

3. Worksheet: Basic R

Erica Nadolski; Z620: Quantitative Biodiversity, Indiana University

18 January, 2023

OVERVIEW

This worksheet introduces some of the basic features of the R computing environment (<http://www.r-project.org>). It is designed to be used along side the **3. RStudio** handout in your binder. You will not be able to complete the exercises without the corresponding handout.

Directions:

1. In the Markdown version of this document in your cloned repo, change “Student Name” on line 3 (above) with your name.
2. Complete as much of the worksheet as possible during class.
3. Use the handout as a guide; it contains a more complete description of data sets along with examples of proper scripting needed to carry out the exercises.
4. Answer questions in the worksheet. Space for your answers is provided in this document and is indicated by the “>” character. If you need a second paragraph be sure to start the first line with “>”. You should notice that the answer is highlighted in green by RStudio (color may vary if you changed the editor theme).
5. Before you leave the classroom today, it is *imperative* that you **push** this file to your GitHub repo, at whatever stage you are. This will enable you to pull your work onto your own computer.
6. When you have completed the worksheet, **Knit** the text and code into a single PDF file by pressing the **Knit** button in the RStudio scripting panel. This will save the PDF output in your ‘3.RStudio’ folder.
7. After Knitting, please submit the worksheet by making a **push** to your GitHub repo and then create a **pull request** via GitHub. Your pull request should include this file (**3.RStudio_Worksheet.Rmd**) with all code blocks filled out and questions answered) and the PDF output of Knitr (**3.RStudio_Worksheet.pdf**).

The completed exercise is due on **Wednesday, January 18th, 2023 before 12:00 PM (noon)**.

1) HOW WE WILL BE USING R AND OTHER TOOLS

You are working in an RMarkdown (.Rmd) file. This allows you to integrate text and R code into a single document. There are two major features to this document: 1) Markdown formatted text and 2) “chunks” of R code. Anything in an R code chunk will be interpreted by R when you *Knit* the document.

When you are done, you will *knit* your document together. However, if there are errors in the R code contained in your Markdown document, you will not be able to knit a PDF file. If this happens, you will need to review your code, locate the source of the error(s), and make the appropriate changes. Even if you are able to knit without issue, you should review the knitted document for correctness and completeness before you submit the Worksheet. Next to the **Knit** button in the RStudio scripting panel there is a spell checker button (ABC) button.

2) SETTING YOUR WORKING DIRECTORY

In the R code chunk below, please provide the code to: 1) clear your R environment, 2) print your current working directory, and 3) set your working directory to your '3.RStudio' folder.

```
rm(list=ls())
getwd()

## [1] "/Users/ericanadolski/GitHub/QB2023_Nadolski/2.Worksheets/3.RStudio"

setwd("/Users/ericanadolski/GitHub/QB2023_Nadolski/2.Worksheets/3.RStudio")
```

3) USING R AS A CALCULATOR

To follow up on the pre-class exercises, please calculate the following in the R code chunk below. Feel free to reference the **1. Introduction to version control and computing tools** handout.

- 1) the volume of a cube with length, l , = 5 (volume = l^3)
- 2) the area of a circle with radius, r , = 2 (area = $\pi * r^2$).
- 3) the length of the opposite side of a right-triangle given that the angle, θ , = $\pi/4$. (radians, a.k.a. 45°) and with hypotenuse length $\sqrt{2}$ (remember: $\sin(\theta) = \text{opposite}/\text{hypotenuse}$).
- 4) the log (base e) of your favorite number.

```
# 1
l = 5
v = l^3
v

## [1] 125

# 2
r = 2
a = pi * r^2
a

## [1] 12.56637

# 3
theta = pi/4
hypotenuse = sqrt(2)
opposite = sin(theta) * hypotenuse
opposite

## [1] 1

# 4
fav_num = 389
log(fav_num)

## [1] 5.963579
```

4) WORKING WITH VECTORS

To follow up on the pre-class exercises, please perform the requested operations in the R-code chunks below.

Basic Features Of Vectors

In the R-code chunk below, do the following: 1) Create a vector **x** consisting of any five numbers. 2) Create a new vector **w** by multiplying **x** by 14 (i.e., “scalar”). 3) Add **x** and **w** and divide by 15.

```
x <- c(3,8,9,17,25)
w <- x * 14
(x + w) / 15
```

```
## [1] 3 8 9 17 25
```

Now, do the following: 1) Create another vector (**k**) that is the same length as **w**. 2) Multiply **k** by **x**. 3) Use the combine function to create one more vector, **d** that consists of any three elements from **w** and any four elements of **k**.

```
k <- c(4,7,1,2,2)
k * x
```

```
## [1] 12 56 9 34 50
```

```
d <- c(w[1:3], k[2:5])
```

Summary Statistics of Vectors

In the R-code chunk below, calculate the **summary statistics** (i.e., maximum, minimum, sum, mean, median, variance, standard deviation, and standard error of the mean) for the vector (**v**) provided.

```
v <- c(16.4, 16.0, 10.1, 16.8, 20.5, NA, 20.2, 13.1, 24.8, 20.2, 25.0, 20.5, 30.5, 31.4, 27.1)
max(na.omit(x))
```

```
## [1] 25
```

```
min(na.omit(x))
```

```
## [1] 3
```

```
sum(na.omit(x))
```

```
## [1] 62
```

```
mean(na.omit(x))
```

```
## [1] 12.4
```

```
median(na.omit(x))
```

```
## [1] 9
```

```
var(na.omit(x))
```

```
## [1] 74.8
```

```
sd(na.omit(x))
```

```
## [1] 8.648699
```

```
sem <- function(x){  
  sd(na.omit(x))/sqrt(length(na.omit(x)))  
}
```

```
sem(v)
```

```
## [1] 1.678435
```

5) WORKING WITH MATRICES

In the R-code chunk below, do the following: Using a mixture of Approach 1 and 2 from the **3. RStudio** handout, create a matrix with two columns and five rows. Both columns should consist of random numbers. Make the mean of the first column equal to 8 with a standard deviation of 2 and the mean of the second column equal to 25 with a standard deviation of 10.

```
c1 <- c(rnorm(5, mean=8, sd=2))  
c2 <- c(rnorm(5, mean=25, sd=10))  
m1 <- cbind(c1,c2)
```

Question 1: What does the `rnorm` function do? What do the arguments in this function specify? Remember to use `help()` or type `?rnorm`.

Answer 1: `rnorm` generates a random set of numbers that fulfill the specified total quantity of numbers, mean, and sd

In the R code chunk below, do the following: 1) Load `matrix.txt` from the **3.RStudio** data folder as matrix `m`. 2) Transpose this matrix. 3) Determine the dimensions of the transposed matrix.

```
m <- as.matrix(read.table("data/matrix.txt", sep = "\t", header=FALSE))  
m <- t(m)  
dim(m)
```

```
## [1] 5 10
```

Question 2: What are the dimensions of the matrix you just transposed?

Answer 2: 5 rows, 10 columns

###Indexing a Matrix

In the R code chunk below, do the following: 1) Index matrix `m` by selecting all but the third column. 2) Remove the last row of matrix `m`.

```
n <- m[,c(1:2,4:10)]
m <- m[1:4,]
```

6) BASIC DATA VISUALIZATION AND STATISTICAL ANALYSIS

Load Zooplankton Data Set

In the R code chunk below, do the following: 1) Load the zooplankton data set from the **3.RStudio** data folder. 2) Display the structure of this data set.

```
zoop <- as.matrix(read.table("data/zoops.txt", sep = "\t", header=TRUE))
str(zoop)
```

```
## chr [1:24, 1:11] " 5" "14" "16" "21" "23" "25" "27" "34" "12" "15" "18" ...
## - attr(*, "dimnames")=List of 2
## ..$ : NULL
## ..$ : chr [1:11] "TANK" "NUTS" "CAL" "DIAP" ...
```

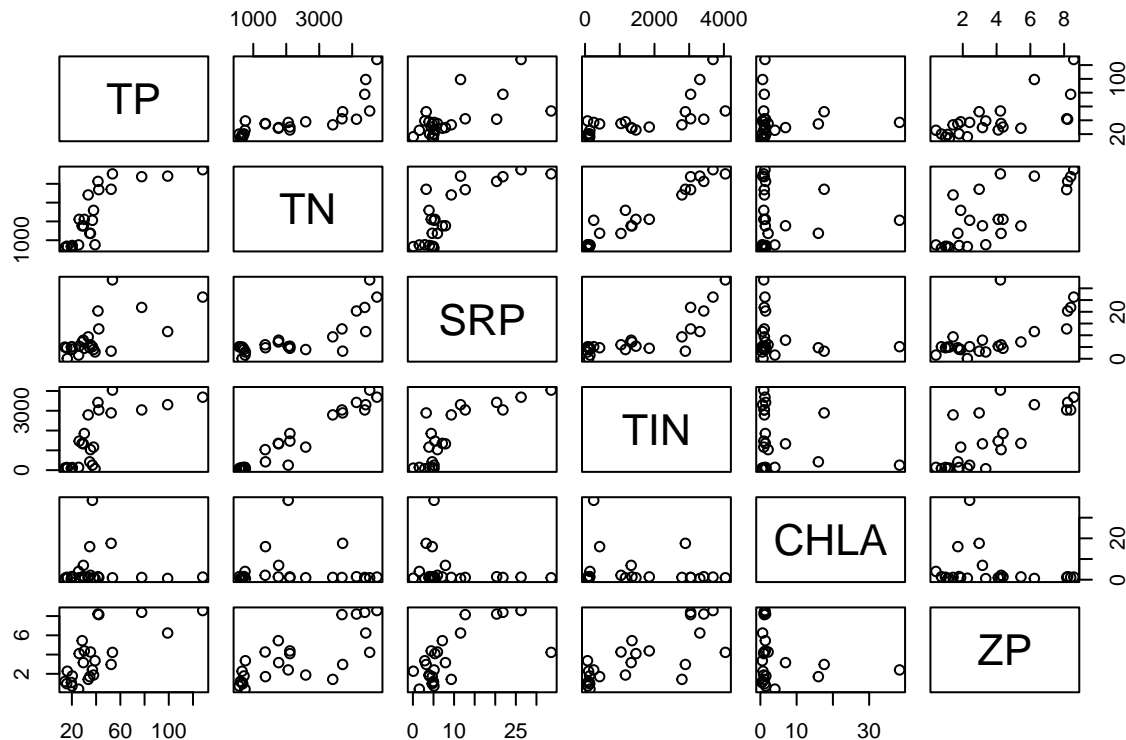
Correlation

In the R-code chunk below, do the following: 1) Create a matrix with the numerical data in the `meso` dataframe. 2) Visualize the pairwise **bi-plots** of the six numerical variables. 3) Conduct a simple **Pearson's correlation** analysis.

```
meso <- as.matrix(read.table("data/zoop_nuts.txt", sep = "\t", header=TRUE))
meso.char <- meso[,3:8]
is.numeric(meso.char)
```

```
## [1] FALSE
```

```
meso.num <- matrix(as.numeric(meso.char),ncol = ncol(meso.char))
help(colnames)
colnames(meso.num) <- c("TP", "TN", "SRP", "TIN", "CHLA", "ZP")
pairs(meso.num)
```



```
cor1 <- cor(meso.num)
```

Question 3: Describe some of the general features based on the visualization and correlation analysis above?

Answer 3: All pairwise variable correlations are fairly high (around .6-.7) except for CHLA does not correlate with the others.

In the R code chunk below, do the following: 1) Redo the correlation analysis using the `corr.test()` function in the `psych` package with the following options: `method = "pearson"`, `adjust = "BH"`. 2) Now, redo this correlation analysis using a non-parametric method. 3) Use the print command from the handout to see the results of each correlation analysis.

```
install.packages('psych', repos="http://cran.rstudio.com/")
```

```
##
## The downloaded binary packages are in
## /var/folders/7j/ntpmcppd2gb8_5v7mlyv0h840000gn/T/RtmpbhV5iI/downloaded_packages
```

```
require('psych')
```

```
## Loading required package: psych
```

```
cor2 <- corr.test(meso.num, method = "pearson", adjust= "BH")
print(cor2, digits=3)
```

```
## Call:corr.test(x = meso.num, method = "pearson", adjust = "BH")
## Correlation matrix
##      TP      TN      SRP      TIN      CHLA      ZP
## TP    1.000  0.787  0.654  0.717 -0.017  0.697
## TN    0.787  1.000  0.784  0.969 -0.004  0.756
## SRP    0.654  0.784  1.000  0.801 -0.189  0.676
## TIN    0.717  0.969  0.801  1.000 -0.157  0.761
## CHLA -0.017 -0.004 -0.189 -0.157  1.000 -0.183
## ZP    0.697  0.756  0.676  0.761 -0.183  1.000
## Sample Size
## [1] 24
## Probability values (Entries above the diagonal are adjusted for multiple tests.)
##      TP      TN      SRP      TIN      CHLA      ZP
## TP    0.000  0.000  0.001  0.000  0.983  0.000
## TN    0.000  0.000  0.000  0.000  0.983  0.000
## SRP    0.001  0.000  0.000  0.000  0.491  0.000
## TIN    0.000  0.000  0.000  0.000  0.536  0.000
## CHLA  0.938  0.983  0.376  0.464  0.000  0.491
## ZP    0.000  0.000  0.000  0.000  0.393  0.000
##
## To see confidence intervals of the correlations, print with the short=FALSE option
```

```
cor.nonp <- corr.test(meso.num, method = "spearman", adjust= "BH")
print(cor.nonp, digits=3)
```

```
## Call:corr.test(x = meso.num, method = "spearman", adjust = "BH")
## Correlation matrix
##      TP      TN      SRP      TIN      CHLA      ZP
## TP    1.000  0.895  0.539  0.761  0.040  0.741
## TN    0.895  1.000  0.647  0.942  0.021  0.748
## SRP    0.539  0.647  1.000  0.726 -0.064  0.627
## TIN    0.761  0.942  0.726  1.000  0.088  0.738
## CHLA  0.040  0.021 -0.064  0.088  1.000 -0.072
## ZP    0.741  0.748  0.627  0.738 -0.072  1.000
## Sample Size
## [1] 24
## Probability values (Entries above the diagonal are adjusted for multiple tests.)
##      TP      TN      SRP      TIN      CHLA      ZP
## TP    0.000  0.000  0.010  0.000  0.914  0.000
## TN    0.000  0.000  0.001  0.000  0.923  0.000
## SRP    0.007  0.001  0.000  0.000  0.884  0.002
## TIN    0.000  0.000  0.000  0.000  0.884  0.000
## CHLA  0.853  0.923  0.767  0.683  0.000  0.884
## ZP    0.000  0.000  0.001  0.000  0.737  0.000
##
## To see confidence intervals of the correlations, print with the short=FALSE option
```

Question 4: Describe what you learned from `corr.test`. Specifically, are the results sensitive to whether you use parametric (i.e., Pearson's) or non-parametric methods? When should one use non-parametric methods instead of parametric methods? With the Pearson's method, is there evidence for false discovery rate due to multiple comparisons? Why is false discovery rate important?

Answer 4: The results are slightly sensitive to parametric vs. nonparametric methods, but overall trends did not change. Nonparametric methods are better for datasets where the distribution of

the data is unknown or the sample size is small. The strong correlations in this data are strong enough that I don't think there was evidence for a false discovery rate, especially because the Benjamini-Hochberg method was applied. Committing a Type 1 error means statistically finding an effect or correlation that doesn't really exist, and false discovery rate is important to keep in mind and avoid because it increases with the number of tests performed.

Linear Regression

In the R code chunk below, do the following: 1) Conduct a linear regression analysis to test the relationship between total nitrogen (TN) and zooplankton biomass (ZP). 2) Examine the output of the regression analysis. 3) Produce a plot of this regression analysis including the following: categorically labeled points, the predicted regression line with 95% confidence intervals, and the appropriate axis labels.

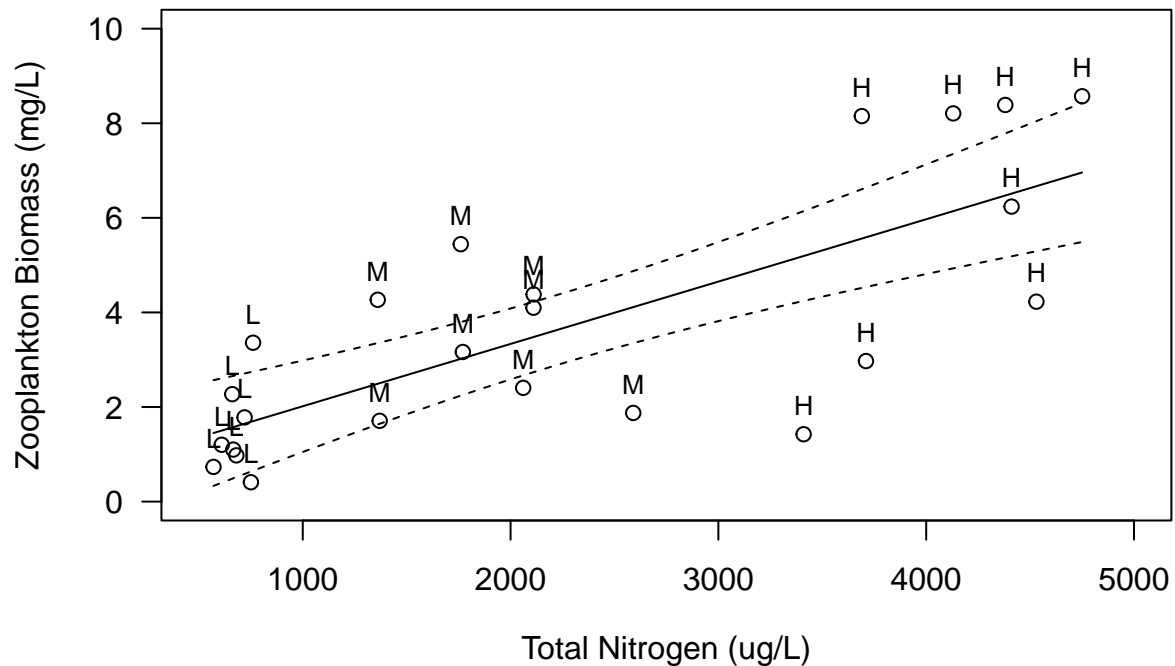
```
mesodf <- as.data.frame(meso)
mesodf$ZP <- as.numeric(as.character(mesodf$ZP))
mesodf$TN <- as.numeric(as.character(mesodf$TN))
fitreg <- lm(as.numeric(as.character(ZP)) ~ as.numeric(as.character(TN)), data=mesodf)
summary(fitreg)

##
## Call:
## lm(formula = as.numeric(as.character(ZP)) ~ as.numeric(as.character(TN)),
##     data = mesodf)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.7690 -0.8491 -0.0709  1.6238  2.5888
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.6977712   0.6496312   1.074   0.294
## as.numeric(as.character(TN)) 0.0013181  0.0002431   5.421 1.91e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.75 on 22 degrees of freedom
## Multiple R-squared:  0.5719, Adjusted R-squared:  0.5525
## F-statistic: 29.39 on 1 and 22 DF,  p-value: 1.911e-05

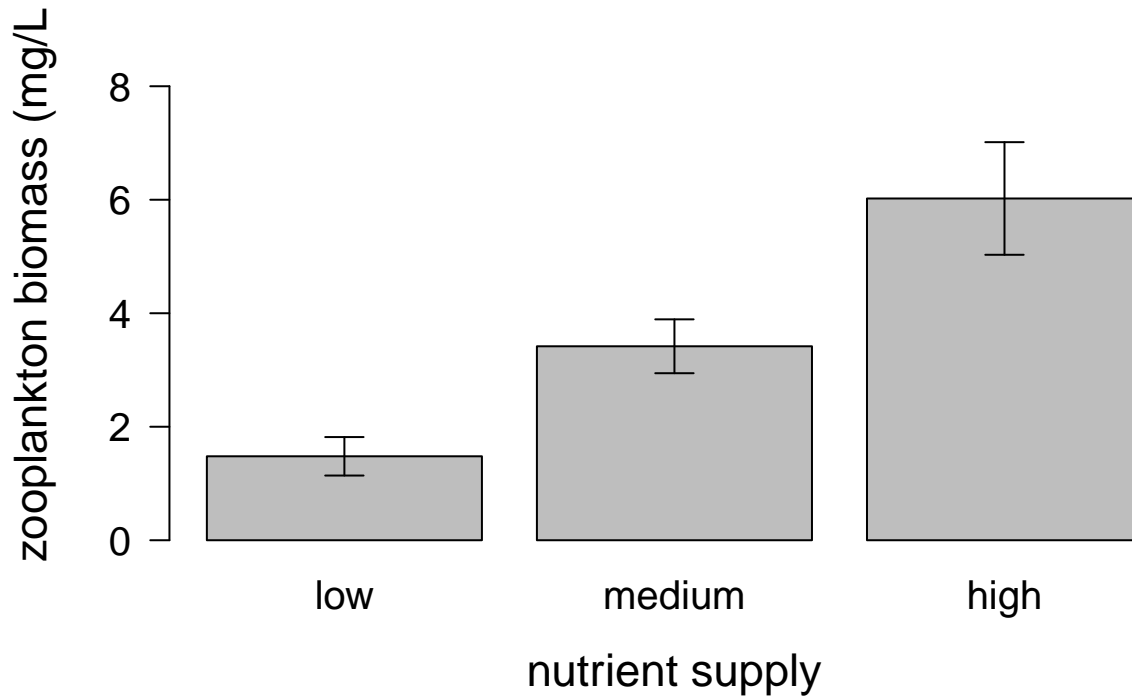
plot(mesodf$TN, mesodf$ZP, ylim = c(0,10), xlim = c(500,5000),
     xlab = expression(paste("Total Nitrogen (ug/L)")),
     ylab = "Zooplankton Biomass (mg/L)", las = 1);
text(mesodf$TN, mesodf$ZP, mesodf$NUTS, pos=3, cex=0.8);
newTN <- seq(min(mesodf$TN), max(mesodf$TN), 10)

regline <- predict(fitreg,newdata=data.frame(TN=newTN))
#error here with predict function "type numeric supplied for TN, but fitted for character"

lines(newTN,regline)
conf95 <- predict(fitreg, newdata = data.frame(TN = newTN),
                  interval = c("confidence"), level=0.95, type="response")
matlines(newTN, conf95[,c("lwr","upr")], type="l", lty=2, lwd=1, col="black")
```

```
bp <- barplot(zp.means, ylim=c(0, round(max(mesodf$ZP), digits=0)),
             pch=15, cex=1.25, las=1, cex.lab=1.4, cex.axis=1.25,
             xlab="nutrient supply",
             ylab="zooplankton biomass (mg/L",
             names.arg=c("low","medium","high"));
arrows(x0=bp, y0=zp.means, y1=zp.means - zp.sem, angle=90,length=0.1, lwd=1);
arrows(x0=bp, y0=zp.means, y1=zp.means + zp.sem, angle=90,length=0.1, lwd=1)
```



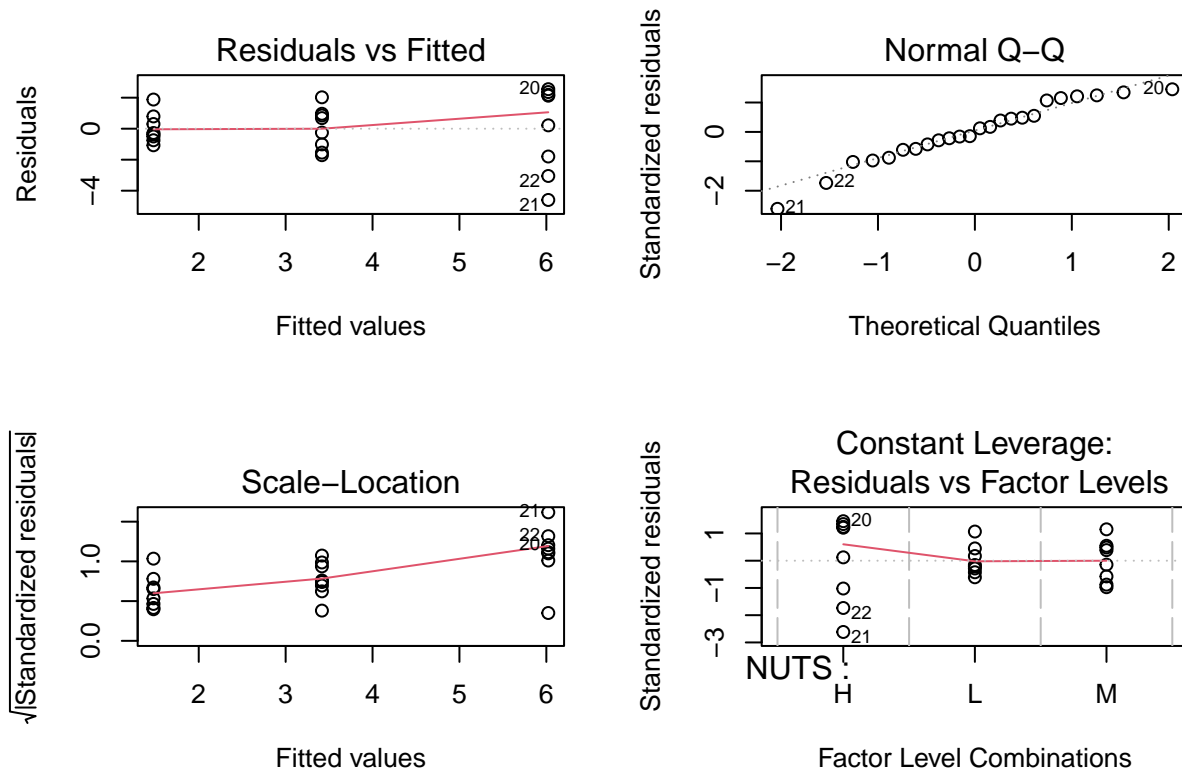
```
fitanova <- aov(ZP ~ NUTS, data=mesodf)
summary(fitanova)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## NUTS        2  83.15   41.58    11.77 0.000372 ***
## Residuals   21  74.16    3.53
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
TukeyHSD(fitanova)
```

```
##      Tukey multiple comparisons of means
##      95% family-wise confidence level
##
## Fit: aov(formula = ZP ~ NUTS, data = mesodf)
##
## $NUTS
##           diff           lwr           upr         p adj
## L-H -4.543175 -6.9115094 -2.1748406 0.0002512
## M-H -2.604550 -4.9728844 -0.2362156 0.0294932
## M-L  1.938625 -0.4297094  4.3069594 0.1220246
```

```
par(mfrow=c(2,2), mar=c(5.1,4.1,4.1,2.1))
plot(fitanova)
```



SYNTHESIS: SITE-BY-SPECIES MATRIX

In the R code chunk below, load the `zoops.txt` data set in your **3.RStudio** data folder. Create a site-by-species matrix (or dataframe) that does *not* include TANK or NUTS. The remaining columns of data refer to the biomass ($\mu\text{g/L}$) of different zooplankton taxa:

- CAL = calanoid copepods
- DIAP = *Diaphanasoma* sp.
- CYL = cyclopoid copepods
- BOSM = *Bosmina* sp.
- SIMO = *Simocephallus* sp.
- CERI = *Ceriodaphnia* sp.
- NAUP = naupuli (immature copepod)
- DLUM = *Daphnia lumholtzi*
- CHYD = *Chydorus* sp.

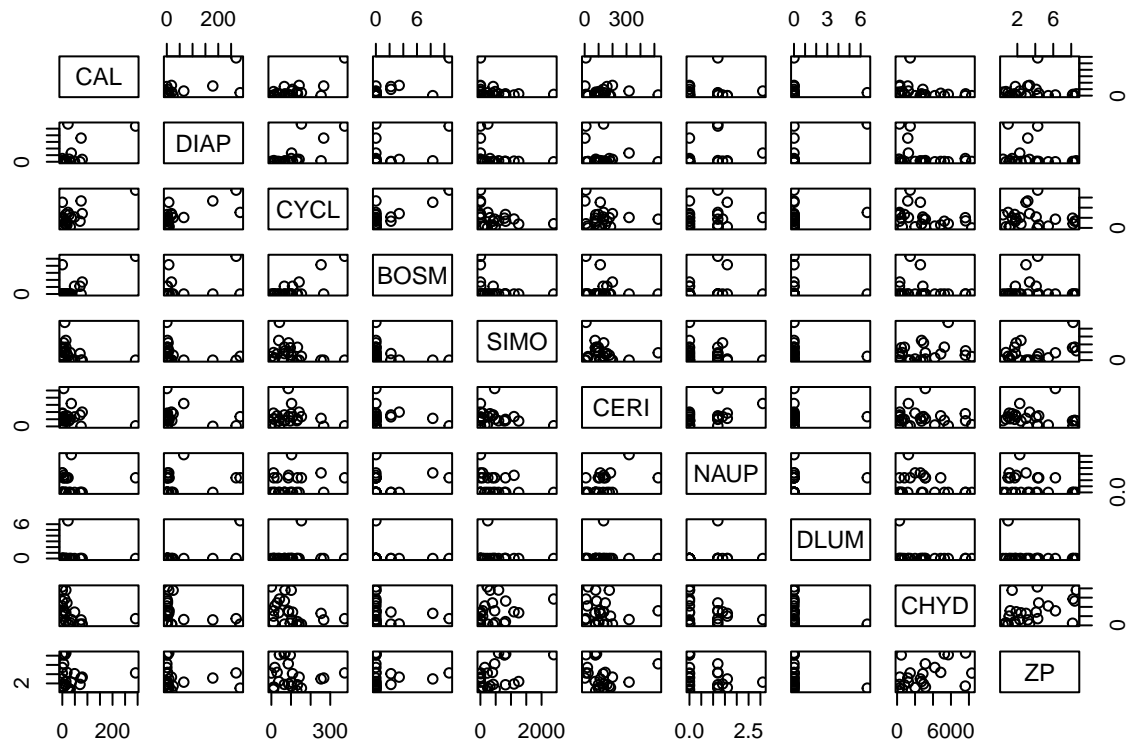
Question 6: With the visualization and statistical tools that we learned about in the **3. RStudio** handout, use the site-by-species matrix to assess whether and how different zooplankton taxa were responsible for the total biomass (ZP) response to nutrient enrichment. Describe what you learned below in the “Answer” section and include appropriate code in the R chunk.

Answer 6: By performing a Pearson's correlation analysis, it appears that *Simocephallus* (0.43) and *Chydorus* (0.46) correlate the most strongly with the overall biomass response to nutrient enrichment; i.e. when these taxa increased relatively more, the overall biomass response increased relatively more. The rest of the taxa showed little correlation with the overall response

```
zoops <- as.matrix(read.table("data/zoops.txt", sep = "\t", header=TRUE))
zoop <- zoops[,3:11]
fulldata <- cbind(zoop,mesodf$ZP)
colnames(fulldata) <- c("CAL","DIAP","CYCL","BOSM","SIMO","CERI","NAUP","DLUM","CHYD","ZP")
fulldata <- as.data.frame(fulldata)
fulldata$ZP <- as.numeric(fulldata$ZP)
fulldata$CAL <- as.numeric(fulldata$CAL)
fulldata$DIAP <- as.numeric(fulldata$DIAP)
fulldata$CYCL <- as.numeric(fulldata$CYCL)
fulldata$BOSM <- as.numeric(fulldata$BOSM)
fulldata$SIMO <- as.numeric(fulldata$SIMO)
fulldata$CERI <- as.numeric(fulldata$CERI)
fulldata$NAUP <- as.numeric(fulldata$NAUP)
fulldata$DLUM <- as.numeric(fulldata$DLUM)
fulldata$CHYD <- as.numeric(fulldata$CHYD)
str(fulldata)
```

```
## 'data.frame': 24 obs. of 10 variables:
## $ CAL : num 70.5 27.1 5.3 79.2 31.4 22.7 0 35.7 74.8 5.3 ...
## $ DIAP: num 0 19.2 8.8 17.9 0 ...
## $ CYCL: num 66.1 129.6 12.7 141.3 11 ...
## $ BOSM: num 2.2 0 0 3.4 0 0 0 0 0 0 ...
## $ SIMO: num 417.8 0 73.1 0 482 ...
## $ CERI: num 159.8 79.4 107.5 199 101.9 ...
## $ NAUP: num 0 0 1.2 0 0 1.2 1.6 3.1 0 1.4 ...
## $ DLUM: num 0 0 0 0 0 6.6 0 0 0 0 ...
## $ CHYD: num 267 159 3158 298 580 ...
## $ ZP : num 1.781 0.409 1.201 3.36 0.733 ...
```

```
pairs(fulldata)
```



```
cor3 <- cor(fulldata)
```

SUBMITTING YOUR WORKSHEET

Use Knitr to create a PDF of your completed **3.RStudio_Worksheet.Rmd** document, push the repo to GitHub, and create a pull request. Please make sure your updated repo include both the PDF and RMarkdown files.

This assignment is due on **Wednesday, January 18th, 2021 at 12:00 PM (noon)**.