

5. Worksheet: Alpha Diversity

Brooke Peckenpaugh; 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()

## [1] "/Users/brooke/GitHub/QB2019_Peckenpaugh/2.Worksheets/5.AlphaDiversity"

setwd("~/GitHub/QB2019_Peckenpaugh/2.Worksheets/5.AlphaDiversity")
#install.packages("vegan")
require("vegan")

## Loading required package: vegan
## 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, ])
```

```
## 1
## 93
```

```
specnumber(BCI[1, ])
```

```
## 1
## 93
```

```
specnumber(BCI[2, ])
```

```
## 2
## 84
```

```
specnumber(BCI[3, ])
```

```
## 3
## 90
```

```
specnumber(BCI[4, ])
```

```
## 4
## 94
```

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. Site 1: 93, Site 2: 84, Site 3: 90, Site 4: 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 - (rowSums(x == 1) / rowSums(x))  
}
```

C(BCI)

```
##          1          2          3          4          5          6          7  
## 0.9308036 0.9287356 0.9200864 0.9468504 0.9287129 0.9174757 0.9326923  
##          8          9         10         11         12         13         14  
## 0.9443155 0.9095355 0.9275362 0.9152120 0.9071038 0.9242054 0.9132420  
##         15         16         17         18         19         20         21  
## 0.9350649 0.9267735 0.8950131 0.9193084 0.8891455 0.9114219 0.8946078  
##         22         23         24         25         26         27         28  
## 0.9066986 0.8705882 0.9030612 0.9095023 0.9115479 0.9088729 0.9198966  
##         29         30         31         32         33         34         35  
## 0.8983516 0.9221053 0.9382423 0.9411765 0.9220183 0.9239374 0.9267887  
##         36         37         38         39         40         41         42  
## 0.9186047 0.9379310 0.9306488 0.9268868 0.9386503 0.8880597 0.9299517  
##         43         44         45         46         47         48         49  
## 0.9140049 0.9168704 0.9234234 0.9348837 0.8847059 0.9228916 0.9086651  
##         50  
## 0.9143519
```

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: The values should range between 0 and 1.

Answer 2b: If n_i equaled N , it would mean that every individual is a singleton species—suggesting that there are many rare species at this site.

Answer 2c: 7% of taxa in site 1 are represented by singletons.

Answer 2d: This was likely good coverage. If the proportion of singletons at your site is that low, it suggests that you were able to find most of the diversity of species at that site.

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("data/soilbac.txt", sep = "\t", header = TRUE, row.names = 1)
soilbac.t <- as.data.frame(t(soilbac))
soilbac1 <- soilbac.t[1,]
```

```
S.obs(soilbac1)
```

```
## T1_1
## 1074
```

```
rowSums(soilbac1)
```

```
## T1_1
## 2119
```

```
C(soilbac1)
```

```
##      T1_1
## 0.6479471
```

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: 45% of samples are represented by singletons. The coverage is not as thorough as the BCI 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
  Q1 = sum(colSums(SbyS.pa) == 1)
  Q2 = sum(colSums(SbyS.pa) == 2)
  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.chao2(1, soilbac.t)

##      T1_1
## 21055.39

S.ace(soilbac1)

## [1] 4465.983

S.chao1(BCI[1, ])

##      1
## 119.6944

S.chao2(1, BCI)

##      1
## 104.6053

S.ace(BCI[1, ])

## [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?

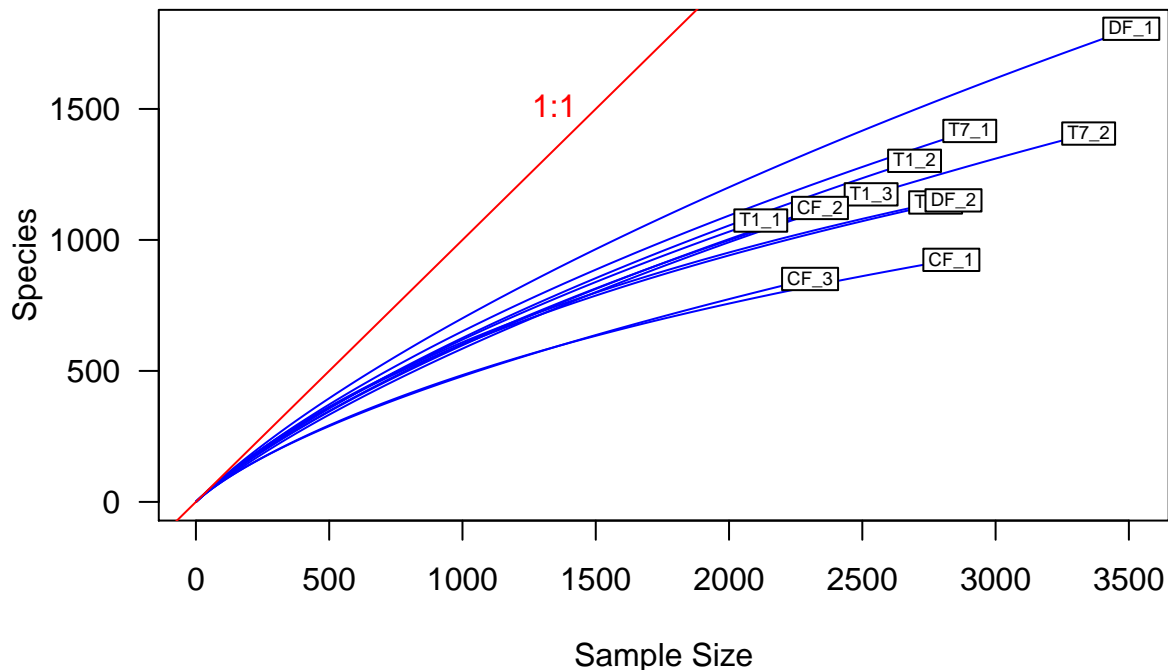
Answer 4: The main difference is that while Chao estimators consider singletons and doubletons, the ACE estimator defines rare species as taxa with 10 or fewer individuals. While the Chao estimators gave fairly consistent results, the ACE estimator is quite different. I would use the Chao estimators for soilbac1, because this dataset has many species where only one individual was sampled in a site. The ACE estimator therefore may be ignoring the majority of sampled species in this case. The estimators are all fairly consistent for the BCI dataset, which makes sense because it is so well sampled.

Rarefaction

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

1. Calculate observed richness for all samples in `soilbac`,
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.

```
soilbac.S <- S.obs(soilbac.t)
min.N <- min(rowSums(soilbac.t))
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

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)"

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: I'm not exactly sure, but I would think that log-transforming the abundance data has the effect of minimizing the range of differences in abundance. This would make it easier to assess the relationship visually, especially if you have many species at the extremes (very low abundance and very high abundance).

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`.

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} .

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:

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"`.

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"`.

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:

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.

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 α takes sampling error into account (the fact that we are not observing every single individual), which the other metrics do not.

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.

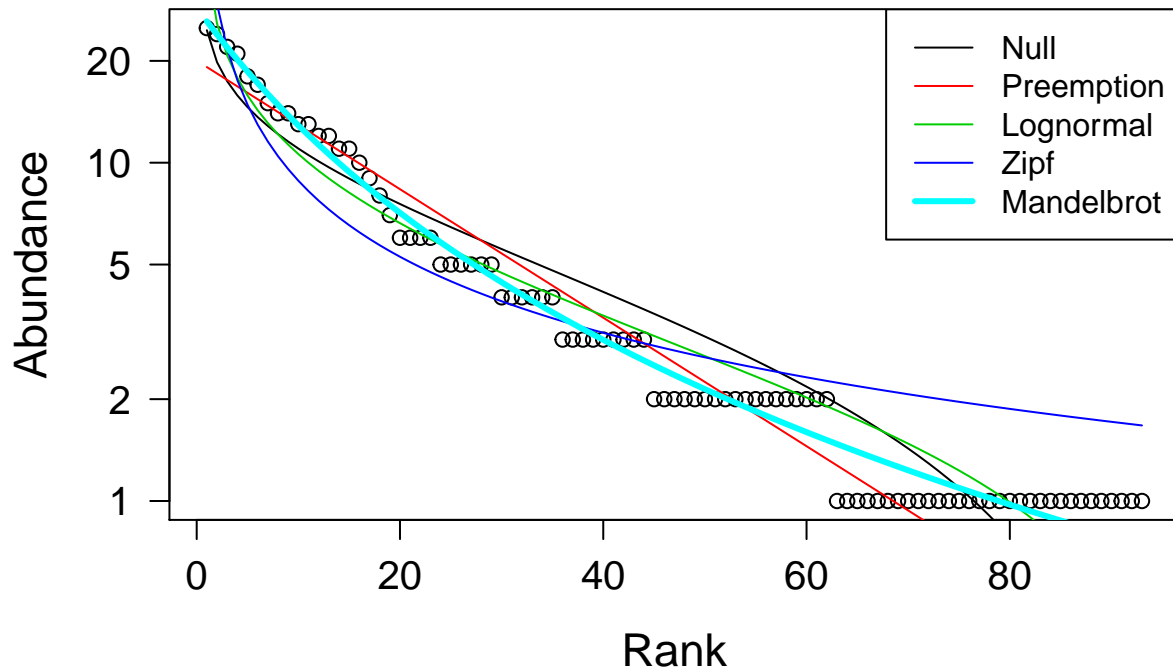
```
RACresults <- radfit(BCI[1, ])  
print(RACresults)
```

```
##  
## 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: Mandelbrot **Answer 9b:** Maybe this would suggest that evenness is an important parameter shaping the structure of this ecological community—specifically, evenness is high among abundant species.

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: The preemption model assumes that total abundance decreases as resources are used. **Answer 10b:** It might look like a straight line because abundance is log-transformed.

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: Increasing the number of parameters used naturally improves the fit of a model to the data, but you also risk over-parameterizing your model, which is not desirable.

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.

```
D <- diversity(BCI[1,], "simp")
D1 <- 1 - D
invD <- 1/D
```

```
print(D)
```

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

```
print(D1)
```

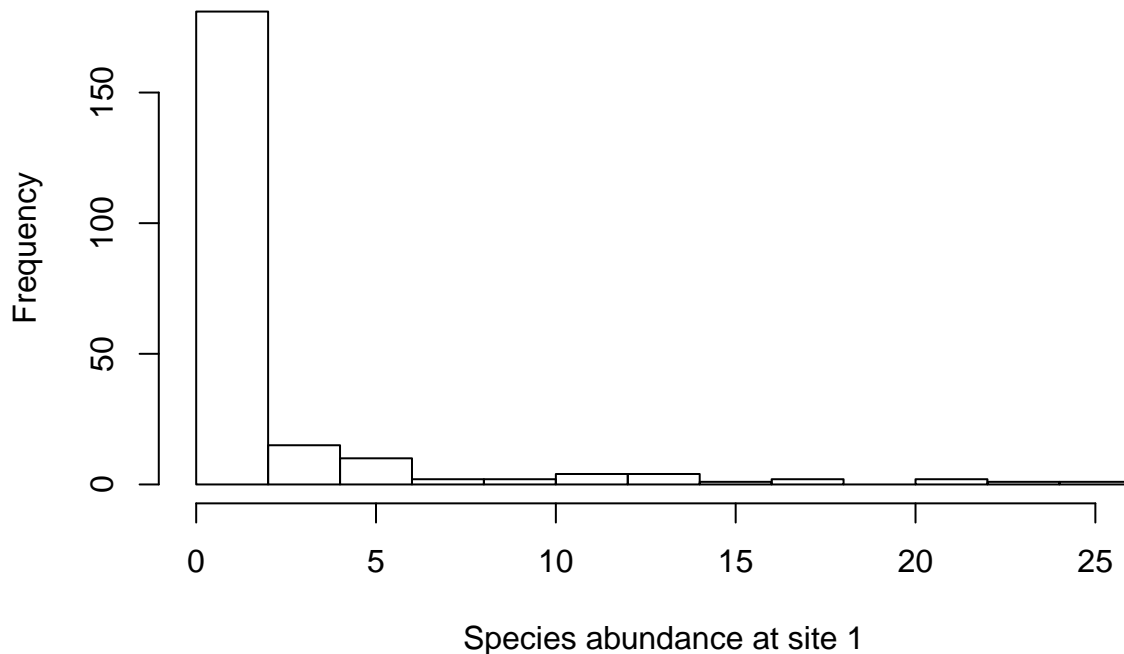
```
## [1] 0.0253707
```

```
print(invD)
```

```
## [1] 1.026031
```

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(as.numeric(BCI[1,]), xlab = "Species abundance at site 1", main="")
```



The data are heavily right skewed. This means that there are many species of low abundance at this site, and a few species of high abundance.

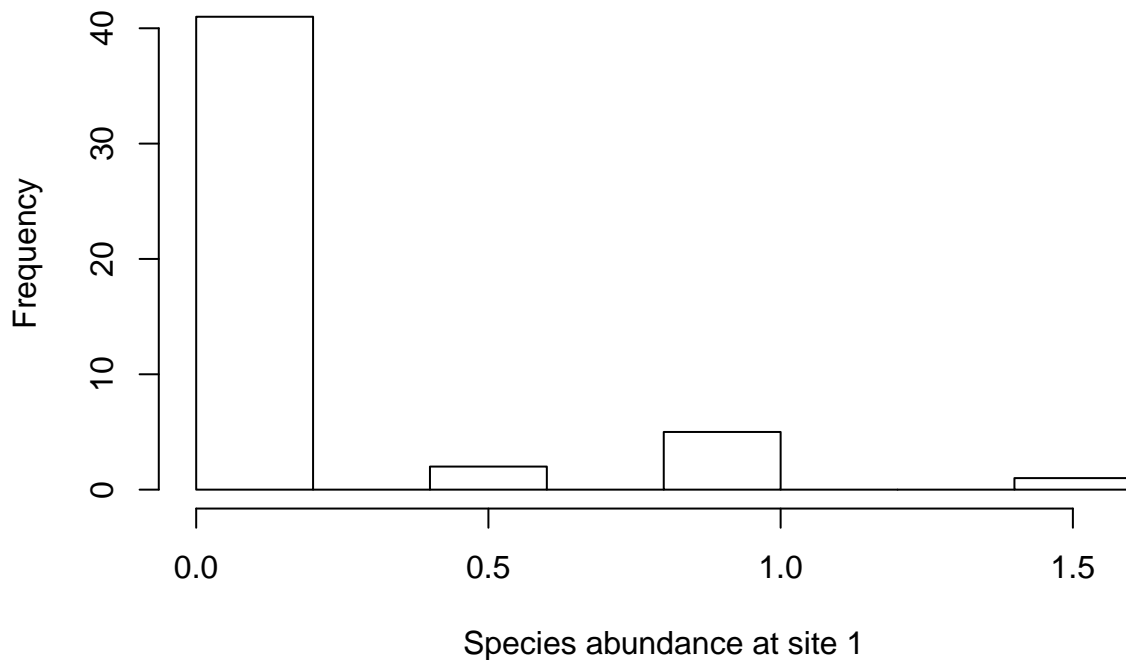
3. We asked you to find a biodiversity dataset with your partner. This data could be one of your own or it could be something that you obtained from the literature. Load that dataset. How many sites are there? How many species are there in the entire site-by-species matrix? Any other interesting observations based on what you learned this week?

```
bird <- read.table("~/Downloads/hf085-01-bird.csv", sep = ",", header=T)
bird[,c(3:ncol(bird))] <- apply(bird[,c(3:ncol(bird))], 1, as.numeric)
```

```
dim(bird)
```

```
## [1] 40 51
```

```
hist(as.numeric(bird[3,c(3:ncol(bird))]), xlab = "Species abundance at site 1", main="")
```



There are 40 sites and 49 species in this site-by-species matrix. Another interesting observation is that the abundance is quantified in a weird way? It seems to be scaled or categorized, which makes the species abundance distribution look weird. So I'll need to figure out what those numbers mean before moving forward with analyzing the data.

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)**.