

5. Worksheet: Alpha Diversity

Diego Rios; Z620: Quantitative Biodiversity, Indiana University

23 January, 2019

OVERVIEW

In this exercise, we will explore aspects of local or site-specific diversity, also known as alpha (α) diversity. First we will quantify two of the fundamental components of (α) diversity: **richness** and **evenness**. From there, we will then discuss ways to integrate richness and evenness, which will include univariate metrics of diversity along with an investigation of the **species abundance distribution (SAD)**.

Directions:

1. In the Markdown version of this document in your cloned repo, change “Student Name” on line 3 (above) to 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 the proper scripting needed to carry out the exercise.
4. Answer questions in the worksheet. Space for your answer is provided in this document and 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, **push** this file to your GitHub repo.
6. For the assignment portion of the worksheet, follow the directions at the bottom of this file.
7. When you are done, **Knit** the text and code into a PDF file.
8. After Knitting, submit the completed exercise by creating a **pull request** via GitHub. Your pull request should include this file `AlphaDiversity_Worskheet.Rmd` and the PDF output of Knitr (`AlphaDiversity_Worskheet.pdf`).

1) R SETUP

In the R code chunk below, please provide the code to: 1) Clear your R environment, 2) Print your current working directory, 3) Set your working directory to your `5.AlphaDiversity` folder, and 4) Load the **vegan** R package (be sure to install first if you haven’t already).

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

## function ()
## .Internal(getwd())
## <bytecode: 0x7ffc6971b1f8>
## <environment: namespace:base>

setwd("~/GitHub/QB2019_Rios/2.Worksheets/5.AlphaDiversity/")
library(vegan)

## Warning: package 'vegan' was built under R version 3.4.4
## Loading required package: permute
## Loading required package: lattice
```

```
## This is vegan 2.5-3
```

2) LOADING DATA

In the R code chunk below, do the following: 1) Load the BCI dataset, and 2) Display the structure of the dataset (if the structure is long, use the `max.level = 0` argument to show the basic information).

```
data(BCI)
str(BCI, max.level=0)

## 'data.frame':    50 obs. of  225 variables:
##  - attr(*, "original.names")= chr  "Abarema.macradenium" "Acacia.melanoceras" "Acalypha.diversifolia"
```

3) SPECIES RICHNESS

Species richness (S) refers to the number of species in a system or the number of species observed in a sample.

Observed richness

In the R code chunk below, do the following:

1. Write a function called `S.obs` to calculate observed richness
2. Use your function to determine the number of species in `site1` of the BCI data set, and
3. Compare the output of your function to the output of the `specnumber()` function in `vegan`.

```
S.obs <- function(x = ""){
  rowSums(x > 0) * 1
}

S.obs(BCI)

##      1      2      3      4      5      6      7      8      9     10     11     12     13     14     15     16     17     18
##  93    84    90    94   101    85    82    88    90    94    87    84    93    98    93    93    93    89
##  19    20    21    22    23    24    25    26    27    28    29    30    31    32    33    34    35    36
## 109   100    99    91    99    95   105    91    99    85    86    97    77    88    86    92    83    92
##   37    38    39    40    41    42    43    44    45    46    47    48    49    50
##   88    82    84    80   102    87    86    81    81    86   102    91    91    93
```

```
specnumber(BCI)

##      1      2      3      4      5      6      7      8      9     10     11     12     13     14     15     16     17     18
##  93    84    90    94   101    85    82    88    90    94    87    84    93    98    93    93    93    89
##  19    20    21    22    23    24    25    26    27    28    29    30    31    32    33    34    35    36
## 109   100    99    91    99    95   105    91    99    85    86    97    77    88    86    92    83    92
##   37    38    39    40    41    42    43    44    45    46    47    48    49    50
##   88    82    84    80   102    87    86    81    81    86   102    91    91    93
```

```
S.obs(BCI) == specnumber(BCI)

##      1      2      3      4      5      6      7      8      9     10     11     12     13     14     15
## TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
##  16    17    18    19    20    21    22    23    24    25    26    27    28    29    30
## TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
```

```
##   31   32   33   34   35   36   37   38   39   40   41   42   43   44   45
## TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
##   46   47   48   49   50
## TRUE TRUE TRUE TRUE TRUE
```

Question 1: Does `specnumber()` from `vegan` return the same value for observed richness in `site1` as our function `S.obs`? What is the species richness of the first four sites (i.e., rows) of the BCI matrix?

Answer 1: Yes it does. 93, 84, 90, and 94

Coverage: How well did you sample your site?

In the R code chunk below, do the following:

1. Write a function to calculate Good's Coverage, and
2. Use that function to calculate coverage for all sites in the BCI matrix.

```
C <- function(x = ""){
  1 - (sum(x == 1) / sum(x))
}
C(BCI)
```

```
## [1] 0.9190474
```

```
min(C(BCI))
```

```
## [1] 0.9190474
```

```
max(C(BCI))
```

```
## [1] 0.9190474
```

Question 2: Answer the following questions about coverage:

- a. What is the range of values that can be generated by Good's Coverage?
- b. What would we conclude from Good's Coverage if n_i equaled N ?
- c. What portion of taxa in `site1` was represented by singletons?
- d. Make some observations about coverage at the BCI plots.

Answer 2a:

it ranges between 0 and 1 **Answer 2b:**

that all the sampled species are singletons **Answer 2c:**

9% **Answer 2d:**

that it has a good coverage as more than 90% of the species were sampled. The remaining 9% might just be cryptic or very infrequent species.

Estimated richness

In the R code chunk below, do the following:

1. Load the microbial dataset (located in the `5.AlphaDiversity/data` folder),
2. Transform and transpose the data as needed (see handout),
3. Create a new vector (`soilbac1`) by indexing the bacterial OTU abundances of any site in the dataset,
4. Calculate the observed richness at that particular site, and
5. Calculate coverage of that site

```

soilbac <- read.table("~/GitHub/QB2019_Rios/2.Worksheets/5.AlphaDiversity/data/soilbac.txt",header = T,
soilbac.t <- as.data.frame(t(soilbac))
soilbac1 <- soilbac.t[1,]

site1 <- BCI[1,]

S.obs(soilbac1)

## T1_1
## 1074

C(soilbac1)

## [1] 0.6479471

S.obs(site1)

## 1
## 93

C(site1)

## [1] 0.9308036

```

Question 3: Answer the following questions about the soil bacterial dataset.

- How many sequences did we recover from the sample `soilbac1`, i.e. N ?
- What is the observed richness of `soilbac1`?
- How does coverage compare between the BCI sample (`site1`) and the KBS sample (`soilbac1`)?

Answer 3a:

2119

Answer 3b: 1074

Answer 3c:

Coverage was higher in the BCI plot (0.93) against 0.65 in the KBS sample

Richness estimators

In the R code chunk below, do the following:

- Write a function to calculate **Chao1**,
- Write a function to calculate **Chao2**,
- Write a function to calculate **ACE**, and
- Use these functions to estimate richness at `site1` and `soilbac1`.

```

S.chao1 <- function(x = ""){
  S.obs(x) + (sum(x == 1)^2) / (2 * sum(x == 2))
}

S.chao2 <- function(site, SbyS){
  SbyS = as.data.frame(SbyS)
  x = SbyS[site,]
  SbyS.pa <- (SbyS > 0) * 1 #convert the SbyS to presence/absence
  Q1 = sum(colSums(SbyS.pa) == 1) #species observed once
  Q2 = sum(colSums(SbyS.pa) == 2) #species observed twice
  S.chao2 = S.obs(x) + (Q1^2)/(2 * Q2)
}

```

```

    return(S.chao2)
}

S.Ace <- function(x = "", thresh = 10){
  x <- x[x>0]
  S.abund <- length(which(x > thresh))
  S.rare <- length(which(x <= thresh))
  singlt <- length(which(x == 1))
  N.rare <- sum(x[which(x <= thresh)])
  C.ace <- 1 - (singlt / N.rare)
  i <- c(1:thresh)
  count <- function(i, y){
    length(y[y == i])
  }
  a.1 <- sapply(i, count, x)
  f.1 <- (i * (i - 1)) * a.1
  G.ace <- (S.rare/C.ace)*(sum(f.1)/(N.rare*(N.rare-1)))
  S.ace <- S.abund + (S.rare/C.ace) + (singlt/C.ace) * max(G.ace,0)
  return(S.ace)
}

S.chao1(soilbac1)

##      T1_1
## 2628.514

S.chao1(site1)

##      1
## 119.6944

S.chao2(SbyS = t(soilbac), site = 1)

##      T1_1
## 21055.39

S.chao2(SbyS = t(BCI), site = 1)

## Abarema.macradenia
##                      NaN

S.Ace(soilbac1)

## [1] 4465.983

S.Ace(site1)

## [1] 159.3404

```

Question 4: What is the difference between ACE and the Chao estimators? Do the estimators give consistent results? Which one would you choose to use and why?

Answer 4: The main difference between the estimators is that the ACE incorporates other rare species (i.e., with less than 10 individuals). Chao's indexes give more importance to singletons and doubletons. The three indexes provide consistent results. In other words, they are of similar magnitude. I would use the ACE index as it synthesizes more information than Chao's index.

Rarefaction

In the R code chunk below, please do the following:

1. Calculate observed richness for all samples in `soilbac.t`,
2. Determine the size of the smallest sample,
3. Use the `rarefy()` function to rarefy each sample to this level,
4. Plot the rarefaction results, and
5. Add the 1:1 line and label.

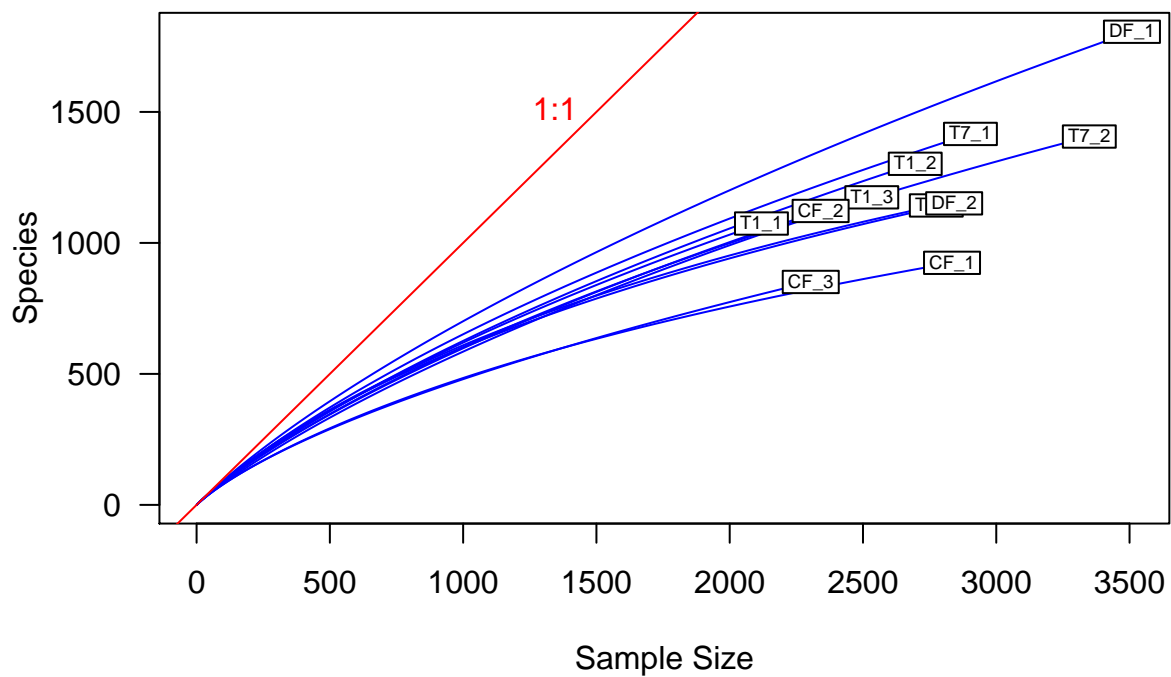
```
S.obs(soilbac.t)

## T1_1 T1_2 T1_3 T7_1 T7_2 T7_3 DF_1 DF_2 CF_1 CF_2 CF_3
## 1074 1302 1174 1416 1406 1143 1806 1151 924 1122 851

min.N <- min(rowSums(soilbac.t))
min.N

## [1] 2119

S.rarefy <- rarefy(x = soilbac.t, sample = min.N, se = TRUE)
rarecurve(x = soilbac.t, step=20, col = "blue", cex = 0.6, las=1)
abline(0, 1, col= "red")
text(1500, 1500, "1:1", pos = 2, col = "red")
```



4) SPECIES EVENNESS

Here, we consider how abundance varies among species, that is, **species evenness**.

Visualizing evenness: the rank abundance curve (RAC)

One of the most common ways to visualize evenness is in a **rank-abundance curve** (sometime referred to as a rank-abundance distribution or Whittaker plot). An RAC can be constructed by ranking species from the most abundant to the least abundant without respect to species labels (and hence no worries about ‘ties’ in abundance).

In the R code chunk below, do the following:

1. Write a function to construct a RAC,
2. Be sure your function removes species that have zero abundances,
3. Order the vector (RAC) from greatest (most abundant) to least (least abundant), and
4. Return the ranked vector

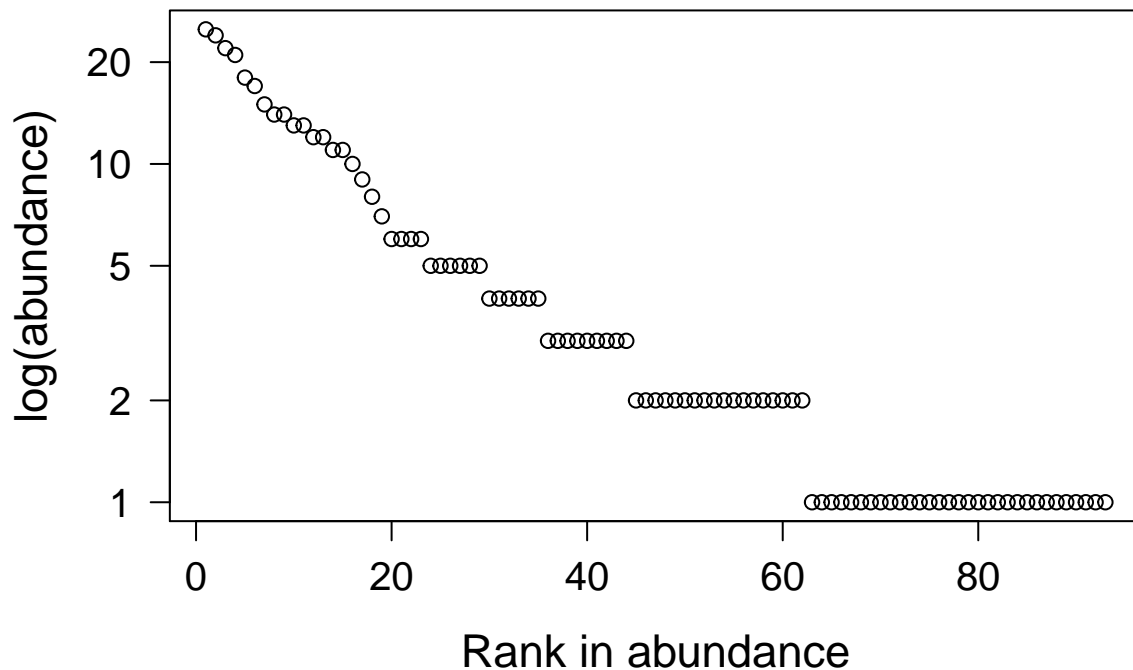
```
RAC <- function(x = ""){  
  x = as.vector(x)  
  x.ab = x[x > 0]  
  x.ab.ranked = x.ab[order(x.ab, decreasing = TRUE)]  
  return(x.ab.ranked)  
}
```

Now, let’s examine the RAC for `site1` of the BCI data set.

In the R code chunk below, do the following:

1. Create a sequence of ranks and plot the RAC with natural-log-transformed abundances,
2. Label the x-axis “Rank in abundance” and the y-axis “log(abundance)”

```
plot.new()  
site1 <- BCI[1,]  
  
rac <- RAC(x = site1)  
ranks <- as.vector(seq(1, length(rac)))  
opar <- par(no.readonly = TRUE)  
par(mar = c(5.1, 5.1, 4.1, 2.1))  
plot(ranks, log(rac), type = "p", axes = F,  
      xlab = "Rank in abundance", ylab = "log(abundance)",  
      las = 1, cex.lab = 1.4, cex.axis = 1.25)  
box()  
axis(side = 1, labels = T, cex.axis = 1.25)  
axis(side = 2, las = 1, cex.axis = 1.25,  
      labels = c(1,2,5,10,20), at = log(c(1,2,5,10,20)))
```



```
par <- opar
```

Question 5: What effect does visualizing species abundance data on a log-scaled axis have on how we interpret evenness in the RAC?

Answer 5: It allows for the visualization -and comparison- of species that have extremely different abundances in a comprehensive way.

Now that we have visualized unevenness, it is time to quantify it using Simpson's evenness ($E_{1/D}$) and Smith and Wilson's evenness index (E_{var}).

Simpson's evenness ($E_{1/D}$)

In the R code chunk below, do the following:

1. Write the function to calculate $E_{1/D}$, and
2. Calculate $E_{1/D}$ for `site1`.

```
SimpE <- function(x = ""){
  S <- S.obs(x)
  x = as.data.frame(x)
  D <- diversity(x, "inv")
  E <- (D)/S
  return(E)
}
SimpE(site1)
```

```
##          1
```



```
## 0.4238232
```

Smith and Wilson's evenness index (E_{var})

In the R code chunk below, please do the following:

1. Write the function to calculate E_{var} ,
2. Calculate E_{var} for `site1`, and
3. Compare $E_{1/D}$ and E_{var} .

```
Evar <- function(x){  
  x <- as.vector(x[x > 0])  
  1 - (2/pi)*atan(var(log(x)))  
}  
Evar(site1)
```

```
## [1] 0.5067211
```

Question 6: Compare estimates of evenness for `site1` of BCI using $E_{1/D}$ and E_{var} . Do they agree? If so, why? If not, why? What can you infer from the results.

Answer 6: Both estimates provide similar values: 0.42 and 0.50 for Simpson's and, Smith and Wilson's indexes, respectively. They differ because of the way both estimates are calculated. Simpson's index gives more importance to highly abundant species, while Smith and Wilson's index corrects this by using the sample variance of the log-transformed abundances.

5) INTEGRATING RICHNESS AND EVENNESS: DIVERSITY METRICS

So far, we have introduced two primary aspects of diversity, i.e., richness and evenness. Here, we will use popular indices to estimate diversity, which explicitly incorporate richness and evenness. We will write our own diversity functions and compare them against the functions in `vegan`.

Shannon's diversity (a.k.a., Shannon's entropy)

In the R code chunk below, please do the following:

1. Provide the code for calculating H' (Shannon's diversity),
2. Compare this estimate with the output of `vegan`'s diversity function using `method = "shannon"`.

```
ShanH <- function(x = ""){  
  H = 0  
  for (n_1 in x){  
    if(n_1 > 0){  
      p = n_1 / sum(x)  
      H = H - p*log(p)  
    }  
  }  
  return(H)  
}  
ShanH(site1)
```

```
## [1] 4.018412
```

```
diversity(site1, index = "shannon")
```

```
## [1] 4.018412
```

Simpson's diversity (or dominance)

In the R code chunk below, please do the following:

1. Provide the code for calculating D (Simpson's diversity),
2. Calculate both the inverse ($1/D$) and $1 - D$,
3. Compare this estimate with the output of **vegan**'s diversity function using method = "simp".

```
SimpD <- function(x = ""){  
  D = 0  
  N = sum(x)  
  for (n_i in x){  
    D = D + (n_i^2)/(N^2)  
  }  
  return(D)  
}
```

```
D.inv <- 1/SimpD(site1)  
D.sub <- 1-SimpD(site1)  
D.inv
```

```
## [1] 39.41555
```

```
D.sub
```

```
## [1] 0.9746293
```

```
diversity(site1, "inv")
```

```
## [1] 39.41555
```

```
diversity(site1, "simp")
```

```
## [1] 0.9746293
```

Question 7: Compare estimates of evenness for **site1** of BCI using E_H' and E_{var} . Do they agree? If so, why? If not, why? What can you infer from the results.

Answer 7: The two values of evenness are similar (Simpson's and Smith and Wilson's). They differ from Shannon's index, in the fact that they measure different aspects of the community. Evenness estimates how similar is the abundance of all species in the community. While the diversity index takes into account identity and abundance, and provides an estimate for the community.

Fisher's α

In the R code chunk below, please do the following:

1. Provide the code for calculating Fisher's α ,
2. Calculate Fisher's α for **site1** of BCI.

```
rac <- as.vector(site1[site1 > 0])
invD <- diversity(rac, "inv")
invD
```

```
## [1] 39.41555
```

```
Fisher <- fisher.alpha(rac)
Fisher
```

```
## [1] 35.67297
```

Question 8: How is Fisher's α different from $E_{H'}$ and E_{var} ? What does Fisher's α take into account that $E_{H'}$ and E_{var} do not?

Answer 8: Fisher's provides an estimate of diversity (i.e., absolute value, from my understanding), not a metric of diversity (relative value). Fisher's alpha takes into account sampling error.

6) MOVING BEYOND UNIVARIATE METRICS OF α DIVERSITY

The diversity metrics that we just learned about attempt to integrate richness and evenness into a single, univariate metric. Although useful, information is invariably lost in this process. If we go back to the rank-abundance curve, we can retrieve additional information – and in some cases – make inferences about the processes influencing the structure of an ecological system.

Species abundance models

The RAC is a simple data structure that is both a vector of abundances. It is also a row in the site-by-species matrix (minus the zeros, i.e., absences).

Predicting the form of the RAC is the first test that any biodiversity theory must pass and there are no less than 20 models that have attempted to explain the uneven form of the RAC across ecological systems.

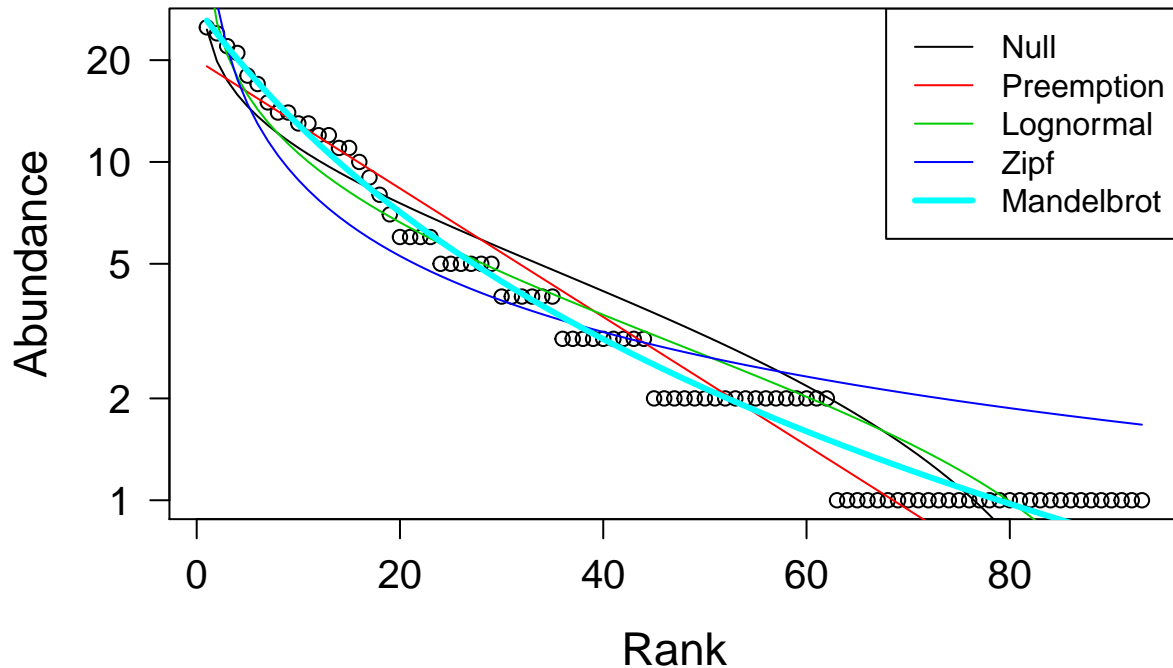
In the R code chunk below, please do the following:

1. Use the `radfit()` function in the `vegan` package to fit the predictions of various species abundance models to the RAC of `site1` in BCI,
2. Display the results of the `radfit()` function, and
3. Plot the results of the `radfit()` function using the code provided in the handout.

```
RACresults <- radfit(site1)
RACresults
```

```
##
## RAD models, family poisson
## No. of species 93, total abundance 448
##
##           par1      par2      par3  Deviance AIC      BIC
## Null                39.5261 315.4362 315.4362
## Preemption 0.042797      21.8939 299.8041 302.3367
## Lognormal  1.0687    1.0186      25.1528 305.0629 310.1281
## Zipf       0.11033 -0.74705      61.0465 340.9567 346.0219
## Mandelbrot 100.52   -2.312    24.084    4.2271 286.1372 293.7350
```

```
plot.new()
plot(RACresults, las = 1, cex.lab = 1.4, cex.axis = 1.25)
```



Question 9: Answer the following questions about the rank abundance curves: a) Based on the output of `radfit()` and plotting above, discuss which model best fits our rank-abundance curve for `site1`? b) Can we make any inferences about the forces, processes, and/or mechanisms influencing the structure of our system, e.g., an ecological community?

Answer 9a: The Mandelbrot model provides the best fit as it has the lowest model deviance, AIC, and BIC. **Answer 9b:** Few species dominate the community, which seemingly impose some sort of restriction on how the rest of species partition the remaining resources.

Question 10: Answer the following questions about the preemption model: a. What does the preemption model assume about the relationship between total abundance (N) and total resources that can be preempted? b. Why does the niche preemption model look like a straight line in the RAD plot?

Answer 10a: that each species in the community requires the same amount of resources from the environment. **Answer 10b:** Because the only fitted parameter is Alpha, which is a constant. Thus, it produces a straight line.

Question 11: Why is it important to account for the number of parameters a model uses when judging how well it explains a given set of data?

Answer 11: Because the number of parameters help explain the adjustment of the model. However, overparameterized models are more difficult to explain, so simpler models are preferred.

SYNTHESIS

1. As stated by Magurran (2004) the $D = \sum p_i^2$ derivation of Simpson's Diversity only applies to communities of infinite size. For anything but an infinitely large community, Simpson's Diversity index

is calculated as $D = \sum \frac{n_i(n_i-1)}{N(N-1)}$. Assuming a finite community, calculate Simpson's D, 1 - D, and Simpson's inverse (i.e. 1/D) for **site 1** of the BCI site-by-species matrix.

```
SimpD2 <- function(x = ""){
  D = 0
  N = sum(x)
  for (n_i in x){
    D = D + (n_i*(n_i - 1))/(N*(N-1))
  }
  return(D)
}

D.inv2 <- 1/SimpD2(site1)
D.sub2 <- 1-SimpD2(site1)
D.inv2
```

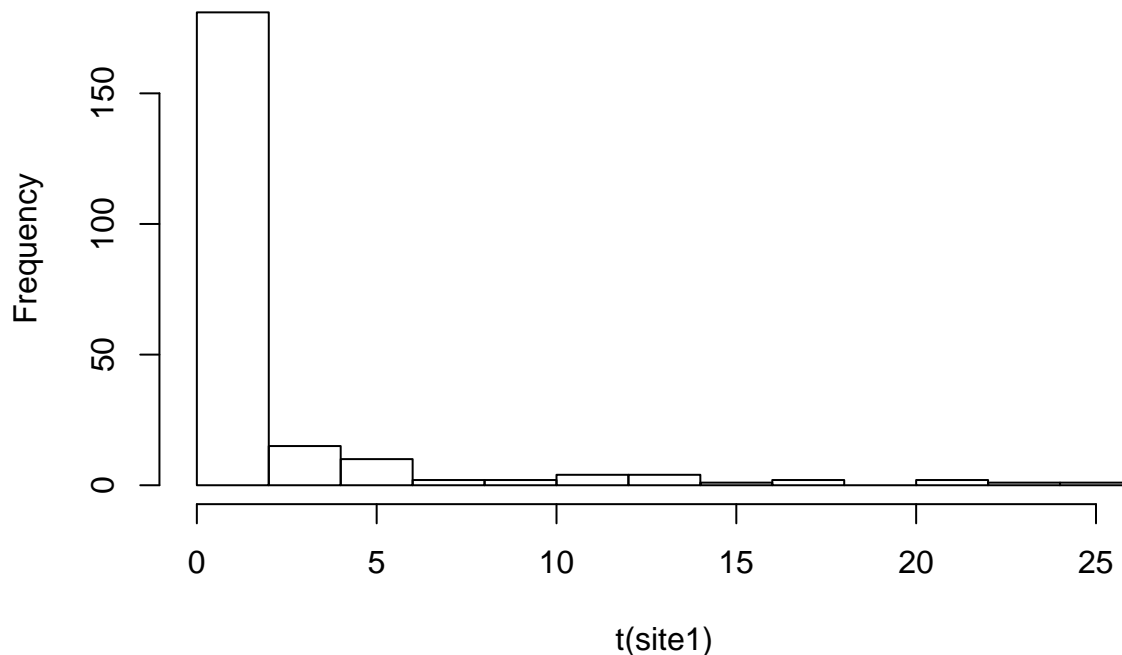
```
## [1] 43.12145
```

```
D.sub2
```

```
## [1] 0.9768097
```

2. Along with the rank-abundance curve (RAC), another way to visualize the distribution of abundance among species is with a histogram (a.k.a., frequency distribution) that shows the frequency of different abundance classes. For example, in a given sample, there may be 10 species represented by a single individual, 8 species with two individuals, 4 species with three individuals, and so on. In fact, the rank-abundance curve and the frequency distribution are the two most common ways to visualize the species-abundance distribution (SAD) and to test species abundance models and biodiversity theories. To address this homework question, use the R function **hist()** to plot the frequency distribution for **site 1** of the BCI site-by-species matrix, and describe the general pattern you see.

```
hist(t(site1), main=NULL)
```



3. We asked you to find a biodiversity dataset with your partner. This data could be one of your own or it could be something that you obtained from the literature. Load that dataset.

```
epiphytes <- read.table("~/GitHub/QB2019_Rios/2.Worksheets/5.AlphaDiversity/data/Epiphytes.txt",header = TRUE)
```

How many sites are there? There are 72 sites

How many species are there in the entire site-by-species matrix? 36 species

Any other interesting observations based on what you learned this week? The communities are epiphytes on trees in an urban landscape in Costa Rica.

SUBMITTING YOUR ASSIGNMENT

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

Unless otherwise noted, this assignment is due on **Wednesday, January 23rd, 2017 at 12:00 PM (noon)**.