

Analysis of Variance (ANOVA) and Multi-way ANOVA

MSDS 660: Statistical Methods and Experimental Design

Week 5

Dr. Ksenia Polson

Recap: Constructing a Hypothesis Test

Step 1. State the null hypothesis: A statement about the true parameter μ

- H_0 : The widespread use of driverless passenger vehicles is a good idea for society (Even better if there is a quantity associated with it).
- H_0 : It's a good idea if 85% of passenger vehicles become driverless.
- $H_0: \mu = \mu_0$, where $\mu_0 = 1.96$.

Step 2. Estimate the true parameter

- Take the sample. 26% of respondents said driverless passenger vehicles is a good idea for society

Step 3. Compute the standard error of the estimate

Step 4. Compute the test statistic (Z-score). This refers to the z-score with associated with the sampling distribution of \bar{x}

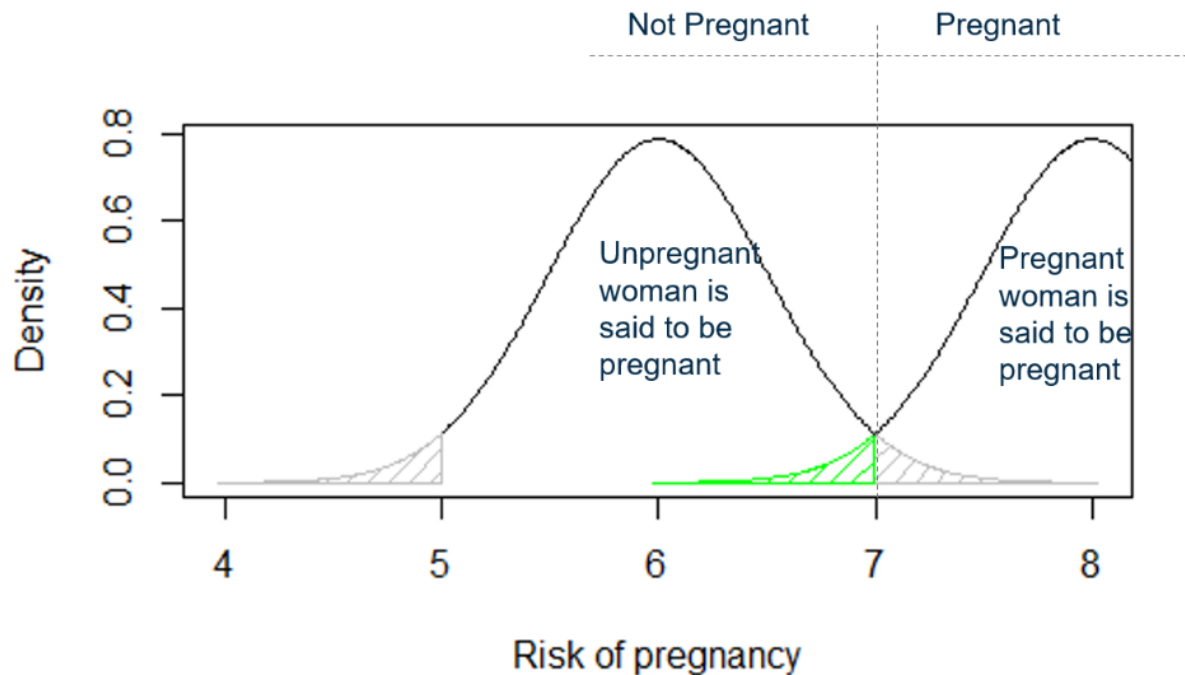
$$Z = \frac{\bar{x} - \mu}{\sigma}$$

Step 5. Make a decision. If $|Z| > 1.96$, reject the null hypothesis. Otherwise, fail to reject.

- Alternatively: Use the value of the test statistic to compute the p-value.
- If p-value is below $\alpha = 0.05$, reject. Otherwise, fail to reject.

Type I and Type II Errors

- If we apply the $\alpha = 0.05$ rejection rule, we control for the false positive rate.
 - When we reject the null hypothesis and it is true, we commit a Type I Error.
 - Say that an unpregnant woman is pregnant...oh oh!
 - At $\alpha = 0.05$, we make Type I Errors only 5% of the time (gray region in the plot)



Agenda

ANOVA (Analysis of Variance)

- Group means comparison simultaneously
- Explain one-way and multi-way ANOVA
- Explain main effects and interaction effects
- Use R to execute fits on data, and interpret the results

Comparing means...there is a problem

If I give you 2 samples, can you test to see if they lead you to believe that the populations they come from have the same mean or not?

What is the meaning of a type I error in this problem?

How about I give you 4 samples. Can you tell if they all have the same mean?

How many tests do I have to do to see if they all have the same mean?

For example, if I have 4 sections of a course and I administer a test to all 4 sections. I am looking to compare if the mean exam score is the same across the 4 sections.

$0.95^6 = 73\%$ is the probability that I got all 6 of them right. The probability, then, that I rejected at least one of the null hypotheses is 26% when I shouldn't have. I have a 26% chance that I reject the null hypothesis in error.

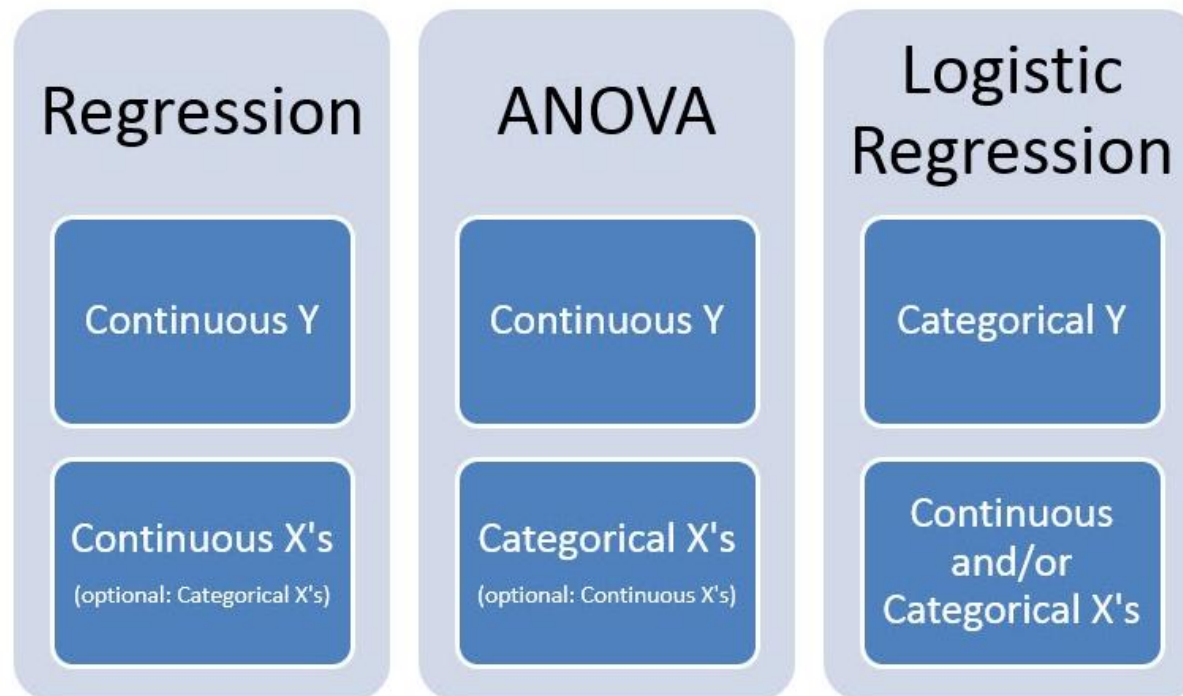
Section	Section	Alpha .05	
1	2	0.95	0.73509
1	3	0.95	0.26491
1	4	0.95	
2	3	0.95	
2	4	0.95	
3	4	0.95	

ANOVA

- ANOVA stands for Analysis of variance
- ANOVA is very similar to linear regression.

In fact, they are the same!

- Typically, we use regression for continuous predictors and ANOVA for categorical predictors.



ANOVA

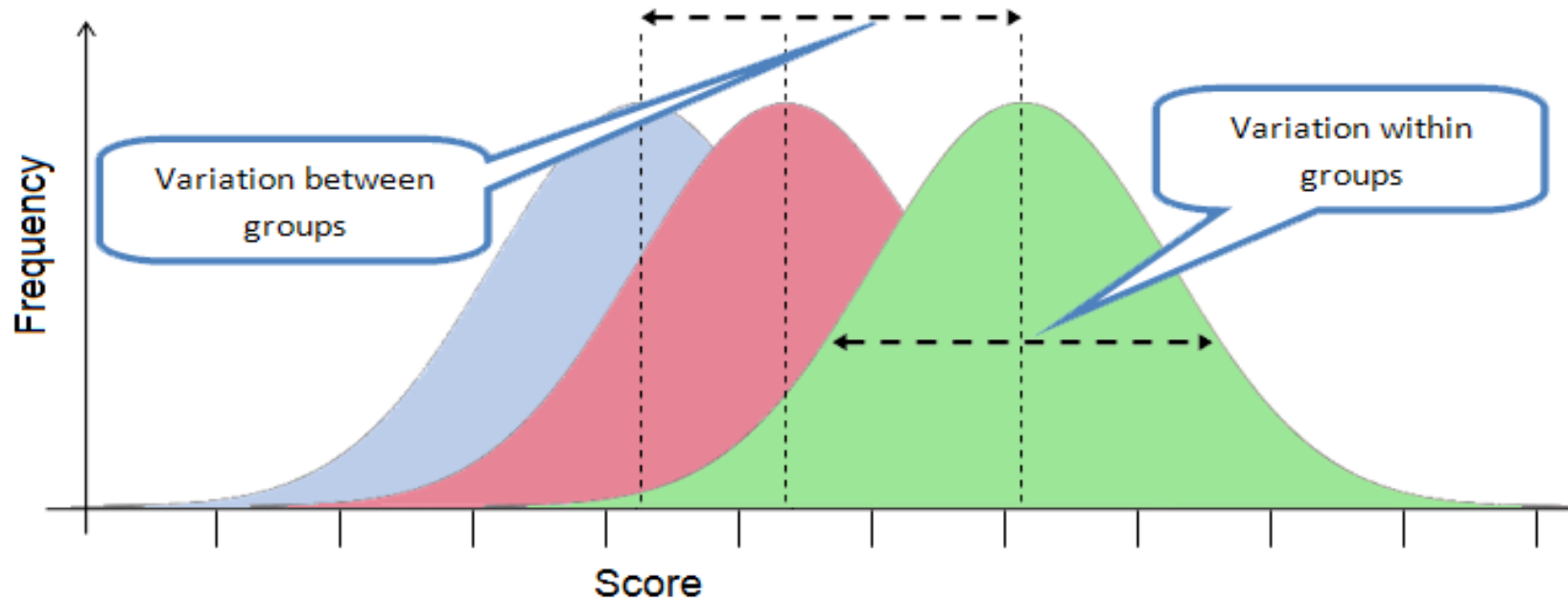
- Used to analyze the differences among group means in a sample.
- ANOVA is a statistical test of whether the population means of 3 or more groups are equal.
- Key Question: Is the variation between the groups bigger than the variation within the groups?
- Hypothesis:
$$H_0: \mu_1 = \mu_2 = \mu_3 = \cdots = \mu_k$$

$$H_a: \text{At least one group mean is different.}$$

Where $k = \#$ of groups
- Be careful, this does not mean that all the means are unequal!

ANOVA

- The p value from ANOVA tells if at least one predictor is correlated to the response – same as regression
- Post hoc tests evaluate which categorical value is significantly different from the others.



If there is more variation *between* the groups than *within*, we say the groups have different population means.

ANOVA Assumptions

1. Observations must be independent
 - There is no relationship between the observations in the different groups and between observations in the same group
2. Homogeneity of variance (aka constant variance)
 - The spread of values in each group of the groups is roughly similar around the mean -> this is called equality of variance
 - This becomes more important if the sample sizes differ between groups
 - Use Levene's Test
3. Normality of data (especially with small sample sizes)
 - Check for skewness and kurtosis values on numerical data
 - Use Q-Q plot: points need to be close to the diagonal line
4. Normally distributed residuals
 - Use Shapiro-Wilk normality test (helpful in linear fits too)
 - H_0 : Distribution of residuals is normal.
 - H_a : Distribution is not normal.

ANOVA

We can run a one-way ANOVA in R:

```
fit <- aov(y ~ factor, data = dt)
```

With a multivariate ANOVA, we just add more factors to the model.

What if we reject the null hypothesis?

The follow-up question is what pairs of samples are unequal?

From here we return to pairwise testing (usually t sample tests)

Remember the one-sample t -test...we now run multiple tests to see which pairs are different.

Post Hoc Tests

- Fisher LSD and Tukey's HSD avoid the problem of inflation of Type I error by increasing the critical value needed to reject the null.
- Post hoc tests make comparisons between the means of groups, but they make it harder to reject the null.

Do we need to do a post hoc if we don't reject the null hypothesis?

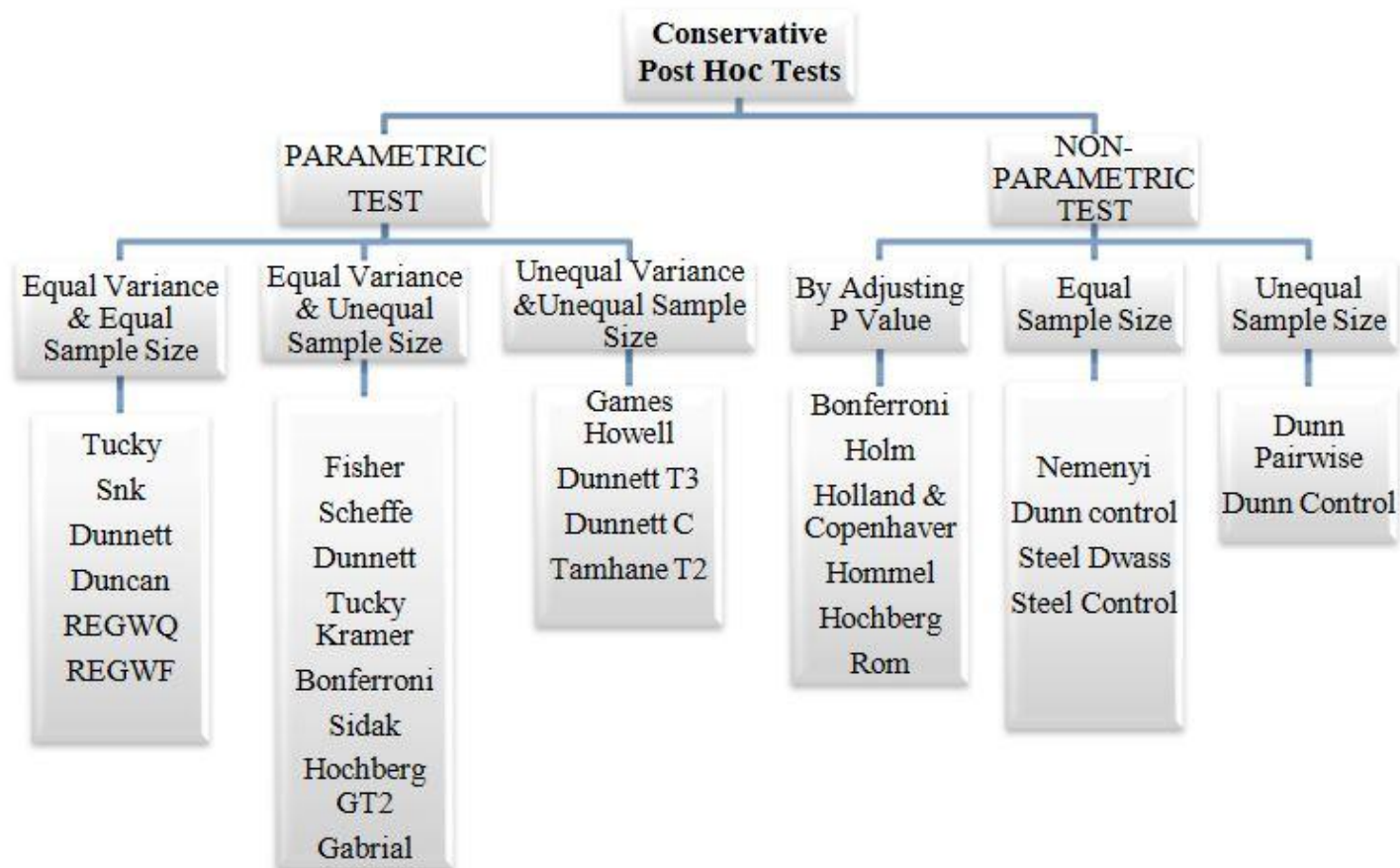
When we test many pairs of groups (they are called multiple comparisons). The Bonferroni correction applies a more stringent significance level for these tests:

$$\alpha^* = \alpha / K$$

where K is the number of comparisons being considered. If there are k groups, then usually all possible pairs are compared and $K = \frac{k(k-1)}{2}$

Post Hoc Tests “Rough” Guidelines

- What were some of the major differences between the Post Hoc tests?



Two-way ANOVA model

- What's the alpha*beta term?

$$y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk}$$

where μ is the grand mean

α_i is the i th level of first factor

β_j is the j th level of the second factor

k refers to the k th observation within the (ij) cell

ε_{ijk} is the error term

Y_{123} = Refers to the 3rd observation in the first level of factor 1 and the second level of factor 2.

Two-way ANOVA model

- What's the alpha*beta term?
 - interaction (synergy)

$$y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk}$$

- Their combined interaction have an effect on the outcome of y.

Two-way ANOVA model

```
fit <- aov (y ~ A * B * C, data = dt)
```

This is equivalent to:

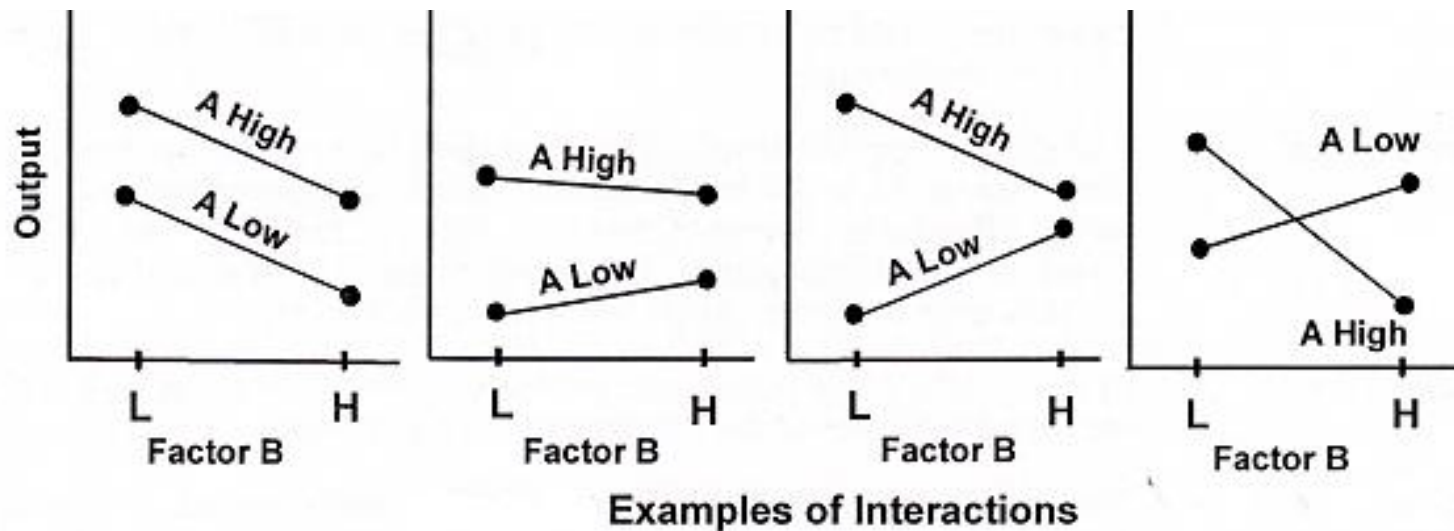
```
fit <- aov (y ~ A + B + C + A:B + A:C + B:C + A:B:C,  
           data = dt)
```

Interaction effects

- When an interaction is present, the impact of one factor depends on the level of the other factor.
- Interaction effects are evident by a low p-value (less than 0.05)
- They will show that at least one factor level is having an interaction on another factor level and giving an interaction.

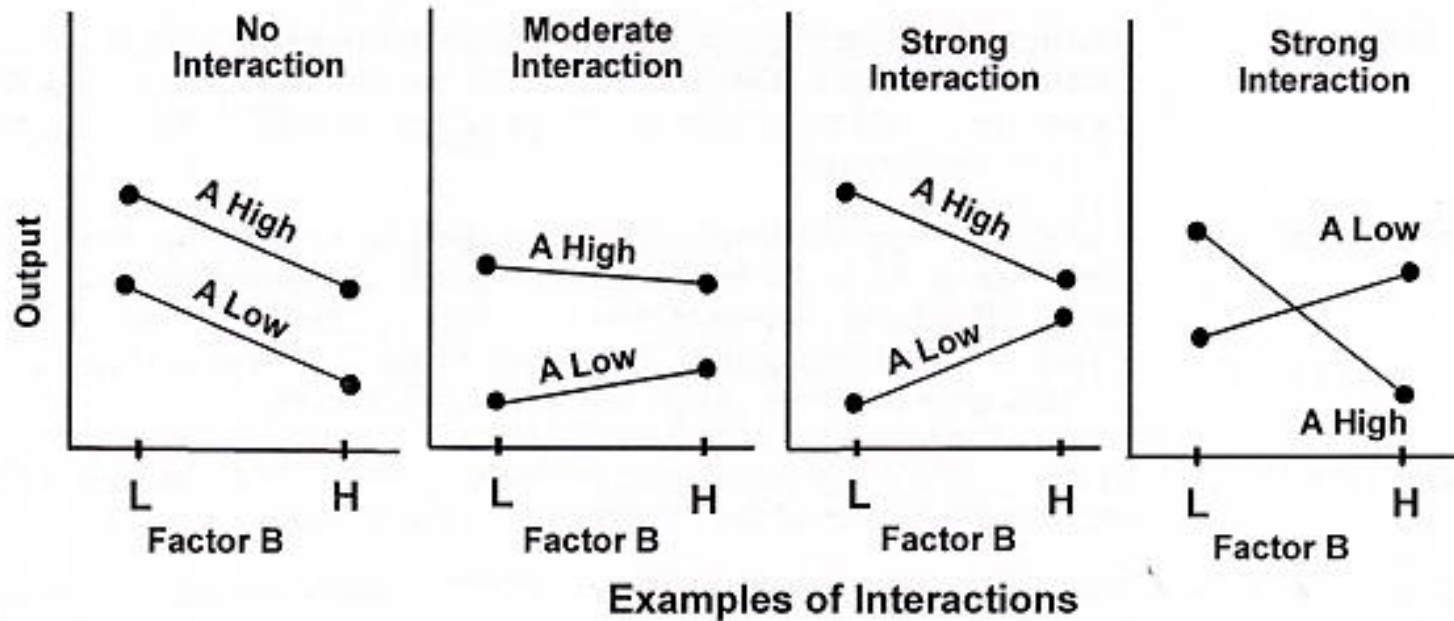
Interaction plots

- Which plots show strong interaction effects?



Interaction plots

- Which plots show strong interaction effects?



Discussion Activity

Continue working with the loan data set.

1. Form a hypothesis for the variables that maybe related. You may have both factors and numerical values in your analysis. You would need factors to create an interaction plot.
2. Run a multi-way ANOVA on loan amount received with at least 2 other variables.
3. Is there a significant interaction effect between the levels of each variable? Please plot at least one interaction plot.
4. Test for ANOVA assumptions. (At least the Levene's test for HOV)
5. Does the analysis support the hypothesis you formed initially?
6. Post your Rmd and responses to the questions to the Week 5 discussion.

Project

Continue working with the marketing data set or a data set of your choice.

1. Form a hypothesis for the variables that maybe related. You may have both factors and numerical values in your analysis. You would need factors to create an interaction plot.
2. Run a multi-way ANOVA on a continuous outcome of your choice with at least 2 other variables.
3. Is there a significant interaction effect between the levels of each variable? Please plot at least one interaction plot.
4. Test for ANOVA assumptions. (At least the Levene's test for HOV)
5. Does the analysis support the hypothesis you formed initially?
6. Post your Rmd file and knitted html file to the assignment dropbox.