

Group Project R

Matthew Farr

November 26, 2018

Parametric and Non-parametric ANOVA analysis in R

First we want to load the dataset into R and omit any missing values using `read.csv()` and `na.omit()`.

```
diet.R = read.csv("./diet.csv", header = TRUE)
diet.R = na.omit(diet.R)
```

The `na.omit()` omitted two data points, Person 25 and 26. These points were missing the value for their gender.

Then we need to tell R that Diet and gender are factors by using `as.factor()` and `factor()`.

```
diet.R$Diet = as.factor(diet.R$Diet)
diet.R$gender = factor(diet.R$gender, c(0,1))
```

Then we need to create the variable of interest ‘weightlost’.

```
diet.R$weightlost = diet.R$pre.weight - diet.R$weight6weeks
```

```
summary(diet.R)
```

```
##      Person      gender      Age      Height      pre.weight
## Min.   : 1.00    0:43   Min.   :16.00   Min.   :141.0   Min.   :58.00
## 1st Qu.:19.75    1:33   1st Qu.:32.50   1st Qu.:163.8   1st Qu.:66.00
## Median :40.50                Median :39.00   Median :169.0   Median :72.00
## Mean   :39.87                Mean   :39.22   Mean   :170.8   Mean   :72.29
## 3rd Qu.:59.25                3rd Qu.:47.25   3rd Qu.:175.2   3rd Qu.:78.00
## Max.   :78.00                Max.   :60.00   Max.   :201.0   Max.   :88.00
## Diet    weight6weeks    weightlost
## 1:24   Min.   :53.00   Min.   : -2.100
## 2:25   1st Qu.:61.95   1st Qu.: 2.300
## 3:27   Median :68.95   Median : 3.700
##        Mean   :68.34   Mean   : 3.946
##        3rd Qu.:73.67   3rd Qu.: 5.650
##        Max.   :84.50   Max.   : 9.200
```

Here is a quick summary of our dataset so you can see what variables are in the dataset and have a brief summary of those variables.

We then split the dataset into subsets by ‘gender’ using `subset()` so we can have separate datasets for the female and male data.

```
Diet.female = subset(diet.R, gender==0)
Diet.male = subset(diet.R, gender==1)
```

Now that we have all the setup work done we can move on to our ANOVA analysis.

Parametric ANOVA

To create our ANOVA model we use 'Diet' as the grouping variable and use `aov()` to create our parametric one-way ANOVA model for the female dataset. We can then use `summary()` to analyze this model.

```
anovaFemale = aov(weightlost~Diet, data = Diet.female)
summary(anovaFemale)
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## Diet           2   92.32    46.16    10.64 0.000197 ***
## Residuals     40  173.53     4.34
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

From the above output of the parametric one-way ANOVA for the female dataset, we can see that the p-value is much smaller than 0.05 so we conclude that there is at least one group that is statistically different from the other groups in the female dataset.

In order to check which group is statistically different, we run the Tukey post hoc test for pairwise comparison following a one-way ANOVA using `TukeyHSD()`.

```
TukeyHSD(anovaFemale)
```

```
##    Tukey multiple comparisons of means
##      95% family-wise confidence level
##
## Fit: aov(formula = weightlost ~ Diet, data = Diet.female)
##
## $Diet
##           diff           lwr           upr           p adj
## 2-1 -0.4428571 -2.3589312  1.473217  0.8406368
## 3-1  2.8300000  0.9461311  4.713869  0.0020846
## 3-2  3.2728571  1.3889883  5.156726  0.0003833
```

From the above result we can see that there is a statistically significant difference in weight loss between the 'Diet 3' group and the 'Diet 1' group and between the 'Diet 3' group and the 'Diet 2' group. Because of this we can determine that the 'Diet 3' group is statistically different from the other two groups.

Now let's do the same thing for the male dataset.

```
anovaMale = aov(weightlost~Diet, data = Diet.male)
summary(anovaMale)
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## Diet           2     2.0    1.001    0.148  0.863
## Residuals     30   202.8    6.760
```

From the above output of the parametric one-way ANOVA for the male dataset, we can see that the p-value is larger than 0.05 so we fail to reject the null hypothesis.

Let's confirm this by looking at the Tukey post hoc test.

```
TukeyHSD(anovaMale)
```

```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = weightlost ~ Diet, data = Diet.male)
##
## $Diet
##      diff      lwr      upr      p adj
## 2-1 0.4590909 -2.341516 3.259698 0.9141710
## 3-1 0.5833333 -2.161144 3.327810 0.8602563
## 3-2 0.1242424 -2.551325 2.799809 0.9928028
```

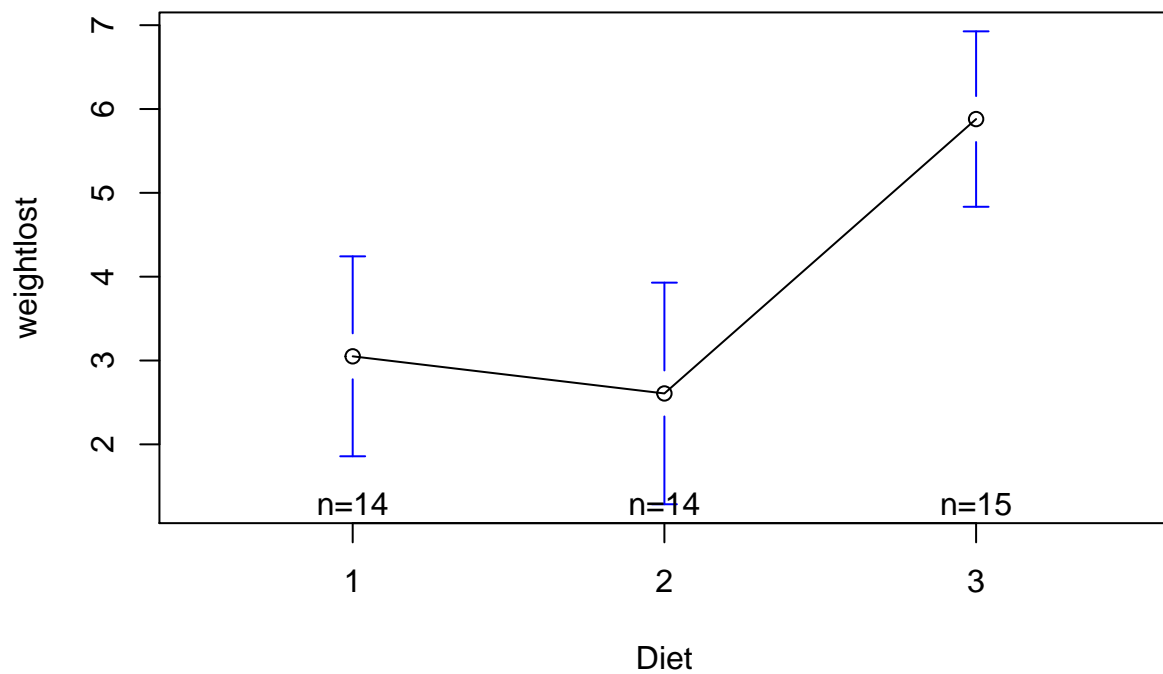
As you can see there is no statistically significant difference between the three groups in the male dataset. The following plots help visualize these two Tukey post hoc tests. Here is the plot for the female data.

```
library(gplots)
```

```
##
## Attaching package: 'gplots'

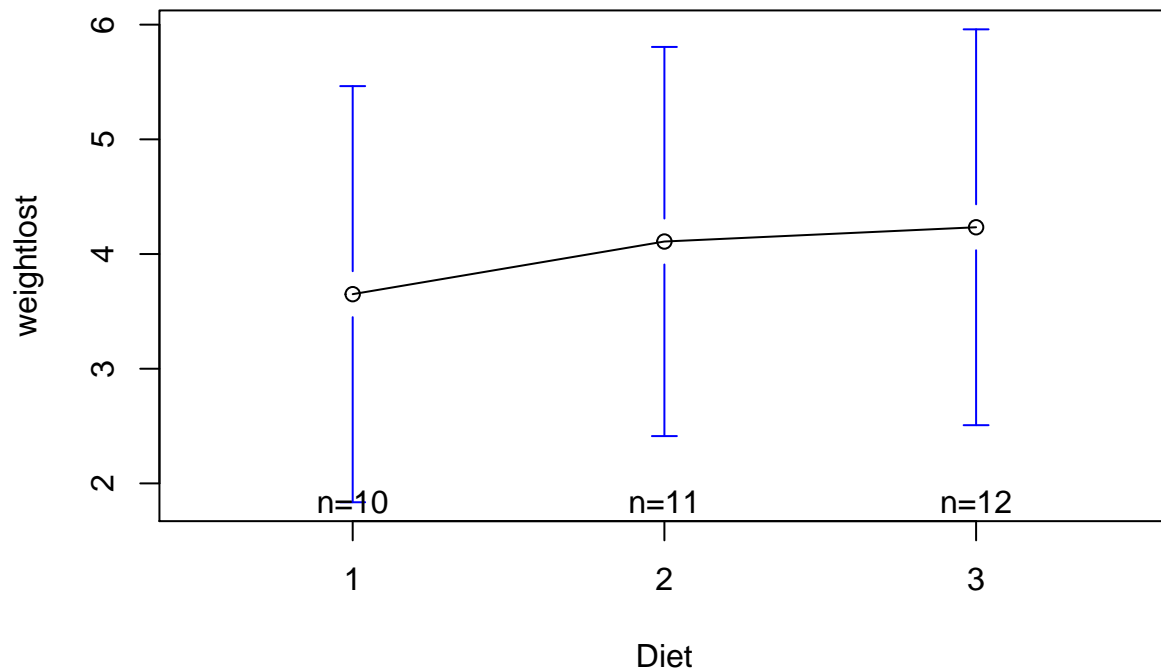
## The following object is masked from 'package:stats':
##
##      lowess
```

```
plotmeans(weightlost~Diet, data = Diet.female)
```



And here is the plot for the male data.

```
plotmeans(weightlost~Diet, data = Diet.male)
```



For both of these plots `n` is the number of people that participated in each Diet.

A similar way we could analysis this data is to do parametric two-way ANOVA for the whole dataset. For this we still use `aov()` but now we regard 'gender' and 'Diet' as two grouping variables.

```
anova2 = aov(weightlost~gender*Diet,data=diet.R)
summary(anova2)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## gender      1     0.3   0.278    0.052 0.82062
## Diet        2    60.4  30.209    5.619 0.00546 **
## gender:Diet  2    33.9  16.952    3.153 0.04884 *
## Residuals   70   376.3    5.376
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The gender:Diet interaction is statistically significant at the $p = 0.04884$ level. There was no statistically significant difference in weight loss between gender ($p = 0.82062$), but there were statistically significant differences between Diet groups ($p = 0.00546$).

Nonparametric ANOVA: Kruskal-Wallis Test

Now let's look at Nonparametric ANOVA using the Kruskal-Wallis Test. We regard 'Diet' as the grouping variable and use `kruskal.test()` to do nonparametric one-way ANOVA, i.e. Kruskal-Wallis test for the female data.

```
kruskal.test(weightlost~Diet, data = Diet.female)
```

```
##  
## Kruskal-Wallis rank sum test  
##  
## data: weightlost by Diet  
## Kruskal-Wallis chi-squared = 14.545, df = 2, p-value = 0.0006945
```

We get a p-value that is much smaller than 0.05 so we can reject the null hypothesis and conclude that there is at least one group statistically different from the other groups in the female dataset. This is the same conclusion we got in the parametric one-way ANOVA for the female data.

Now we do the same thing for the male dataset.

```
kruskal.test(weightlost~Diet, data = Diet.male)
```

```
##  
## Kruskal-Wallis rank sum test  
##  
## data: weightlost by Diet  
## Kruskal-Wallis chi-squared = 0.62218, df = 2, p-value = 0.7326
```

We get a p-value that is larger than 0.05 so there is no statistically significant difference in weight loss between the three groups in the male dataset. This is the same result we got from the parametric one-way ANOVA for the male data.

Finally, after all of this analysis we can conclude that for this dataset that you can either analyze the data using either parametric or nonparametric ANOVA and receive the same results. Our conclusion is that one group, the Diet 3 group, had a significant difference in weight loss compared to the other groups in the female dataset. However, for the male dataset there was no significant difference in weight loss between the different Diet groups.