

Homework Assignment 4

Brendan Smith

March 3, 2016

Objective Statement: The purpose of this lab is to determine if the linear models we are developing can adequately characterize the biomass found in the studied habitats. Further, we must compare sites and genera, determining which vary more than others. We are building upon previous labs, and developing a full linear model that relates tree height as a function of DBH and site or genus. A new report will be created that shows the model results for the height as a function of DBH and a factor, being site or genus.

Methods: We begin by importing the “cleaned” riparian data frame from the previous homework. We are to first determine if there is a systematic bias in the height variation by site, then to use genus as a desired factor. The systematic bias is studied by creating a summary subset and analyzing.

Data: The data utilized in this lab is taken from the previous two labs. These data are of several genera of trees’ diameter at breast height (DBH) and height, which has been reduced to be of only five most common genera found in the original dataset. Additionally, the height was measured in meters and has been scaled to centimeters for data analysis purposes. DBH was measured in centimeters.

Code: The new code introduced and used in this assignment are `lapply`, `for` loop, `do.call()`, `with()`, `aov()` and `TukeyHSD()`.

Results: We begin our analysis as usual by importing the dataset from the comma separated value file:

```
# Load libraries
#library(stats)
#library(HH)

# Load Ripdata and place into a dataframe
rip <- read.csv("newripdata_survey.csv",sep = ",",header = TRUE)
# Add an object that scales the value of height from meters to centimeters
rip$htcm <- rip$Woody_Height_m*100
```

We proceed to make necessary adjustments to the dataframe for data manipulation purposes:

```
#Concatenate ProjCode and Plot.Name using the paste() function, typecast as a factor, then place these
rip$projplot <- as.factor(paste(rip$ProjCode,rip$Plot.Name))

#use tapply() to cycle through each project plot and generate stats
#where 'htcm' is height in cm
ripsum <- data.frame(cbind(tapply(rip$htcm,rip$projplot,mean),tapply(rip$htcm,rip$projplot,sd),tapply(r
#add column names
#(height mean, height standard deviation, number of plots)
colnames(ripsum) <- c("htcmmn","htcmsd","plot.n")
#add a projplot column (from row names) to ripsum
ripsum$projplot <- as.factor(rownames(ripsum))

#subset for plots with more than one measurement
ripsum <- ripsum[ripsum$plot.n > 1,]
#Add a proj column and populate with the first five letters of projplot (the ProjCode) via the substr()
ripsum$proj <- as.factor(substr(ripsum$projplot,1,5))
```

We compare the for loop operation to the list apply (lapply) in order to demonstrate that although the functions have similar outcomes, the list apply is more efficient for evaluating an array of values simultaneously.

```
#create list of project sites
projlevels <- levels(ripsum$proj) #compare a 'for' loop of summary
for (p in 1:length(projlevels)) print(summary(ripsum[ripsum$proj == projlevels[p],]))
```

```
##      htcmmn      htcmsd      plot.n      projplot
## Min.   : 204.1   Min.    : 42.43   Min.    : 2.00   COSRP 1: 1
## 1st Qu.: 769.0   1st Qu.: 277.50   1st Qu.: 9.50   COSRP 2: 1
## Median :1000.0   Median : 436.05   Median : 18.00   COSRP 3: 1
## Mean   :1042.3   Mean    : 459.94   Mean    : 25.23   COSRP 4: 1
## 3rd Qu.:1269.2   3rd Qu.: 624.36   3rd Qu.: 32.50   COSRP 5: 1
## Max.   :3330.0   Max.    :1031.60   Max.    :111.00   COSRP 6: 1
##                                     (Other):81
##
##      proj
## COSRP:87
## HEROW: 0
## NAPSO: 0
## SACTO: 0
##
##
##      htcmmn      htcmsd      plot.n      projplot
## Min.   : 421.3   Min.    : 68.18   Min.    : 8.00   HEROW RIP01:1
## 1st Qu.: 505.6   1st Qu.: 141.83   1st Qu.:11.50   HEROW RIP02:1
## Median : 883.0   Median : 218.21   Median :27.00   HEROW RIP03:1
## Mean   : 978.1   Mean    : 401.83   Mean    :39.14   HEROW RIP04:1
## 3rd Qu.:1332.6   3rd Qu.: 546.12   3rd Qu.:62.00   HEROW RIP05:1
## Max.   :1866.0   Max.    :1150.50   Max.    :92.00   HEROW RIP06:1
##                                     (Other)   :1
##
##      proj
## COSRP:0
## HEROW:7
## NAPSO:0
## SACTO:0
##
##
##      htcmmn      htcmsd      plot.n      projplot
## Min.   : 671.0   Min.    :171.6   Min.    : 4.00   NAPSO Crp2013_509:1
## 1st Qu.: 745.4   1st Qu.:472.4   1st Qu.: 6.75   NAPSO S02013_200 :1
## Median : 814.6   Median :556.3   Median :17.50   NAPSO S02013_202 :1
## Mean   : 977.8   Mean    :538.2   Mean    :24.00   NAPSO S02013_203 :1
## 3rd Qu.:1016.4   3rd Qu.:649.9   3rd Qu.:29.25   NAPSO S02013_207 :1
## Max.   :1888.3   Max.    :819.5   Max.    :81.00   NAPSO S02013_211 :1
##                                     (Other)      :2
##
##      proj
## COSRP:0
## HEROW:0
## NAPSO:8
## SACTO:0
##
##
```

```
##
##      htcmmn      htcmsd      plot.n      projplot
## Min.   : 443.2   Min.    : 93.04   Min.    : 2.00   SACTO EW2013_100: 1
## 1st Qu.: 678.7   1st Qu.: 254.57   1st Qu.: 4.00   SACTO EW2013_101: 1
## Median : 977.0   Median : 584.23   Median : 7.00   SACTO EW2013_102: 1
## Mean   :1071.0   Mean    : 623.96   Mean    :11.82   SACTO EW2013_103: 1
## 3rd Qu.:1256.6   3rd Qu.: 882.65   3rd Qu.:17.00   SACTO EW2013_106: 1
## Max.   :2532.9   Max.    :1668.77   Max.    :59.00   SACTO EW2013_110: 1
##                                     (Other)      :38
##      proj
## COSRP: 0
## HEROW: 0
## NAPSO: 0
## SACTO:44
##
##
##
```

```
#with a summary using lapply() (known as list apply)
lapply(projlevels, function(x) summary(riplesum[riplesum$proj == x,]))
```

```
## [[1]]
##      htcmmn      htcmsd      plot.n      projplot
## Min.   : 204.1   Min.    : 42.43   Min.    : 2.00   COSRP 1: 1
## 1st Qu.: 769.0   1st Qu.: 277.50   1st Qu.: 9.50   COSRP 2: 1
## Median :1000.0   Median : 436.05   Median : 18.00   COSRP 3: 1
## Mean   :1042.3   Mean    : 459.94   Mean    : 25.23   COSRP 4: 1
## 3rd Qu.:1269.2   3rd Qu.: 624.36   3rd Qu.: 32.50   COSRP 5: 1
## Max.   :3330.0   Max.    :1031.60   Max.    :111.00   COSRP 6: 1
##                                     (Other):81
##      proj
## COSRP:87
## HEROW: 0
## NAPSO: 0
## SACTO: 0
##
##
##
## [[2]]
##      htcmmn      htcmsd      plot.n      projplot
## Min.   : 421.3   Min.    : 68.18   Min.    : 8.00   HEROW RIPO1:1
## 1st Qu.: 505.6   1st Qu.: 141.83   1st Qu.:11.50   HEROW RIPO2:1
## Median : 883.0   Median : 218.21   Median :27.00   HEROW RIPO3:1
## Mean   : 978.1   Mean    : 401.83   Mean    :39.14   HEROW RIPO4:1
## 3rd Qu.:1332.6   3rd Qu.: 546.12   3rd Qu.:62.00   HEROW RIPO5:1
## Max.   :1866.0   Max.    :1150.50   Max.    :92.00   HEROW RIPO6:1
##                                     (Other)    :1
##      proj
## COSRP:0
## HEROW:7
## NAPSO:0
## SACTO:0
##
```

```
##
##
##
## [[3]]
##      htcmmn      htcmsd      plot.n      projplot
## Min.   : 671.0   Min.   :171.6   Min.   : 4.00   NAPS0 Crp2013_509:1
## 1st Qu.: 745.4   1st Qu.:472.4   1st Qu.: 6.75   NAPS0 S02013_200 :1
## Median : 814.6   Median :556.3   Median :17.50   NAPS0 S02013_202 :1
## Mean   : 977.8   Mean   :538.2   Mean   :24.00   NAPS0 S02013_203 :1
## 3rd Qu.:1016.4   3rd Qu.:649.9   3rd Qu.:29.25   NAPS0 S02013_207 :1
## Max.   :1888.3   Max.   :819.5   Max.   :81.00   NAPS0 S02013_211 :1
##                                     (Other)      :2
##      proj
## COSRP:0
## HEROW:0
## NAPS0:8
## SACTO:0
##
##
##
##
## [[4]]
##      htcmmn      htcmsd      plot.n      projplot
## Min.   : 443.2   Min.   : 93.04   Min.   : 2.00   SACTO EW2013_100: 1
## 1st Qu.: 678.7   1st Qu.: 254.57   1st Qu.: 4.00   SACTO EW2013_101: 1
## Median : 977.0   Median : 584.23   Median : 7.00   SACTO EW2013_102: 1
## Mean   :1071.0   Mean   : 623.96   Mean   :11.82   SACTO EW2013_103: 1
## 3rd Qu.:1256.6   3rd Qu.: 882.65   3rd Qu.:17.00   SACTO EW2013_106: 1
## Max.   :2532.9   Max.   :1668.77   Max.   :59.00   SACTO EW2013_110: 1
##                                     (Other)      :38
##      proj
## COSRP: 0
## HEROW: 0
## NAPS0: 0
## SACTO:44
##
##
##
```

It can be seen that while the outcome is the same, the setup and application of the for loop is somewhat inefficient in that we must manually indicate the start and stop indices in order to iterate through the entire array/vector individually with the for loop, whereas with the list apply, all elements are evaluated by the function automatically.

We then use `lapply` to randomly select six sample plot summaries from each project site. This is done by first introducing a variable that is set to the integer 6, the number of samples desired. The function `sample` is then used to output the desired number of samples (6) randomly from each project site. To execute this for all project sites, the `lapply` function is utilized. These values are stored and combined by row using the `rbind()` function along with `do.call()`. Finally, the summary is output.

```
nsamples <- 6 #Set the number of samples
ripres <- lapply(projlevels, function(x) ripsum[which(ripsum$proj == x),][sample(nrow(ripsum[which(ripsum$proj == x),]), nsamples)])
# combine samples by row using rbind()
# and by calling ripres lapply function from do.call()
```

```
ripsample <- do.call(rbind,ripres)
summary(ripsample$proj)
```

```
## COSRP HEROW NAPSO SACTO
##      6      6      6      6
```

The coefficient of variation (CV) then added to the summary table by calculating the CV by means of the `with()` function. The `with()` function evaluates an R expression (second input term) of the input data (first term). In this case, we are evaluating the coefficient of variation, which is the standard deviation divided by the mean. This value is then stored into our dataframe.

```
#calculate CV using with(data,calc)
ripsample$cv <- with(ripsample, htcmsd / htcmmn)
```

Equipped with the coefficient of variation, we can now run a one-way analysis of variation (ANOVA) on the data frame. This is done by utilizing the `aov` function, which we input the CV as a function of the project code. The output is the sum of squares, degrees of freedom and the residual standard error, all of which are stored in a new variable followed by a summary output.

```
rip.proj.cv.aov = aov(cv~proj,data=ripsample)
summary(rip.proj.cv.aov)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## proj              3  1.306  0.4353   0.626  0.607
## Residuals        20 13.908  0.6954
```

```
#compare it against
summary.lm(rip.proj.cv.aov)
```

```
##
## Call:
## aov(formula = cv ~ proj, data = ripsample)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.9748 -0.2473 -0.0856  0.1520  3.2953
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   0.9940     0.3404   2.920  0.00847 **
## projHEROW    -0.6357     0.4815  -1.320  0.20164
## projNAPSO    -0.3074     0.4815  -0.638  0.53041
## projSACTO    -0.4583     0.4815  -0.952  0.35246
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8339 on 20 degrees of freedom
## Multiple R-squared:  0.08584,    Adjusted R-squared:  -0.05128
## F-statistic: 0.626 on 3 and 20 DF,  p-value: 0.6066
```

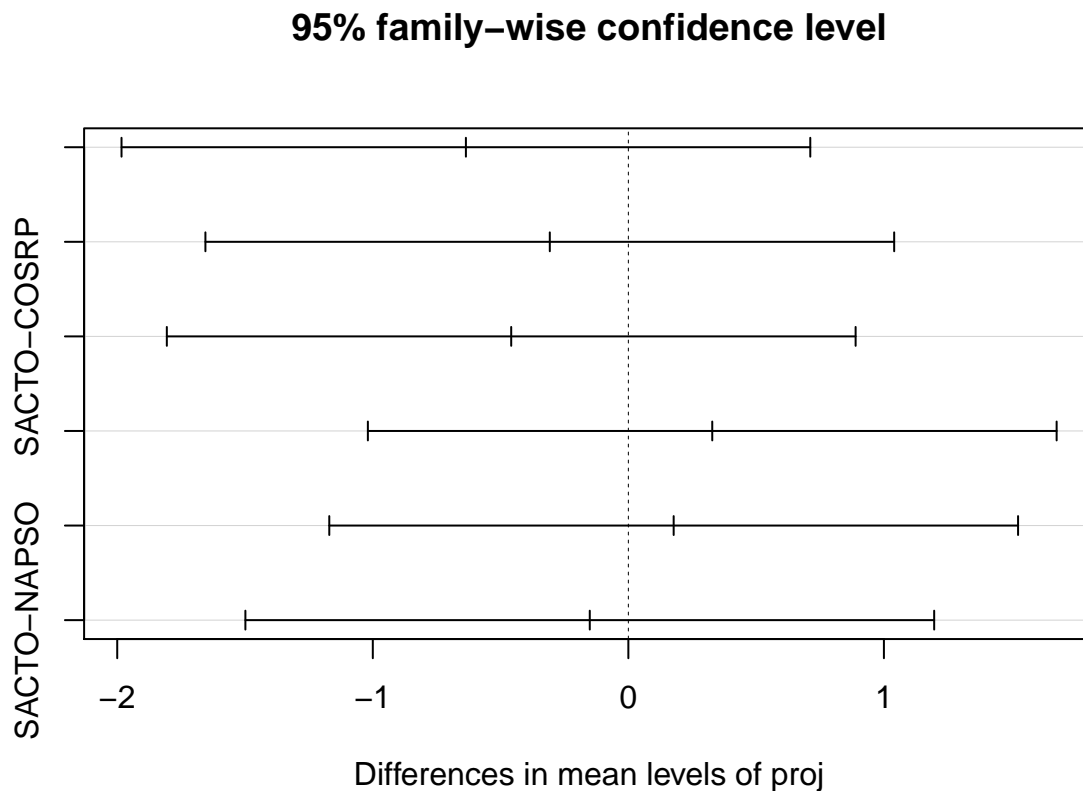
The main differences between the summary of the ANOVA results and the summary of the linear model of the ANOVA is that the lm version yields information regarding the individual project sites and the residuals, whereas the summary of the ANOVA yields information regarding all project sites and residuals, limited to DOF, sum of squares, mean squared, F value and probability.

A Tukey test is performed to check for significant differences between sites, and print out the results. We can see from the print out of the Tukey test and the plot that there is not a significant difference between sites.

```
rip.aov.hsd <- TukeyHSD(rip.proj.cv.aov)
rip.aov.hsd
```

```
##    Tukey multiple comparisons of means
##      95% family-wise confidence level
##
## Fit: aov(formula = cv ~ proj, data = ripsample)
##
## $proj
##              diff          lwr          upr      p adj
## HEROW-COSRP -0.6356749 -1.983250  0.7118999 0.5612779
## NAPSO-COSRP -0.3073971 -1.654972  1.0401777 0.9182967
## SACTO-COSRP -0.4583468 -1.805922  0.8892281 0.7775100
## NAPSO-HEROW  0.3282778 -1.019297  1.6758527 0.9027786
## SACTO-HEROW  0.1773281 -1.170247  1.5249030 0.9824389
## SACTO-NAPSO -0.1509497 -1.498525  1.1966252 0.9889948
```

```
plot(rip.aov.hsd)
```



In step two we are evaluating the analysis of covariance (ANCOVA), in which we are using the height as a

function of DBH and genus as a factor. We begin by creating a general linear model where the height of the tree is a function of the DBH and a function of genus separately.

```
# Create a linear model for height versus DBH
rip.cov.htdbh <-glm(rip$htcm~rip$Woody_DBH_cm)
# Generate summary for this model
summary(rip.cov.htdbh)
```

```
##
## Call:
## glm(formula = rip$htcm ~ rip$Woody_DBH_cm)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -3425.7   -260.8    -71.8    198.0   9815.7
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    555.214     10.546   52.64  <2e-16 ***
## rip$Woody_DBH_cm  20.786       0.387   53.71  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 197786.6)
##
##      Null deviance: 1201446793  on 3191  degrees of freedom
## Residual deviance:  630939390  on 3190  degrees of freedom
## AIC: 47989
##
## Number of Fisher Scoring iterations: 2
```

```
summary.lm(rip.cov.htdbh)
```

```
##
## Call:
## glm(formula = rip$htcm ~ rip$Woody_DBH_cm)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3425.7   -260.8    -71.8    198.0   9815.7
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    555.214     10.546   52.64  <2e-16 ***
## rip$Woody_DBH_cm  20.786       0.387   53.71  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 444.7 on 3190 degrees of freedom
## Multiple R-squared:  0.4749, Adjusted R-squared:  0.4747
## F-statistic: 2884 on 1 and 3190 DF,  p-value: < 2.2e-16
```

```
# Create a linear model for height versus Genus
rip.cov.htg <-glm(rip$hgtcm~rip$Genus)
# Generate summary for this model
summary(rip.cov.htg)
```

```
##
## Call:
## glm(formula = rip$hgtcm ~ rip$Genus)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1264.3   -323.4    -85.1    212.7   9789.5
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      707.32      25.47  27.773 < 2e-16 ***
## rip$GenusFraxinus  -96.81      35.76  -2.707  0.00683 **
## rip$GenusPopulus    676.99      31.69  21.361 < 2e-16 ***
## rip$GenusQuercus    258.19      32.85   7.860 5.22e-15 ***
## rip$GenusSalix      37.82      32.08   1.179  0.23859
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 293171.6)
##
##      Null deviance: 1201446793  on 3191  degrees of freedom
## Residual deviance:  934337819  on 3187  degrees of freedom
## AIC: 49248
##
## Number of Fisher Scoring iterations: 2
```

```
summary.lm(rip.cov.htg)
```

```
##
## Call:
## glm(formula = rip$hgtcm ~ rip$Genus)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1264.3   -323.4    -85.1    212.7   9789.5
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      707.32      25.47  27.773 < 2e-16 ***
## rip$GenusFraxinus  -96.81      35.76  -2.707  0.00683 **
## rip$GenusPopulus    676.99      31.69  21.361 < 2e-16 ***
## rip$GenusQuercus    258.19      32.85   7.860 5.22e-15 ***
## rip$GenusSalix      37.82      32.08   1.179  0.23859
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 541.5 on 3187 degrees of freedom
```



```
## Multiple R-squared:  0.2223, Adjusted R-squared:  0.2213
## F-statistic: 227.8 on 4 and 3187 DF,  p-value: < 2.2e-16
```

```
# Create a linear model for height as a function of DBH and Genus
rip.cov.htdbhg <-glm(rip$htcm~rip$Woody_DBH_cm*rip$Genus)
# Generate summary for this model
summary(rip.cov.htdbhg)
```

```
##
## Call:
## glm(formula = rip$htcm ~ rip$Woody_DBH_cm * rip$Genus)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2721.1   -231.5    -25.1    193.9   10030.7
##
## Coefficients:
##                                Estimate Std. Error t value Pr(>|t|)
## (Intercept)                   303.993     29.714   10.231 < 2e-16
## rip$Woody_DBH_cm                30.549       1.735   17.605 < 2e-16
## rip$GenusFraxinus              28.923     41.305    0.700 0.483838
## rip$GenusPopulus              661.299     35.487   18.635 < 2e-16
## rip$GenusQuercus              143.344     37.167    3.857 0.000117
## rip$GenusSalix                210.982     39.513    5.340 9.97e-08
## rip$Woody_DBH_cm:rip$GenusFraxinus  -4.592       2.677   -1.715 0.086381
## rip$Woody_DBH_cm:rip$GenusPopulus  -14.620       1.809   -8.083 8.86e-16
## rip$Woody_DBH_cm:rip$GenusQuercus  -10.035       1.849   -5.427 6.17e-08
## rip$Woody_DBH_cm:rip$GenusSalix    -8.593       2.696   -3.187 0.001452
##
## (Intercept)                  ***
## rip$Woody_DBH_cm             ***
## rip$GenusFraxinus
## rip$GenusPopulus             ***
## rip$GenusQuercus             ***
## rip$GenusSalix              ***
## rip$Woody_DBH_cm:rip$GenusFraxinus .
## rip$Woody_DBH_cm:rip$GenusPopulus ***
## rip$Woody_DBH_cm:rip$GenusQuercus ***
## rip$Woody_DBH_cm:rip$GenusSalix  **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 161859.5)
##
##      Null deviance: 1201446793  on 3191  degrees of freedom
## Residual deviance:  515037025  on 3182  degrees of freedom
## AIC: 47357
##
## Number of Fisher Scoring iterations: 2
```

```
summary.lm(rip.cov.htdbhg)
```

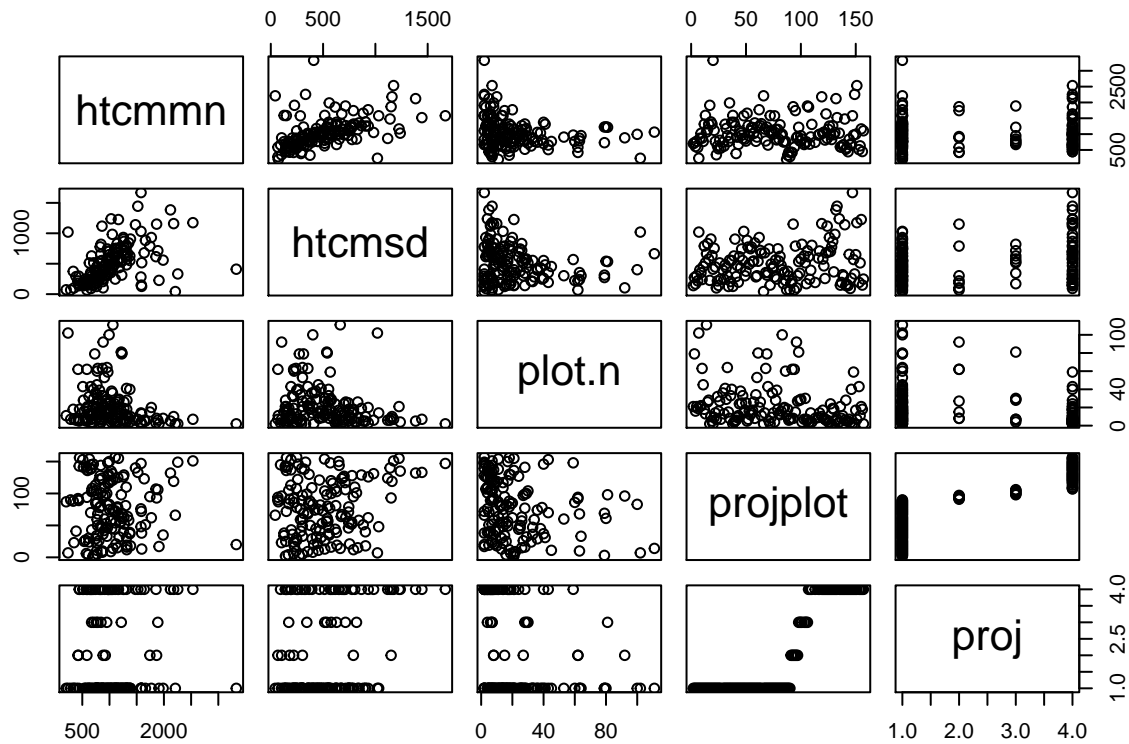
```
##
```

```
## Call:
## glm(formula = rip$h_tcm ~ rip$Woody_DBH_cm * rip$Genus)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2721.1  -231.5   -25.1   193.9 10030.7
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      303.993     29.714  10.231 < 2e-16
## rip$Woody_DBH_cm      30.549      1.735  17.605 < 2e-16
## rip$GenusFraxinus     28.923     41.305   0.700 0.483838
## rip$GenusPopulus     661.299     35.487  18.635 < 2e-16
## rip$GenusQuercus     143.344     37.167   3.857 0.000117
## rip$GenusSalix       210.982     39.513   5.340 9.97e-08
## rip$Woody_DBH_cm:rip$GenusFraxinus  -4.592      2.677  -1.715 0.086381
## rip$Woody_DBH_cm:rip$GenusPopulus  -14.620      1.809  -8.083 8.86e-16
## rip$Woody_DBH_cm:rip$GenusQuercus  -10.035      1.849  -5.427 6.17e-08
## rip$Woody_DBH_cm:rip$GenusSalix    -8.593      2.696  -3.187 0.001452
##
## (Intercept)          ***
## rip$Woody_DBH_cm      ***
## rip$GenusFraxinus
## rip$GenusPopulus      ***
## rip$GenusQuercus      ***
## rip$GenusSalix        ***
## rip$Woody_DBH_cm:rip$GenusFraxinus .
## rip$Woody_DBH_cm:rip$GenusPopulus ***
## rip$Woody_DBH_cm:rip$GenusQuercus ***
## rip$Woody_DBH_cm:rip$GenusSalix    **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 402.3 on 3182 degrees of freedom
## Multiple R-squared:  0.5713, Adjusted R-squared:  0.5701
## F-statistic: 471.2 on 9 and 3182 DF, p-value: < 2.2e-16
```

```
#ancovaplot(rip$h_tcm~rip$Woody_DBH_cm*rip$Genus,data=survey)
```

The next step would be to generate an `ancovaplot()` for the two model formulations, with an without the interaction; however, the `ancovaplot()` function will always through a “subset” error for this data set for some reason.

```
plot(riptest)
```



Discussion: Through the use of the coefficient of variation, we were able to run a one-way ANOVA, followed by checking for significant differences between sites using the Tukey test. By analyzing the summary and plot of these results, we can see that there is not significant differences in the variation of the height at these project codes. This is to say that height variability does not differ much between project codes. By utilizing the Tukey test, we are analyzing the differences between the means of the levels of the factor created by the one-way ANOVA. We can see that the difference between the means are relatively low, indicating that there is not a significant difference. Furthermore, if we look at the plot generated utilizing the Tukey data, we can see that all the project codes' differences in mean levels fall around zero, reinforcing our notion.

Limitations: A limitation that was encountered in this homework assignment was the inability to utilize the `ancovaplot()` function. Though the input parameters seemed to be accurate, an error was always thrown.