)} \sim F(k-1, N-k)$$

MSM Biases:

$$MSM = SSM/(k-1)$$

under  $H_0$ :  $E(MSM) = E\left(\frac{SSM}{k-1}\right) = \frac{E(SSM)}{k-1} = \frac{\sigma^2(k-1)}{k-1} = \sigma^2$  because  $SSM \sim \chi^2_{(k-1)}$  under  $H_0$

Otherwise:  $E(SSM)$  is greater than  $\sigma^2(k-1)$  because variation between groups will no longer be due to pure random sampling.

$$\text{Thus } E(MSM) > \sigma^2$$

Also under  $H_0$ :  $Y_i \sim N(\mu, \frac{\sigma^2}{n})$

$$E\left(\frac{\sum (Y_i - \bar{Y})^2}{k-1}\right) = \frac{\sigma^2}{n}$$

✓ c) Define  $R^2 = \frac{SSM}{SST} = 1 - \frac{SSE}{SST}$  as the proportion of variability explained by the groups. Show that the null sampling distribution of  $R^2$  is Beta( $k-1, N-k$ ). With  $k=5$  and  $n=10$ , what values of the  $F$ -statistic would lead you to reject  $H_0$ ? What values of  $R^2$  would lead you to reject? What conclusion could you make?

$$R^2 = \frac{SSM}{SST} = 1 - \frac{SSE}{SST} \quad \text{proportion of variability explained by the groups.}$$

Under  $H_0$ :

All group means are equal, so the model does not explain any variability beyond what is expected by random error.  
 $\sigma^2 \sim X^2_{k-1}$  and  $\sigma^2 \sim X^2_{N-k}$  also  $\frac{SSM}{SST} \sim X^2_{k-1}$  (part b)

$$R^2 = \frac{SSM}{SST} = 1 - \frac{SSE}{SST}$$

$$SSM = \sigma^2 X \quad \text{where } X \sim X^2_{k-1} \sim \text{Gamma}\left(\frac{k-1}{2}, \frac{1}{2\sigma^2}\right) \quad \text{under } H_0$$

$$\text{so } SST \sim \text{Gamma}\left(\frac{N-1}{2}, \frac{1}{2\sigma^2}\right)$$

$$SSE = \sigma^2 Y \quad \text{where } Y \sim X^2_{N-k} \sim \text{Gamma}\left(\frac{N-k}{2}, \frac{1}{2\sigma^2}\right) \quad \text{always}$$

$$SST = \sigma^2 W \quad \text{where } W \sim X^2_{N-1}$$

$$\text{Thus } R^2 = \frac{SSM}{SST} = \frac{X}{W} = \frac{X}{X+Y} \quad \text{Let } \frac{X}{X+Y} = Z$$

$$\frac{X}{Z} = X+Y \quad Y=M$$

$$\text{Inverses: } X = \frac{ZW}{1-Z}, \quad Y = M$$

$$J = \begin{vmatrix} \frac{\partial X}{\partial Z} & \frac{\partial X}{\partial M} \\ \frac{\partial Y}{\partial Z} & \frac{\partial Y}{\partial M} \end{vmatrix} = \begin{bmatrix} \frac{M}{(1-Z)^2} & \frac{Z}{(1-Z)} \\ 0 & 1 \end{bmatrix} \quad |J| = \frac{M}{(1-Z)^2}$$

$X$  and  $Y$  are independent (in part b  $SSM$  and  $SSE$  indep.)

$$\text{So, } f_{X,Y}(X,Y) = f_X(X) f_Y(Y)$$

$$\text{Then, } f_{Z,M}(Z,M) = f_{X,Y}(X,Y) \cdot |J|$$

$$f_{Z,M}(Z,M) = f_X\left(\frac{ZW}{1-Z}\right) f_Y(M) \cdot |J|$$

$$= \frac{1}{\Gamma\left(\frac{k-1}{2}\right)} Z^{k/2-1} \left(\frac{ZW}{1-Z}\right)^{\frac{k-1}{2}-1} e^{-\frac{ZW}{2(1-Z)}} \cdot \frac{1}{\Gamma\left(\frac{N-k}{2}\right)} Z^{\frac{N-k}{2}} e^{-\frac{M}{2}} \cdot \frac{M}{(1-Z)^2}$$

$$= \frac{1}{\Gamma\left(\frac{k-1}{2}\right) \Gamma\left(\frac{N-k}{2}\right) Z^{\frac{(k-1)(N-k)}{2}}} Z^{\frac{k-1}{2}-1} (1-Z)^{-\frac{k-1}{2}-1} M^{\frac{N-k}{2}-1} e^{-\frac{1}{2}\left(\frac{M}{1-Z}\right)}$$

$$f_Z(Z) = \int_0^\infty f_{Z,M}(Z,M) dM$$

$$= \frac{1}{\Gamma\left(\frac{k-1}{2}\right) \Gamma\left(\frac{N-k}{2}\right) Z^{\frac{(k-1)(N-k)}{2}}} Z^{\frac{k-1}{2}-1} (1-Z)^{-\frac{k-1}{2}-1} \int_0^\infty M^{\frac{N-k}{2}} e^{-\frac{1}{2}\left(\frac{M}{1-Z}\right)} dM$$

$$= \frac{1}{\Gamma\left(\frac{k-1}{2}\right) \Gamma\left(\frac{N-k}{2}\right) Z^{\frac{(k-1)(N-k)}{2}}} Z^{\frac{k-1}{2}-1} (1-Z)^{-\frac{k-1}{2}-1} \left(\frac{\Gamma\left(\frac{N-k}{2}\right)}{\frac{1}{2(1-Z)}^{\frac{N-k}{2}}}\right) \int_0^\infty \frac{\left(\frac{1}{2(1-Z)}\right)^{\frac{N-k}{2}}}{\Gamma\left(\frac{N-k}{2}\right)} \cdot M^{\frac{N-k}{2}-1} e^{-\frac{1}{2(1-Z)}M} dM$$

$$\downarrow \text{Gamma } a = \frac{N-1}{2} \quad b = \frac{1}{2(1-Z)}$$

know pdf integrates to 1 so

$$f_Z(Z) = \frac{1}{\Gamma\left(\frac{k-1}{2}\right) \Gamma\left(\frac{N-k}{2}\right) Z^{\frac{(k-1)(N-k)}{2}}} Z^{\frac{k-1}{2}-1} (1-Z)^{-\frac{k-1}{2}-1} \cdot \frac{\Gamma\left(\frac{N-1}{2}\right)}{\left(\frac{1}{2(1-Z)}\right)^{\frac{N-1}{2}}}$$

$$= \frac{\Gamma\left(\frac{N-1}{2}\right)}{\Gamma\left(\frac{k-1}{2}\right) \Gamma\left(\frac{N-k}{2}\right)} Z^{\frac{k-1}{2}-1} (1-Z)^{\frac{N-k}{2}-1}$$

$$\text{Beta}\left(k-1/2, N-k/2\right)$$

$$k=5 \text{ and } n=10 \quad N=30$$

$$k-1=4 \text{ and } N-k=45 \quad R^2 = \text{Beta}\left(2, \frac{50-3}{2} = 22.5\right)$$

$$F = \frac{MSM}{MSE} = \frac{SSM/(k-1)}{SSE/(N-k)} \quad \alpha = .05$$

$$F(k-1, N-k) \sim F(4, 43)$$

$$\text{at } \alpha = .05 \quad F \text{ critical value} = 2.578 \quad \text{qf(.95, 4, 43, lower.tail=TRUE)}$$

Reject for  $F_{\text{stat}} > 2.58$

$$R^2 \sim \text{Beta}(2, 22.5)$$

$$R^2_{.05} = 0.1752 \quad \text{qbeta(.95, 2, 22.5, lower.tail=TRUE)}$$

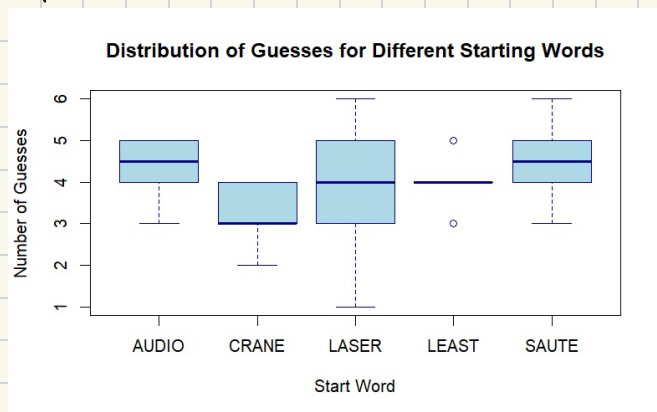
Reject for  $R^2 > 0.18647$

A higher  $R^2$  is equivalent to rejecting for a large  $F$ -stat. However  $R^2$  is more intuitive because it explains that a large  $R^2$  means a large variability btwn groups. meaning  $H_0$  seems likely untrue.

Link to  $R^2 \sim \text{Beta}(\frac{m}{2}, \frac{n}{2})$  proof: <https://statproofbook.github.io/R/beta-chi2.html>

- d) Use my Wordle data to demonstrate, with  $k = 5$  different start words and  $n = 10$  games played with each start word. The responses are the numbers of attempts to guess the word (assume these are approximately Normal with a constant variance but possibly different means). Show a graph of the data, fill in the ANOVA table and compute the  $F$ -statistic and  $R^2$ .

Graph



ANOVA Table

Source of variation	Df	Sum of Squares	Mean Square	F-Stat	Pr(>F)
Between Groups	4	9.92	2.480	2.79	.0374
Within Groups	45	40.0	0.8889		
Total	49	49.92			

$$R^2 = \frac{SSM}{SST} = 1 - \frac{SSE}{SST} = \frac{9.92}{49.92} = 0.1987$$

Both  $F$ -stat and  $R^2$  allow us to reject  $H_0$ . Meaning Each individual word does not yield the same average guesses.