# 12.Phylogenetic Diversity - Communities

*Andrea Phillips; Z620: Quantitative Biodiversity, Indiana University*

*23 February, 2019*

## OVERVIEW

Complementing taxonomic measures of $\alpha$- and $\beta$-diversity with evolutionary information yields insight into a broad range of biodiversity issues including conservation, biogeography, and community assembly. In this worksheet, you will be introduced to some commonly used methods in phylogenetic community ecology.

After completing this assignment you will know how to:

1. incorporate an evolutionary perspective into your understanding of community ecology
2. quantify and interpret phylogenetic $\alpha$- and $\beta$-diversity
3. evaluate the contribution of phylogeny to spatial patterns of biodiversity

## 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 '8.BetaDiversity' 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 *12.PhyloCom_Worksheet.Rmd* and the PDF output of `Knitr` (*12.PhyloCom_Worksheet.pdf*).

## 1) SETUP

Typically, the first thing you will do in either an R script or an RMarkdown file is setup your environment. This includes things such as setting the working directory and loading any packages that you will need.

In the R code chunk below, provide the code to:
1. clear your R environment,
2. print your current working directory,
3. set your working directory to your `/Week7-PhyloCom` folder,
4. load all of the required R packages (be sure to install if needed), and
5. load the required R source file.

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

## [1] "C:/Users/andjr/Github2/QB2019_Phillips/2.Worksheets/12.PhyloCom"

```
#setwd("~/GitHub/QB2019_Phillips/1.HandOuts/12.PhyloCom")
package.list <- c('picante', 'ape', 'seqinr', 'vegan', 'fossil', 'reshape', 'simba')
for (package in package.list) {
  if (!require(package, character.only = TRUE, quietly = TRUE)) {
    install.packages(package, repos = 'http://cran.us.r-project.org')
    library(package, character.only = TRUE)
  }
}
```

## This is vegan 2.5-4

##
## Attaching package: 'seqinr'

## The following object is masked from 'package:nlme':
##
##     gls

## The following object is masked from 'package:permute':
##
##     getType

## The following objects are masked from 'package:ape':
##
##     as.alignment, consensus

##
## Attaching package: 'shapefiles'

## The following objects are masked from 'package:foreign':
##
##     read.dbf, write.dbf

## This is simba 0.3-5

##
## Attaching package: 'simba'

## The following object is masked from 'package:picante':
##
##     mpd

## The following object is masked from 'package:stats':
##
##     mad

```
source("./bin/MothurTools.R")
```

## 2) DESCRIPTION OF DATA

**need to discuss data set from spatial ecology!**

In 2013 we sampled > 50 forested ponds in Brown County State Park, Yellowwood State Park, and Hoosier National Forest in southern Indiana. In addition to measuring a suite of geographic and environmental variables, we characterized the diversity of bacteria in the ponds using molecular-based approaches. Specifically, we amplified the 16S rRNA gene (i.e., the DNA sequence) and 16S rRNA transcripts (i.e., the RNA transcript of the gene) of bacteria. We used a program called `mothur` to quality-trim our data set and assign sequences to operational taxonomic units (OTUs), which resulted in a site-by-OTU matrix.
In this module we will focus on taxa that were present (i.e., DNA), but there will be a few steps where we need to parse out the transcript (i.e., RNA) samples. See the handout for a further description of this week's dataset.

## 3) LOAD THE DATA

In the R code chunk below, do the following:
1. load the environmental data for the Brown County ponds (*20130801_PondDataMod.csv*),
2. load the site-by-species matrix using the `read.otu()` function,
3. subset the data to include only DNA-based identifications of bacteria,
4. rename the sites by removing extra characters,
5. remove unnecessary OTUs in the site-by-species, and
6. load the taxonomic data using the `read.tax()` function from the source-code file.

```
env <- read.table("data/20130801_PondDataMod.csv", sep = ",", header = TRUE)
env <- na.omit(env)
comm <- read.otu(shared = "./data/INPonds.final.rdp.shared", cutoff = "1")
comm <- comm[grep("*-DNA", rownames(comm)), ]
rownames(comm) <- gsub("\\-DNA", "", rownames(comm))
rownames(comm) <- gsub("\\_", "", rownames(comm))
comm <- comm[rownames(comm) %in% env$Sample_ID, ]
comm <- comm[ , colSums(comm) > 0]
tax <- read.tax(taxonomy = "./data/INPonds.final.rdp.1.cons.taxonomy")
```

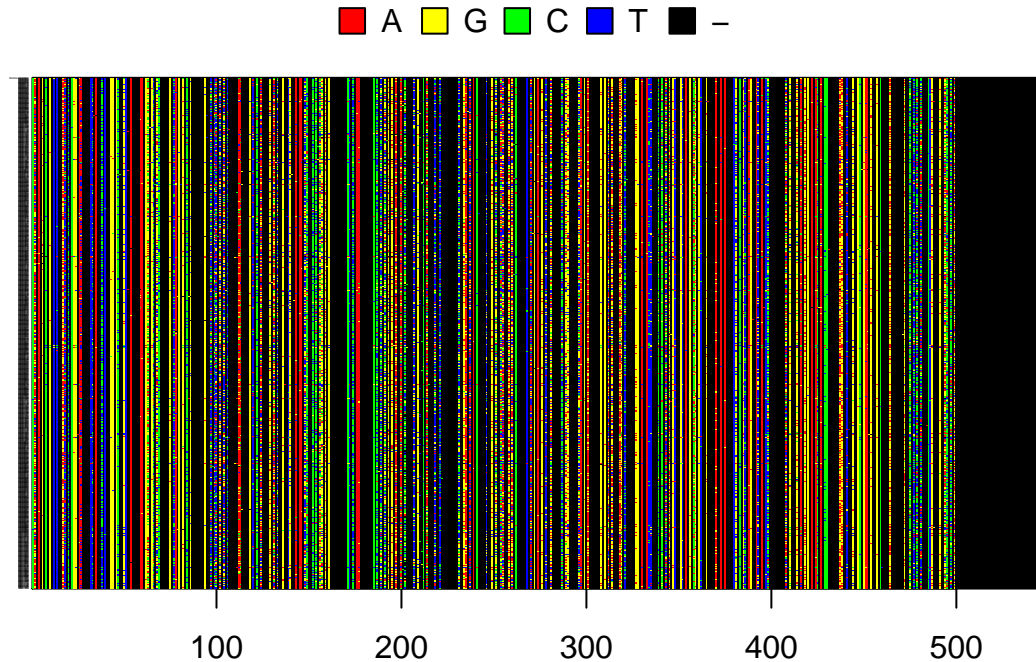Next, in the R code chunk below, do the following:
1. load the FASTA alignment for the bacterial operational taxonomic units (OTUs),
2. rename the OTUs by removing everything before the tab (\t) and after the bar (|),
3. import the *Methanosarcina* outgroup FASTA file,
4. convert both FASTA files into the DNAbin format and combine using `rbind()`,
5. visualize the sequence alignment,
6. using the alignment (with outgroup), pick a DNA substitution model, and create a phylogenetic distance matrix,
7. using the distance matrix above, make a neighbor joining tree,
8. remove any tips (OTUs) that are not in the community data set,
9. plot the rooted tree.

```
ponds.cons <- read.alignment(file = "./data/INPonds.final.rdp.1.rep.fasta",
                            format = "fasta")
ponds.cons$nam <- gsub("\\|.*$", "", gsub("^.*?\t", "", ponds.cons$nam))
```

```
outgroup <- read.alignment(file = "./data/methanosarcina.fasta", format = "fasta")
DNAbin <- rbind(as.DNAbin(outgroup), as.DNAbin(ponds.cons))
image.DNAbin(DNAbin, show.labels = T, cex.lab = 0.05, las = 1)
```



```
seq.dist.jc <- dist.dna(DNAbin, model = "JC", pairwise.deletion = FALSE)
phy.all <- bionj(seq.dist.jc)
phy <- drop.tip(phy.all, phy.all$tip.label[!phy.all$tip.label %in%
                                            c(colnames(comm), "Methanosarcina")])
outgroup <- match("Methanosarcina", phy$tip.label)
phy <- root(phy, outgroup, resolve.root = TRUE)
par(mar = c(1, 1, 2, 1) + 0.1)
plot.phylo(phy, main = "Neighbor Joining Tree", "phylogram", show.tip.label = FALSE, use.edge.length = 
```

# Neighbor Joining Tree



# 4) PHYLOGENETIC ALPHA DIVERSITY

**A. Faith's Phylogenetic Diversity (PD)**

In the R code chunk below, do the following:
1. calculate Faith's D using the **pd()** function.

```r
pd <- pd(comm, phy, include.root = FALSE)
```

In the R code chunk below, do the following:
1. plot species richness (S) versus phylogenetic diversity (PD),
2. add the trend line, and
3. calculate the scaling exponent.

```r
par(mar = c(5, 5, 4, 1) + 0.1)
plot(log(pd$S), log(pd$PD),
     pch = 20, col = "red", las = 1,
     xlab = "ln(S)", ylab = "ln(PD)", cex.main = 1,
     main = "Phylodiversity (PD) vs. Taxonomic richness (S)")
fit <- lm('log(pd$PD) ~ log(pd$S)')
abline(fit, col = "red", lw = 2)
exponent <- round(coefficients(fit)[2], 2)
legend("topleft", legend = paste("Scaling exponent = ", exponent, sep = ""),
       bty = "n", lw = 2, col = "red")
```

## Phylodiversity (PD) vs. Taxonomic richness (S)



*Question 1*: Answer the following questions about the PD-S pattern.
a. Based on how PD is calculated, why should this metric be related to taxonmic richness? b. Describe the relationship between taxonomic richness and phylodiversity. c. When would you expect these two estimates of diversity to deviate from one another? d. Interpret the significance of the scaling PD-S scaling exponent.

> *Answer 1a*: PD is calculated by summing up the branch lengths for each species in the sample–the more taxa that exist in the tree, the higher Faith's Diversity should be, at least how I understand it.
>
> *Answer 1b*: As shown in the figure, as taxonomic richness increases, so does PD; they are positively correlated. *Answer 1c*: These should differ from each other as the branch lengths for taxa increase. If the taxa are all really closely related, the PD value might be low, even though the species diversity is technically high. *Answer 1d*: I think this is just showing the slope of the relationship...? So instead of 1:1 it's basically 3:4, or 0.75:1.

### i. Randomizations and Null Models

In the R code chunk below, do the following:
1. estimate the standardized effect size of PD using the `richness` randomization method.

```
ses.pd <- ses.pd(comm[1:2,], phy, null.model = "richness", runs = 25,
                 include.root = FALSE)
```

*Question 2*: Using `help()` and the table above, run the `ses.pd()` function using two other null models and answer the following questions:

> a. What are the null and alternative hypotheses you are testing via randomization when calculating `ses.pd`?

b. How did your choice of null model influence your observed ses.pd values? Explain why this choice affected or did not affect the output.

**Answer 2a**: The alternative hypothesis we are testing is that the sample is significantly more phylogenetically diverse than what would be expected by chance compared to the null model, which says the phylogenetic diversity is not significantly different. **Answer 2b**: The richness model randomizes community data matrix abundance withint samples, and maintains the sample species richness. Keeping the sample species richness the same when randomizing means the richness of a species is going to be taken into account as we compare the phylogenetic diversity between two sites.

## B. Phylogenetic Dispersion Within a Sample

Another way to assess phylogenetic $\alpha$-diversity is to look at dispersion within a sample.

### i. Phylogenetic Resemblance Matrix

In the R code chunk below, do the following:
1. calculate the phylogenetic resemblance matrix for taxa in the Indiana ponds data set.

```
phydist <- cophenetic.phylo(phy)
```

### ii. Net Relatedness Index (NRI)

In the R code chunk below, do the following:
1. Calculate the NRI for each site in the Indiana ponds data set.

```
ses.mpd <- ses.mpd(comm, phydist, null.model = "taxa.labels",
                   abundance.weighted = TRUE, runs = 25)
NRI <- as.matrix(-1 * ((ses.mpd[,2] - ses.mpd[,3]) / ses.mpd[,4]))
rownames(NRI) <- row.names(ses.mpd)
colnames(NRI) <- "NRI"
```

### iii. Nearest Taxon Index (NTI)

In the R code chunk below, do the following: 1. Calculate the NTI for each site in the Indiana ponds data set.

```
ses.mntd <- ses.mntd(comm, phydist, null.model = "taxa.labels",
                     abundance.weighted = TRUE, runs = 25)
NTI <- as.matrix(-1 * ((ses.mntd[,2] - ses.mntd[,3]) / ses.mntd[,4]))
rownames(NTI) <- row.names(ses.mntd)
colnames(NTI <- "NTI")
```

```
## NULL
```

*Question 3*:

a. In your own words describe what you are doing when you calculate the NRI.
b. In your own words describe what you are doing when you calculate the NTI.
c. Interpret the NRI and NTI values you observed for this dataset.
d. In the NRI and NTI examples above, the arguments "abundance.weighted = FALSE" means that the indices were calculated using presence-absence data. Modify and rerun the code so that NRI and NTI are calculated using abundance data. How does this affect the interpretation of NRI and NTI?

7

***Answer 3a***: When we calculate the NRI we are seeing if the taxa are more or less similar between sites than what is expected by chance when we compare all the way back to the root. ***Answer 3b***: When we calculate the NTI we are seeing it the taxa are more or less similar between sites than what is expected by chance when we compare the closest neighbor, not all the way back to the root like in the NRI. ***Answer 3c***: The NRI values are negative, which means that when we go all the back to the root, it shows that the taxa are overdispersed between sites. For the NTI output I keep getting "NULL", whether the abundance is weighted or not.

***Answer 3d***: The NRI values are negative when abundance is weight. This changes the interpretation to say that the there is more phylogenetic clustering.

## 5) PHYLOGENETIC BETA DIVERSITY

### A. Phylogenetically Based Community Resemblance Matrix

In the R code chunk below, do the following:
1. calculate the phylogenetically based community resemblance matrix using Mean Pair Distance, and
2. calculate the phylogenetically based community resemblance matrix using UniFrac distance.

```
dist.mp <- comdist(comm, phydist)
```

```
## [1] "Dropping taxa from the distance matrix because they are not present in the community data:"
## [1] "Methanosarcina"
```

```
dist.uf <- unifrac(comm, phy)
```
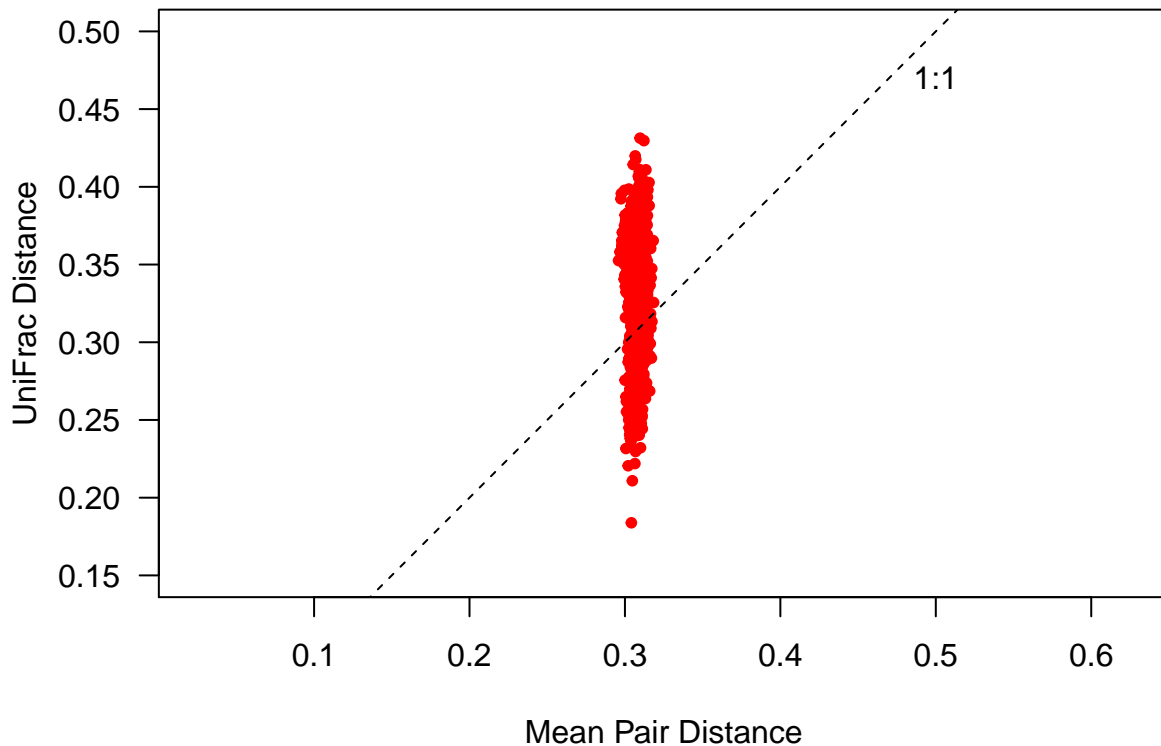
In the R code chunk below, do the following:
1. plot Mean Pair Distance versus UniFrac distance and compare.

```
par(mar = c(5, 5, 2, 1) + 0.1)
plot(dist.mp, dist.uf,
     pch = 20, col = "red", las = 1, asp = 1, xlim = c(0.15, 0.5), ylim = c(0.15, 0.5),
     xlab = "Mean Pair Distance", ylab = "UniFrac Distance")
abline(b = 1, a = 0, lty = 2)
text(0.5, 0.47, "1:1")
```

*Question 4*:

a. In your own words describe Mean Pair Distance, UniFrac distance, and the difference between them.
b. Using the plot above, describe the relationship between Mean Pair Distance and UniFrac distance. Note: we are calculating unweighted phylogenetic distances (similar to incidence based measures). That means that we are not taking into account the abundance of each taxon in each site.
c. Why might MPD show less variation than UniFrac?

**Answer 4a**: Mean Pairwise Distance is basically the phylogenetic distance between two taxa. The UniFrac distance is calculated by taking the number of unshared branches and dividing it by the total number of branches, giving us a measure of difference between two taxa. **Answer 4b**: The Mean Pair Distance in this case does not have a ton of variation, it all sits around 0.3. This means the two taxa are not very phylogenetically distant when looking at means. There is more variation in the UniFrac Distance, meaning just looking at the raw number of shared and unshared branches shows variation in phylogenetic diversity. **Answer 4c**: This difference in variation may exist because even if there is variation in the number of shared branches, it may not change the phylogenetic diversity overall? There is a variety of ways to put together any given tree depending on the sample, so overall the mean pair distance may remain the same even as UniFrac Distance varies.

## B. Visualizing Phylogenetic Beta-Diversity

Now that we have our phylogenetically based community resemblance matrix, we can visualize phylogenetic diversity among samples using the same techniques that we used in the $\beta$-diversity module from earlier in the course.

In the R code chunk below, do the following:
1. perform a PCoA based on the UniFrac distances, and
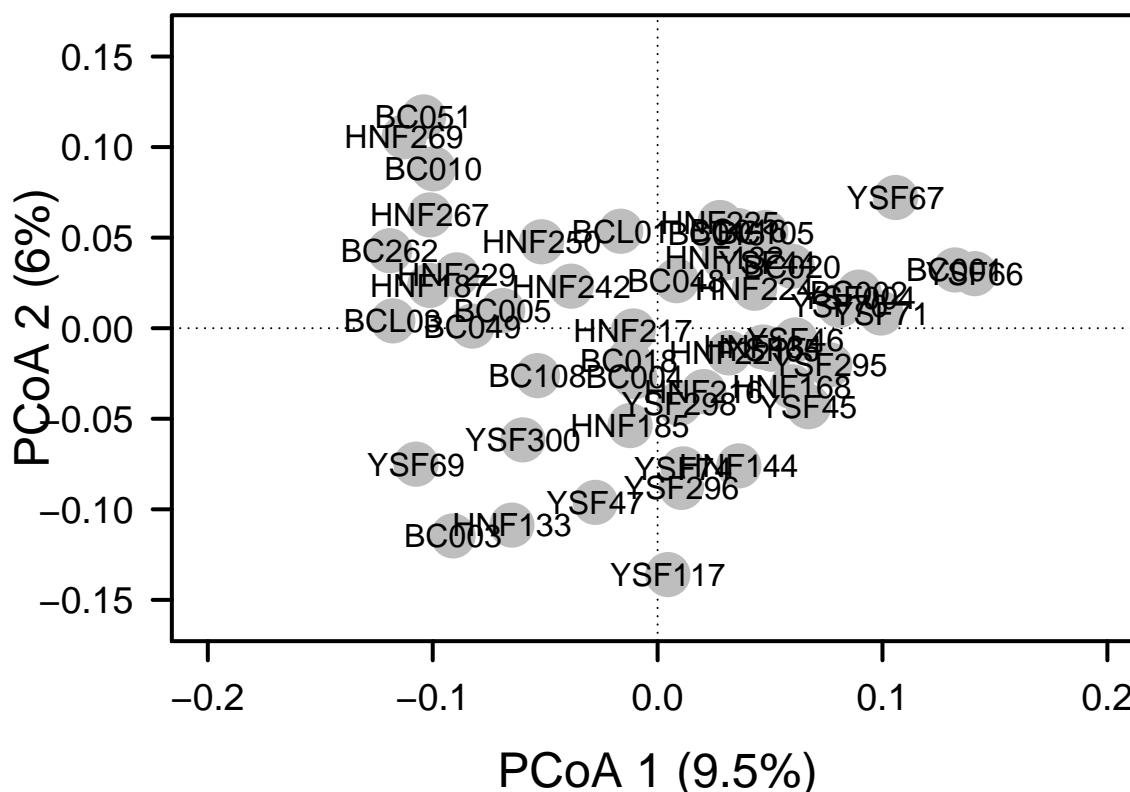2. calculate the explained variation for the first three PCoA axes.

```r
pond.pcoa <- cmdscale(dist.uf, eig = T, k = 3)
explainvar1 <- round(pond.pcoa$eig[1] / sum(pond.pcoa$eig), 3) * 100
explainvar2 <- round(pond.pcoa$eig[2] / sum(pond.pcoa$eig), 3) * 100
explainvar3 <- round(pond.pcoa$eig[3] / sum(pond.pcoa$eig), 3) * 100
sum.eig <- sum(explainvar1, explainvar2, explainvar3)
```

Now that we have calculated our PCoA, we can plot the results.

In the R code chunk below, do the following:
1. plot the PCoA results using either the R base package or the `ggplot` package,
2. include the appropriate axes,
3. add and label the points, and
4. customize the plot.

```r
par(mar = c(5, 5, 1, 2) +0.1)
plot(pond.pcoa$points[ ,1], pond.pcoa$points[ ,2],
     xlim = c(-0.2, 0.2), ylim = c(-.16, 0.16),
     xlab = paste("PCoA 1 (", explainvar1, "%)", sep = ""),
     ylab = paste("PCoA 2 (", explainvar2, "%)", sep = ""),
     pch = 16, cex = 2.0, type = "n", cex.lab = 1.5, cex.axis = 1.2, axes = FALSE)
axis(side = 1, labels = T, lwd.ticks = 2, cex.axis = 1.2, las = 1)
axis(side = 2, labels = T, lwd.ticks = 2, cex.axis = 1.2, las = 1)
abline(h = 0, v = 0, lty = 3)
box(lwd = 2)
points(pond.pcoa$points[ ,1], pond.pcoa$points[ ,2],
       pch = 19, cex = 3, bg = "gray", col = "gray")
text(pond.pcoa$points[ ,1], pond.pcoa$points[ ,2],
     labels = row.names(pond.pcoa$points))
```

In the following R code chunk: 1. perform another PCoA on taxonomic data using an appropriate measure of dissimilarity, and 2. calculate the explained variation on the first three PCoA axes.

```
tax.pcoa <- cmdscale(dist.mp, eig = T, k = 3)
explainvar1 <- round(tax.pcoa$eig[1] / sum(tax.pcoa$eig), 3) * 100
explainvar2 <- round(tax.pcoa$eig[2] / sum(tax.pcoa$eig), 3) * 100
explainvar3 <- round(tax.pcoa$eig[3] / sum(tax.pcoa$eig), 3) * 100
sum.tax <- sum(explainvar1, explainvar2, explainvar3)

par(mar = c(5, 5, 1, 2) +0.1)
plot(tax.pcoa$points[ ,1], tax.pcoa$points[ ,2],
     xlim = c(-0.2, 0.2), ylim = c(-.16, 0.16),
     xlab = paste("PCoA 1 (", explainvar1, "%)", sep = ""),
     ylab = paste("PCoA 2 (", explainvar2, "%)", sep = ""),
     pch = 16, cex = 2.0, type = "n", cex.lab = 1.5, cex.axis = 1.2, axes = FALSE)
axis(side = 1, labels = T, lwd.ticks = 2, cex.axis = 1.2, las = 1)
axis(side = 2, labels = T, lwd.ticks = 2, cex.axis = 1.2, las = 1)
abline(h = 0, v = 0, lty = 3)
box(lwd = 2)
points(tax.pcoa$points[ ,1], tax.pcoa$points[ ,2],
       pch = 19, cex = 3, bg = "gray", col = "gray")
text(tax.pcoa$points[ ,1], tax.pcoa$points[ ,2],
     labels = row.names(tax.pcoa$points))
```
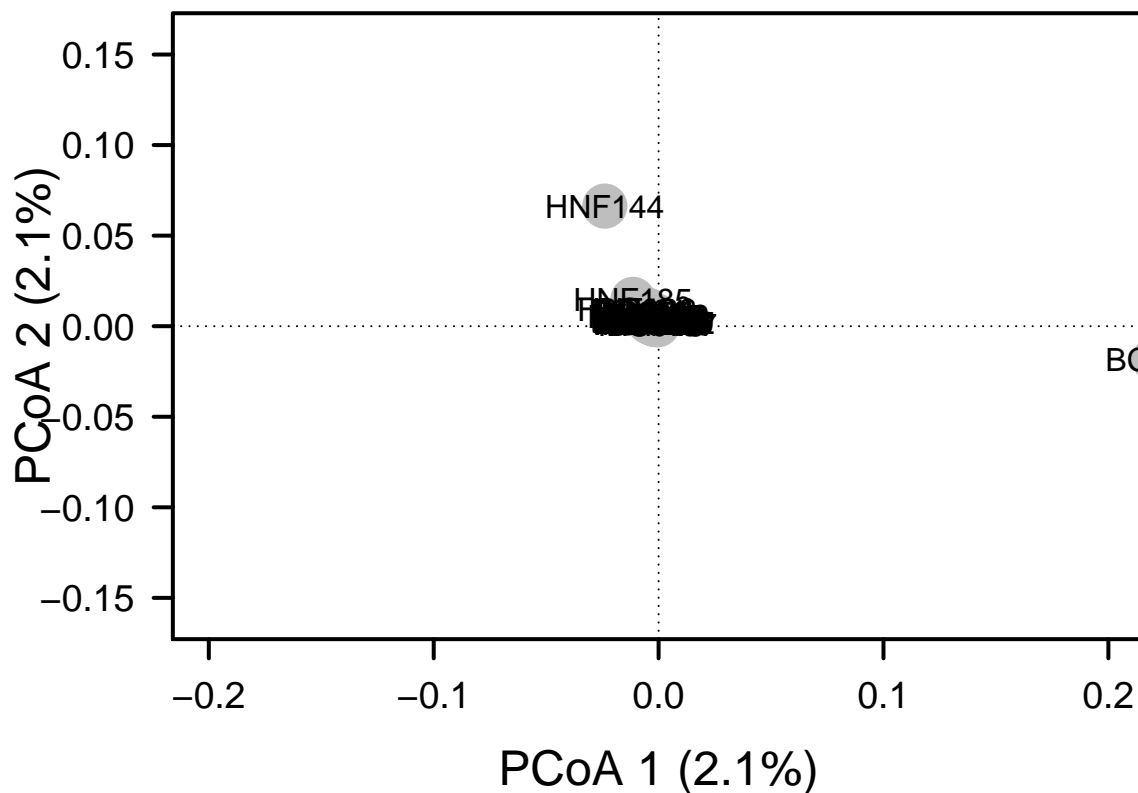
**Question 5**: Using a combination of visualization tools and percent variation explained, how does the phylogenetically based ordination compare or contrast with the taxonomic ordination? What does this tell you about the importance of phylogenetic information in this system?

> **Answer 5**: They are quite different, the taxonomic ordination alone hardly shows any variation in the system, while the phylogenetic ordination clearly shows quite a bit more variation. This shows that phylogenetic information tells a lot more about this system than taxonomic data alone.

## C. Hypothesis Testing

### i. Categorical Approach

In the R code chunk below, do the following:
1. test the hypothesis that watershed has an effect on the phylogenetic diversity of bacterial communities.

```
watershed <- env$Location
adonis(dist.uf ~ watershed, permutations =999)
```

```
##
## Call:
## adonis(formula = dist.uf ~ watershed, permutations = 999)
##
## Permutation: free
## Number of permutations: 999
```

```
## 
## Terms added sequentially (first to last)
## 
##            Df SumsOfSqs  MeanSqs F.Model      R2 Pr(>F)
## watershed  2   0.13316 0.066579  1.2679 0.0492  0.026 *
## Residuals 49   2.57305 0.052511         0.9508
## Total     51   2.70621                  1.0000
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
adonis(
  vegdist(
    decostand(comm, method = "log"),
    method = "bray") ~ watershed,
  permutations = 999)
```

```
## 
## Call:
## adonis(formula = vegdist(decostand(comm, method = "log"), method = "bray") ~      watershed, permuta
## 
## Permutation: free
## Number of permutations: 999
## 
## Terms added sequentially (first to last)
## 
##            Df SumsOfSqs  MeanSqs F.Model      R2 Pr(>F)
## watershed  2   0.16601 0.083003  1.5689 0.06018  0.005 **
## Residuals 49   2.59229 0.052904         0.93982
## Total     51   2.75829                  1.00000
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

### ii. Continuous Approach

In the R code chunk below, do the following: 1. from the environmental data matrix, subset the variables related to physical and chemical properties of the ponds, and
2. calculate environmental distance between ponds based on the Euclidean distance between sites in the environmental data matrix (after transforming and centering using `scale()`).

```
envs <- env[, 5:19]
envs <- envs[, -which(names(envs) %in% c("TDS", "Salinity", "Cal_Volume"))]
env.dist <- vegdist(scale(envs), method = "euclid")
```

In the R code chunk below, do the following:
1. conduct a Mantel test to evaluate whether or not UniFrac distance is correlated with environmental variation.

```
mantel(dist.uf, env.dist)
```

```
## 
## Mantel statistic based on Pearson's product-moment correlation
## 
```

```
## Call:
## mantel(xdis = dist.uf, ydis = env.dist)
##
## Mantel statistic r: 0.1604
##       Significance: 0.063
##
## Upper quantiles of permutations (null model):
##    90%   95% 97.5%    99%
## 0.129 0.179 0.211 0.232
## Permutation: free
## Number of permutations: 999
```

Last, conduct a distance-based Redundancy Analysis (dbRDA).

In the R code chunk below, do the following:
1. conduct a dbRDA to test the hypothesis that environmental variation effects the phylogenetic diversity of bacterial communities,
2. use a permutation test to determine significance, and 3. plot the dbRDA results

```
ponds.dbrda <- vegan::dbrda(dist.uf ~ ., data = as.data.frame(scale(envs)))
anova(ponds.dbrda, by = "axis")
```

```
## Permutation test for dbrda under reduced model
## Forward tests for axes
## Permutation: free
## Number of permutations: 999
##
## Model: vegan::dbrda(formula = dist.uf ~ Elevation + Diameter + Depth + ORP + Temp + SpC + DO + pH + (
##           Df SumOfSqs      F Pr(>F)
## dbRDA1     1  0.10566 2.0152  0.441
## dbRDA2     1  0.09258 1.7658  0.607
## dbRDA3     1  0.07555 1.4409  0.973
## dbRDA4     1  0.06677 1.2735  0.994
## dbRDA5     1  0.05666 1.0807  1.000
## dbRDA6     1  0.05293 1.0095  1.000
## dbRDA7     1  0.04750 0.9059  1.000
## dbRDA8     1  0.03941 0.7517  1.000
## dbRDA9     1  0.03775 0.7201  1.000
## dbRDA10    1  0.03280 0.6256  1.000
## dbRDA11    1  0.02876 0.5485  1.000
## dbRDA12    1  0.02501 0.4770  1.000
## Residual 39  2.04482
```

```
ponds.fit <- envfit(ponds.dbrda, envs, perm = 999)
ponds.fit
```

```
##
## ***VECTORS
##
##             dbRDA1    dbRDA2     r2 Pr(>r)
## Elevation  0.77670  0.62986 0.0959  0.100 .
## Diameter  -0.27972 -0.96008 0.0541  0.240
## Depth     -0.63137  0.77548 0.1756  0.009 **
```

```
## ORP          0.41879 -0.90808 0.1437  0.027 *
## Temp        -0.98250  0.18628 0.1523  0.025 *
## SpC         -0.77101  0.63682 0.2087  0.004 **
## DO          -0.39318 -0.91946 0.0464  0.332
## pH          -0.96210 -0.27270 0.1756  0.016 *
## Color        0.06353  0.99798 0.0464  0.287
## chla        -0.60392 -0.79704 0.2626  0.008 **
## DOC          0.99847 -0.05526 0.0382  0.378
## DON         -0.91633  0.40042 0.0339  0.434
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## Permutation: free
## Number of permutations: 999
```

```r
dbrda.explainvar1 <- round(ponds.dbrda$CCA$eig[1] /
                            sum(c(ponds.dbrda$CCA$eig, ponds.dbrda$CA$eig)), 3) * 100
dbrda.explainvar2 <- round(ponds.dbrda$CCA$eig[2] /
                            sum(c(ponds.dbrda$CCA$eig, ponds.dbrda$CA$eig)), 3) * 100
par(mar = c(5, 5, 4, 4) + 0.1)
plot(scores(ponds.dbrda, display = "wa"), xlim = c(-2, 2), ylim = c(-2, 2),
     xlab = paste("dbRDA 1 (", dbrda.explainvar1, "%)", sep = ""),
     ylab = paste("dbRDA 2 (", dbrda.explainvar2, "%)", sep = ""),
     pch = 16, cex = 2.0, type = "n", cex.lab = 1.5, cex.axis =1.2, axes = FALSE)
axis(side = 1, labels = T, lwd.ticks = 2, cex.axis = 1.2, las = 1)
axis(side = 2, labels = T, lwd.ticks = 2, cex.axis = 1.2, las = 1)
abline(h = 0, v = 0, lty = 3)
box(lwd = 2)
points(scores(ponds.dbrda, display = "wa"),
       pch = 19, cex = 3, bg = "gray", col = "gray")
text(scores(ponds.dbrda, display = "wa"),
     labels = row.names(scores(ponds.dbrda, display = "wa")), cex = 0.5)
vectors <- scores(ponds.dbrda, display = "bp")
arrows(0, 0, vectors[,1] * 2, vectors[, 2] * 2,
       lwd = 2, lty = 1, length = 0.2, col = "red")
text(vectors[,1] * 2, vectors[, 2] * 2, pos = 3,
     labels = row.names(vectors))
axis(side = 3, lwd.ticks = 2, cex.axis = 1.2, las = 1, col = "red", lwd = 2.2,
     at = pretty(range(vectors[, 1])) * 2, labels = pretty(range(vectors[, 1])))
axis(side = 4, lwd.ticks = 2, cex.axis = 1.2, las= 1, col = "red", lwd = 2.2,
     at = pretty(range(vectors[, 2])) * 2, labels = pretty(range(vectors[, 2])))
```

**Question 6**: Based on the multivariate procedures conducted above, describe the phylogenetic patterns of β-diversity for bacterial communities in the Indiana ponds.

> **Answer 6**: It appears that the environmental factors are each a contributor to the probable appearance of bacterial species in the communities. To give a more extreme example, BC051 (top left) is strongly negatively correlated with ORP, DOC, and Elevation, and strongly positively correlated with Temp, SpC, and Depth. BC049 (bottom left) is strongly negatively correlated with Color, Elevation, and DOC, and strongly positively correlated with chla and pH. Overall, the environmental features of each pond are going to fairly significantly impact the community composition of the bacteria.

## 6) SPATIAL PHYLOGENETIC COMMUNITY ECOLOGY

### A. Phylogenetic Distance-Decay (PDD)

A distance decay (DD) relationship reflects the spatial autocorrelation of community similarity. That is, communities located near one another should be more similar to one another in taxonomic composition than distant communities. (This is analagous to the isolation by distance (IBD) pattern that is commonly found when examining genetic similarity of a populations as a function of space.) Historically, the two most common explanations for the taxonomic DD are that it reflects spatially autocorrelated environmental variables and the influence of dispersal limitation. However, if phylogenetic diversity is also spatially autocorrelated, then evolutionary history may also explain some of the taxonomic DD pattern. Here, we will construct the phylogenetic distance-decay (PDD) relationship

First, calculate distances for geographic data, taxonomic data, and phylogenetic data among all unique pair-wise combinations of ponds.

In the R code chunk below, do the following:
1. calculate the geographic distances among ponds,
2. calculate the taxonomic similarity among ponds,
3. calculate the phylogenetic similarity among ponds, and
4. create a dataframe that includes all of the above information.

```r
long.lat <- as.matrix(cbind(env$long, env$lat))
coord.dist <- earth.dist(long.lat, dist = TRUE)
bray.curtis.dist <- 1 - vegdist(comm)
unifrac.dist <- 1 - dist.uf
unifrac.dist.ls <- liste(unifrac.dist, entry = "unifrac")
bray.curtis.dist.ls <- liste(bray.curtis.dist, entry = "bray.curtis")
coord.dist.ls <- liste(coord.dist, entry = "geo.dist")
env.dist.ls <- liste(env.dist, entry = "env.dist")
df <- data.frame(coord.dist.ls, bray.curtis.dist.ls[, 3], unifrac.dist.ls[, 3],
                 env.dist.ls[, 3])
names(df)[4:6] <- c("bray.curtis", "unifrac", "env.dist")
```

Now, let's plot the DD relationships:
In the R code chunk below, do the following:
1. plot the taxonomic distance decay relationship,
2. plot the phylogenetic distance decay relationship, and
3. add trend lines to each.

```r
par(mfrow=c(2, 1), mar = c(1, 5, 2, 1) + 0.1, oma = c(2, 0, 0, 0))
plot(df$geo.dist, df$bray.curtis, xlab = "", xaxt = "n", las = 1, ylim = c(0.1, 0.9),
     ylab = "Bray-Curtis Similarity",
     main = "Distance Decay", col = "SteelBlue")
DD.reg.bc <- lm(df$bray.curtis ~ df$geo.dist)
summary(DD.reg.bc)
```

```
##
## Call:
## lm(formula = df$bray.curtis ~ df$geo.dist)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.31151 -0.08843  0.00315  0.09121  0.43817
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.4463453  0.0066883  66.735   <2e-16 ***
## df$geo.dist -0.0013051  0.0005864  -2.226   0.0262 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1303 on 1324 degrees of freedom
## Multiple R-squared:  0.003728,   Adjusted R-squared:  0.002975
## F-statistic: 4.954 on 1 and 1324 DF,  p-value: 0.0262
```

```r
abline(DD.reg.bc, col = "red4", lwd = 2)
par(mar = c(2, 5, 1, 1) + 0.1)
plot(df$geo.dist, df$unifrac, xlab = "", las = 1, ylim = c(0.1, 0.9),
```

```
      ylab = "Unifrac Similarity", col = "darkorchid4")
DD.reg.uni <- lm(df$unifrac ~ df$geo.dist)
summary(DD.reg.uni)
```

```
##
## Call:
## lm(formula = df$unifrac ~ df$geo.dist)
##
## Residuals:
##       Min        1Q     Median        3Q       Max
## -0.105629 -0.027107 -0.000077  0.026761  0.140215
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept) 0.6735186  0.0019206 350.677   <2e-16 ***
## df$geo.dist 0.0002976  0.0001684   1.767   0.0774 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.03741 on 1324 degrees of freedom
## Multiple R-squared:  0.002354,   Adjusted R-squared:  0.0016
## F-statistic: 3.124 on 1 and 1324 DF,  p-value: 0.07738
```
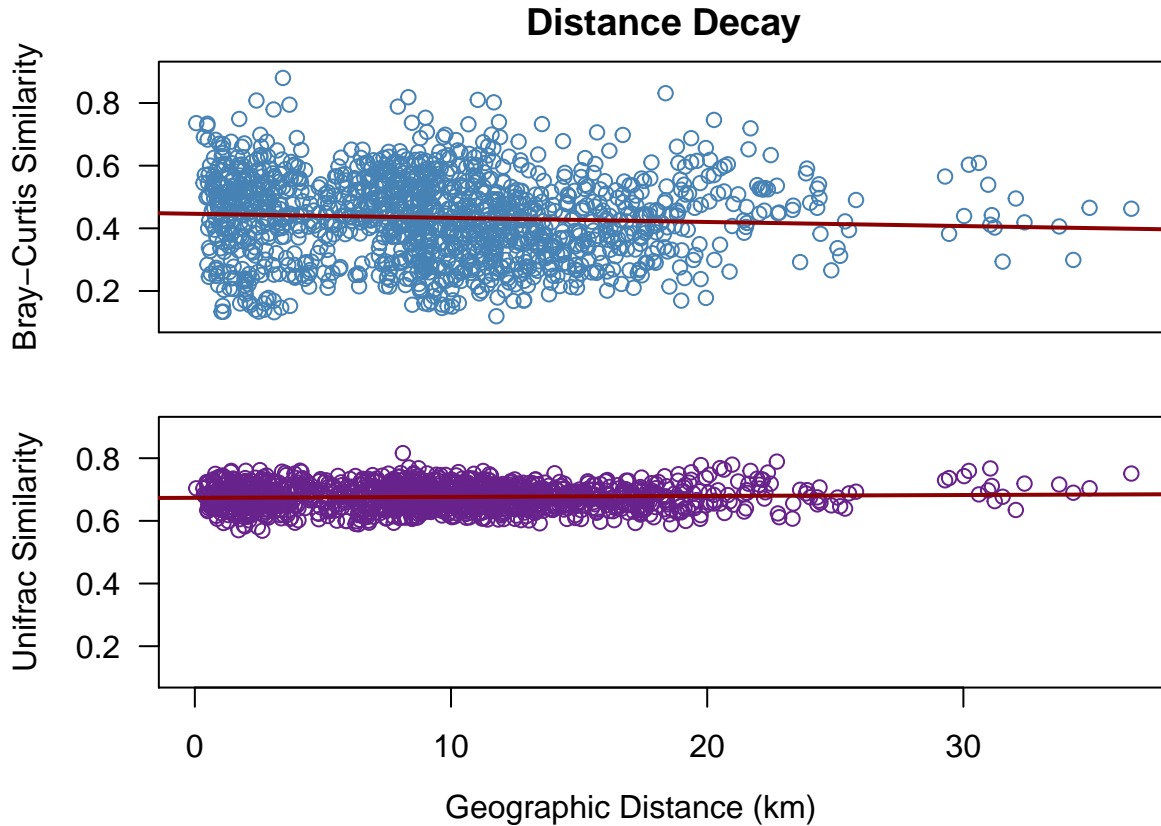
```
abline(DD.reg.uni, col = "red4", lwd = 2)
mtext("Geographic Distance (km)", side = 1, adj = 0.55,
      line = 0.5, outer = TRUE)
```

**Distance Decay**

In the R code chunk below, test if the trend lines in the above distance decay relationships are different from one another.

```
diffslope(df$geo.dist, df$unifrac, df$geo.dist, df$bray.curtis)
```

```
##
## Is difference in slope significant?
## Significance is based on 1000 permutations
##
## Call:
## diffslope(x1 = df$geo.dist, y1 = df$unifrac, x2 = df$geo.dist,     y2 = df$bray.curtis)
##
## Difference in Slope: 0.001603
## Significance: 0.005
##
## Empirical upper confidence limits of r:
##      90%      95%    97.5%      99%
## 0.000739 0.000954 0.001220 0.001476
```

***Question 7***: Interpret the slopes from the taxonomic and phylogenetic DD relationships. If there are differences, hypothesize why this might be.

> ***Answer 7***: The difference in slope is small, only .002, but this is significant with a p-value of .004. This may be because evolutionary history is also impacting the community composition, and not just the change in environment and dispersal limitations as distance increases between the bogs.

# SYNTHESIS

Ignoring technical or methodological constraints, discuss how phylogenetic information could be useful in your own research. Specifically, what kinds of phylogenetic data would you need? How could you use it to answer important questions in your field? In your response, feel free to consider not only phylogenetic approaches related to phylogenetic community ecology, but also those we discussed last week in the PhyloTraits module, or any other concepts that we have not covered in this course.

>  ***Answer Synthesis***: Well, seeing as my field is education, the answer to this question is a little different! Phylogenetics is a really important tool in biology education when it comes to tree thinking–understanding how to read/interpret phylogenetic relationships through trees is an important lense for also understanding the mechanisms and applications of evolution. Additionally, considering the depth of understanding we are continuing to gain about the molecular mechanisms of evolution, understanding of phylogenetics is a vital part of preparing to be a 21st-century scientist, and even just being a scientifically literate citizen, a major goal of science education. Just last week I was reading about a looming issue in science education where the current questions being asked by scientists are at a level of technicality that would be extremely difficult to address at a secondary or even an undergraduate level. The focus of education needs to shift away from teaching a lot of content to teaching students how to reason scientifically and interpret data. Tree thinking is a great example of this. Research has shown that content retention from class to class over time is minimal anyways, so we may as well focus on developing scientific reasoning (identifying and controlling variables, proportional reasoning, hypothetic-deductive reasoning, etc.) that is more likely to be retained over time and create scientifically literate citizens, as well as students who are prepared to tackle their science interests upon starting college.