

## 6. Worksheet: Among Site (Beta) Diversity – Part 1

Yongsoo Choi; Z620: Quantitative Biodiversity, Indiana University

05 February, 2025

### OVERVIEW

In this worksheet, we move beyond the investigation of within-site  $\alpha$ -diversity. We will explore  $\beta$ -diversity, which is defined as the diversity that occurs among sites. This requires that we examine the compositional similarity of assemblages that vary in space or time.

After completing this exercise you will know how to:

1. formally quantify  $\beta$ -diversity
2. visualize  $\beta$ -diversity with heatmaps, cluster analysis, and ordination
3. test hypotheses about  $\beta$ -diversity using multivariate statistics

### 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, **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 (**6.BetaDiversity\_1\_Worksheet.Rmd**) with all code blocks filled out and questions answered) and the PDF output of Knitr (**6.BetaDiversity\_1\_Worksheet.pdf**).

The completed exercise is due on **Wednesday, February 5<sup>th</sup>, 2025 before 12:00 PM (noon)**.

### 1) R 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, please provide the code to:

- 1) Clear your R environment,
- 2) Print your current working directory,
- 3) Set your working directory to your **Week3-Beta/** folder folder, and
- 4) Load the **vegan** R package (be sure to install first if you have not already).

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

## [1] "/cloud/project/QB2025_Choi/Week3-Beta"
setwd("/cloud/project/QB2025_Choi/Week3-Beta/")
package.list <- c('vegan', 'ade4', 'viridis', 'gplots', 'BiodiversityR', 'indicspecies')
for (package in package.list) {
  if (!require(package, character.only = TRUE, quietly = TRUE)){
    install.packages(package)
    library(package, character.only = TRUE)
  }
}

## This is vegan 2.6-8
##
## Attaching package: 'gplots'
## The following object is masked from 'package:stats':
##
##      lowess
## Warning in fun(libname, pkgname): couldn't connect to display ":0"
## BiodiversityR 2.17-1.1: Use command BiodiversityRGUI() to launch the Graphical User Interface;
## to see changes use BiodiversityRGUI(changeLog=TRUE, backward.compatibility.messages=TRUE)
```

## 2) LOADING DATA

### Load dataset

In the R code chunk below, do the following:

1. load the `doubs` dataset from the `ade4` package, and
2. explore the structure of the dataset.

```
data(doubs)
str(doubs, max.level = 1)

## List of 4
## $ env      : 'data.frame': 30 obs. of  11 variables:
## $ fish      : 'data.frame': 30 obs. of  27 variables:
## $ xy         : 'data.frame': 30 obs. of  2 variables:
## $ species    : 'data.frame': 27 obs. of  4 variables:
head(doubs$env)
```

```
##   dfs alt   slo flo pH har pho nit amm oxy bdo
## 1    3 934 6.176 84 79 45   1 20   0 122 27
## 2   22 932 3.434 100 80 40   2 20  10 103 19
## 3   102 914 3.638 180 83 52   5 22   5 105 35
## 4   185 854 3.497 253 80 72  10 21   0 110 13
## 5   215 849 3.178 264 81 84  38 52  20  80 62
## 6   324 846 3.497 286 79 60  20 15   0 102 53
```

**Question 1:** Describe some of the attributes of the `doubs` dataset.

- a. How many objects are in `doubs`?
- b. How many fish species are there in the `doubs` dataset?

- c. How many sites are in the `doubs` dataset?

**Answer 1a:** `doubs` data has 4 different objects which are `env`, `fish`, `xy`, and `species`. **Answer 1b:** It has 27 fish species. **Answer 1c:** It has total 30 sites.

## Visualizing the Doubs River Dataset

**Question 2:** Answer the following questions based on the spatial patterns of richness (i.e.,  $\alpha$ -diversity) and Brown Trout (*Salmo trutta*) abundance in the Doubs River.

- How does fish richness vary along the sampled reach of the Doubs River? `#reach = length(?)`
- How does Brown Trout (*Salmo trutta*) abundance vary along the sampled reach of the Doubs River?
- What do these patterns say about the limitations of using richness when examining patterns of biodiversity?

**Answer 2a:** It seems like downstream has more richness than upstream. **Answer 2b:** In contrast to richness, here the upstream has more abundance than downstream. **Answer 2c:** If we only use richness to evaluate the biodiversity, it might lead to the wrong conclusion. Since these two values lead to two different conclusion, we should check not only richness but also the abundance.

## 3) QUANTIFYING BETA-DIVERSITY

In the R code chunk below, do the following:

- write a function (`beta.w()`) to calculate Whittaker's  $\beta$ -diversity (i.e.,  $\beta_w$ ) that accepts a site-by-species matrix with optional arguments to specify pairwise turnover between two sites, and
- use this function to analyze various aspects of  $\beta$ -diversity in the Doubs River.

```
beta.w <- function(site.by.species = "", sitenum1 = "", sitenum2 = "",
  pairwise = FALSE){
  if (pairwise == TRUE){
    if (sitenum1 == "" | sitenum2 == ""){
      print("Error: please specify sites to compare")
      return(NA)}
    site1 = site.by.species[sitenum1, ]
    site2 = site.by.species[sitenum2, ]
    site1 = subset(site1, select = site1 > 0)
    site2 = subset(site2, select = site2 > 0)
    gamma = union(colnames(site1), colnames(site2))
    s      = length(gamma)
    a.bar = mean(c(specnumber(site1), specnumber(site2)))
    b.w   = round(s/a.bar - 1, 3)
    return(b.w)
  }
  else{
    SbyS.pa <- decostand(site.by.species, method = "pa")
    S <- ncol(SbyS.pa[,which(colSums(SbyS.pa) > 0)])
    a.bar <- mean(specnumber(SbyS.pa))
    b.w <- round(S/a.bar, 3)
    return(b.w)
  }
}

beta.w(doubs$fish, 1, 2, pairwise = TRUE)

## [1] 0.5
```

```
beta.w(doubs$fish, 1, 10, pairwise = TRUE)
```

```
## [1] 0.714
```

**Question 3:** Using your `beta.w()` function above, answer the following questions:

- Describe how local richness ( $\alpha$ ) and turnover ( $\beta$ ) contribute to regional ( $\gamma$ ) fish diversity in the Doubs.
- Is the fish assemblage at site 1 more similar to the one at site 2 or site 10?
- Using your understanding of the equation  $\beta_w = \gamma/\alpha$ , how would your interpretation of  $\beta$  change if we instead defined beta additively (i.e.,  $\beta = \gamma - \alpha$ )?

**Answer 3a:** alpha diversity represents the number of species in one site and beta diversity represents the differences among the sites. Thus, high alpha diversity and high beta diversity lead to high gamma diversity. However, there is a possibility that alpha diversity is high, but beta diversity is low, suggesting many species are overlapped. In this case, gamma diversity might not be as high as in the cases where the beta diversity is also high.

**Answer 3b:** beta.w from site 1 and 2 is 0.5 and from 1 and 10 is 0.714, suggesting that the site 1 and 10 is more dissimilar than site 1 and 2. **Answer 3c:** If we change beta diversity additively, we can interpret beta diversity as how many more species in the regional pool than in local sites, rather than how many times more diverse the region is than local sites.

## The Resemblance Matrix

In order to quantify  $\beta$ -diversity for more than two samples, we need to introduce a new primary ecological data structure: the **Resemblance Matrix**.

**Question 4:** How do incidence- and abundance-based metrics differ in their treatment of rare species?

**Answer 4:** In incidence-based metrics, rare species have equal weight as common species, so it can be useful if we need to detect and compare the presence or absence of rare species (Here, a site I found only 1 individual of rare species and another site that I found 30 of rare species considered same. However, abundance based matrices reflect the relative population size, so it can be less sensitive to compare rare species.

In the R code chunk below, do the following:

- make a new object, `fish`, containing the fish abundance data for the Doubs River,
- remove any sites where no fish were observed (i.e., rows with sum of zero),
- construct a resemblance matrix based on Sørensen's Similarity ("`fish.ds`"), and
- construct a resemblance matrix based on Bray-Curtis Distance ("`fish.db`").

```
fish <- doubs$fish
fish <- fish[rowSums(fish) > 0, ]
fish.ds <- vegdist(fish, method = "bray", binary = TRUE)
fish.db <- vegdist(fish, method = "bray", upper = TRUE, diag = TRUE)
```

**Question 5:** Using the distance matrices from above, answer the following questions:

- Does the resemblance matrix (`fish.db`) represent similarity or dissimilarity? What information in the resemblance matrix led you to arrive at your answer?
- Compare the resemblance matrices (`fish.db` or `fish.ds`) you just created. How does the choice of the Sørensen or Bray-Curtis distance influence your interpretation of site (dis)similarity?

**Answer 5a:** It shows "Bray-Curtis Dissimilarity," which indicates the percentage difference. Each value in the matrix shows how paired sites differ from each other. **Answer 5b:** I think they show quite similar trend. However, `fish.ds` doesn't reflect the relative abundance, so its dissimilarity value is little lower than `fish.db` which reflects relativeness.

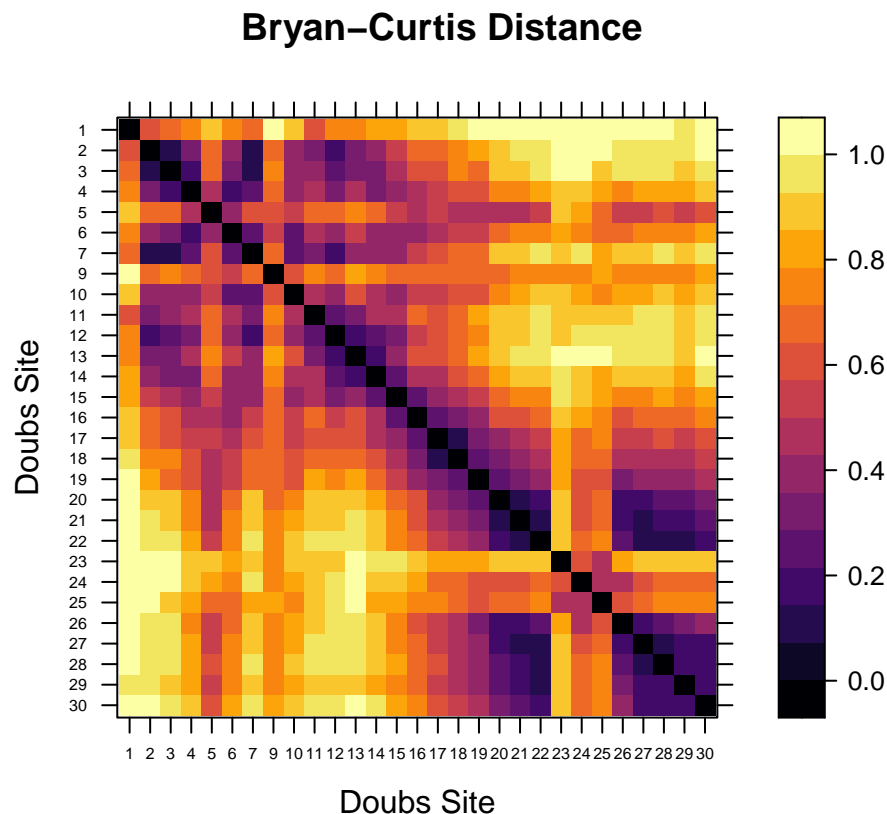
## 4) VISUALIZING BETA-DIVERSITY

### A. Heatmaps

In the R code chunk below, do the following:

1. define a color palette,
2. define the order of sites in the Doubs River, and
3. use the `levelplot()` function to create a heatmap of fish abundances in the Doubs River.

```
order <- rev(attr(fish.db, "Labels"))
levelplot(as.matrix(fish.db)[, order], aspect = "iso", col.regions = inferno,
          xlab = "Doubs Site", ylab = "Doubs Site", scales = list(cex = 0.5),
          main = "Bryan-Curtis Distance")
```



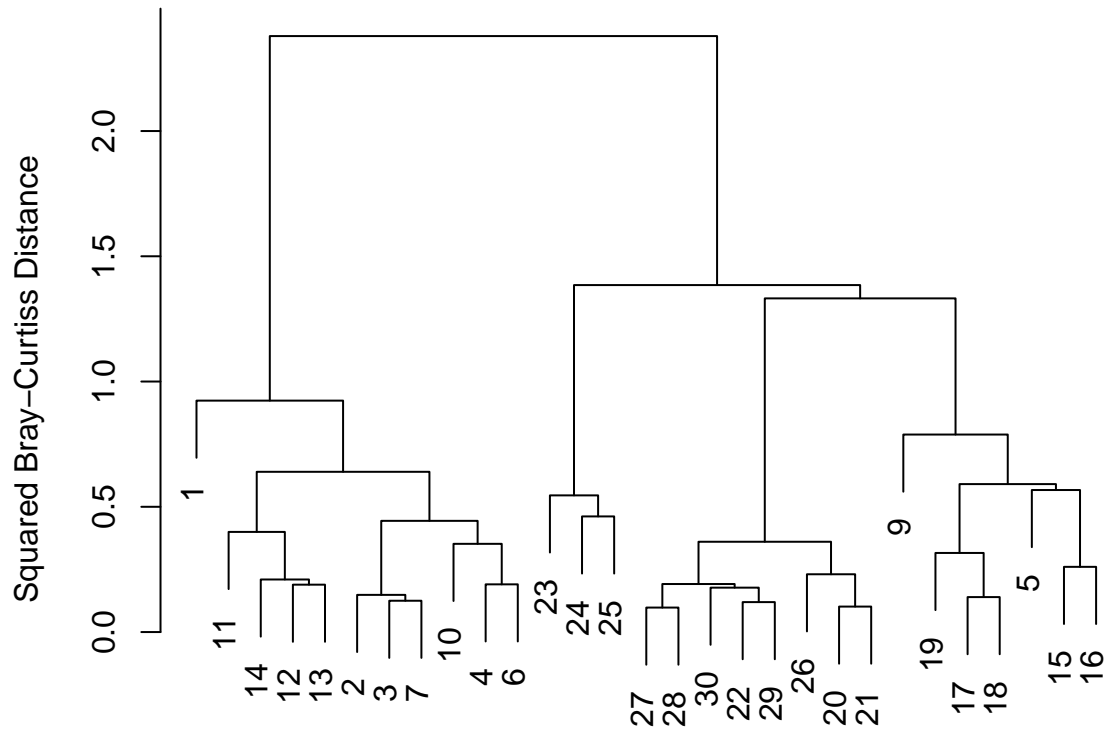
### B. Cluster Analysis

In the R code chunk below, do the following:

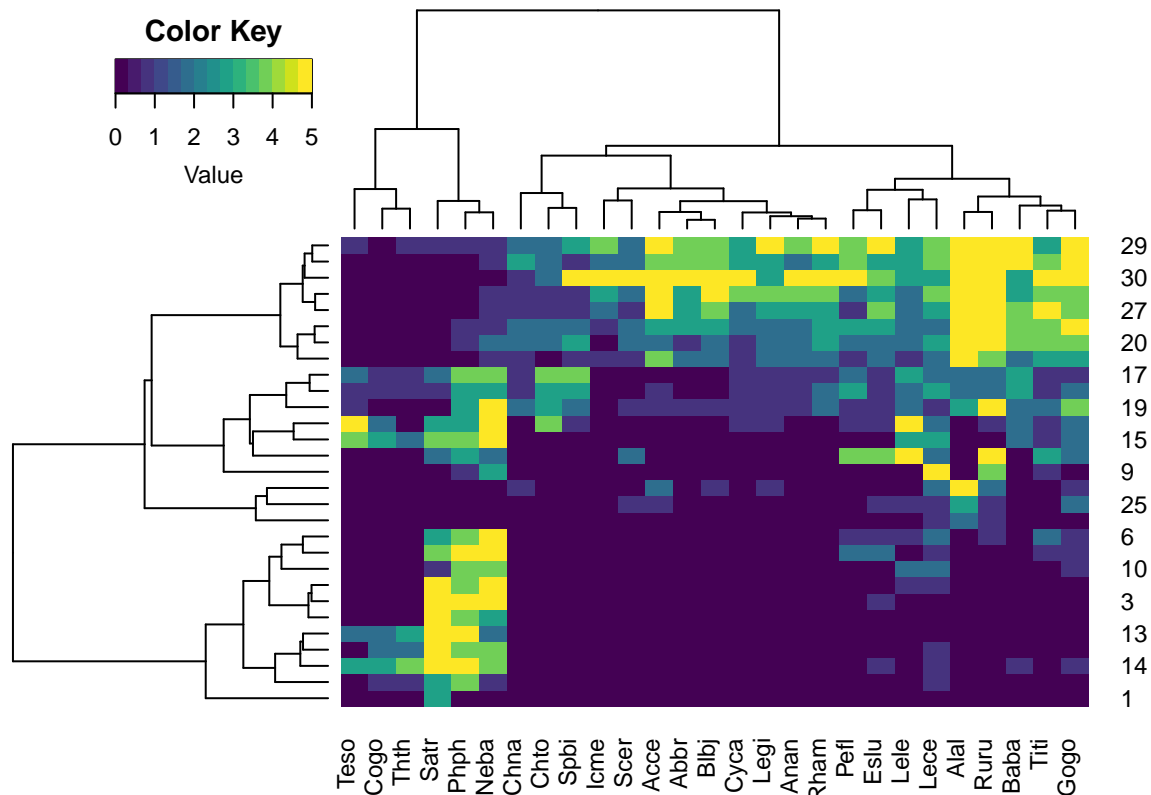
1. perform a cluster analysis using Ward's Clustering, and
2. plot your cluster analysis (use either `hclust` or `heatmap.2`).

```
fish.ward <- hclust(fish.db, method = "ward.D2")
par(mar = c(1, 5, 2, 2) + 0.1)
plot(fish.ward, main = "Doubs River Fish: Ward's Clustering",
     ylab = "Squared Bray-Curtiss Distance")
```

## Doubs River Fish: Ward's Clustering



```
gplots::heatmap.2(as.matrix(fish),
  distfun = function(x) vegdist(x, method = "bray"),
  hclustfun = function(x) hclust(x, method = "ward.D2"),
  col = viridis, trace = "none", density.info = "none")
```



**Question 6:** Based on cluster analyses and the introductory plots that we generated after loading the data, develop an ecological hypothesis for fish diversity the Doubs data set?

**Answer 6:** Wald's clustering shows the sites of Doubs' river are divided into two big clusters. Then, I can hypothesize the fish within each cluster are likely to exhibit strong interactions, such as being part of the same food chain or engaging in symbiotic relationships.

## C. Ordination

### Principal Coordinates Analysis (PCoA)

In the R code chunk below, do the following:

1. perform a Principal Coordinates Analysis to visualize beta-diversity
2. calculate the variation explained by the first three axes in your ordination
3. plot the PCoA ordination,
4. label the sites as points using the Doubs River site number, and
5. identify influential species and add species coordinates to PCoA plot.

```
add.spec.scores.class <-
function(ordi,comm,method="cor.scores",multi=1,Rscale=F,scaling="1") {
  ordiscores <- scores(ordi,display="sites")
  n <- ncol(comm)
  p <- ncol(ordiscores)
  specscores <- array(NA,dim=c(n,p))
  rownames(specscores) <- colnames(comm)
  colnames(specscores) <- colnames(ordiscores)
  if (method == "cor.scores") {
    for (i in 1:n) {
      for (j in 1:p) {specscores[i,j] <- cor(comm[,i],ordiscores[,j],method="pearson")}
    }
  }
}
```

```

}
if (method == "wa.scores") {specscores <- wascores(ordiscores,comm)}
if (method == "pcoa.scores") {
  rownames(ordiscores) <- rownames(comm)
  eigenv <- ordi$eig
  accounted <- sum(eigenv)
  tot <- 2*(accounted/ordi$GOF[2])-(accounted/ordi$GOF[1])
  eigen.var <- eigenv/(nrow(comm)-1)
  neg <- length(eigenv[eigenv<0])
  pos <- length(eigenv[eigenv>0])
  tot <- tot/(nrow(comm)-1)
  eigen.percen <- 100*eigen.var/tot
  eigen.cumpercen <- cumsum(eigen.percen)
  constant <- ((nrow(comm)-1)*tot)^0.25
  ordiscores <- ordiscores * (nrow(comm)-1)^-0.5 * tot^-0.5 * constant
  p1 <- min(p, pos)
  for (i in 1:n) {
    for (j in 1:p1) {
      specscores[i,j] <- cor(comm[,i],ordiscores[,j])*sd(comm[,i])/sd(ordiscores[,j])
      if(is.na(specscores[i,j])) {specscores[i,j]<-0}
    }
  }
  if (Rscale==T && scaling=="2") {
    percen <- eigen.var/tot
    percen <- percen^0.5
    ordiscores <- sweep(ordiscores,2,percen,"/")
    specscores <- sweep(specscores,2,percen,"*")
  }
  if (Rscale==F) {
    specscores <- specscores / constant
    ordiscores <- ordi$points
  }
  ordi$points <- ordiscores
  ordi$eig <- eigen.var
  ordi$eig.percen <- eigen.percen
  ordi$eig.cumpercen <- eigen.cumpercen
  ordi$eigen.total <- tot
  ordi$R.constant <- constant
  ordi$Rscale <- Rscale
  ordi$scaling <- scaling
}
specscores <- specscores * multi
ordi$cproj <- specscores
return(ordi)
}

```

```

fish.pcoa <- cmdscale(fish.db, eig = TRUE, k = 3)
explainvar1 <- round(fish.pcoa$eig[1] / sum(fish.pcoa$eig), 3) * 100
explainvar2 <- round(fish.pcoa$eig[2] / sum(fish.pcoa$eig), 3) * 100
explainvar3 <- round(fish.pcoa$eig[3] / sum(fish.pcoa$eig), 3) * 100

```

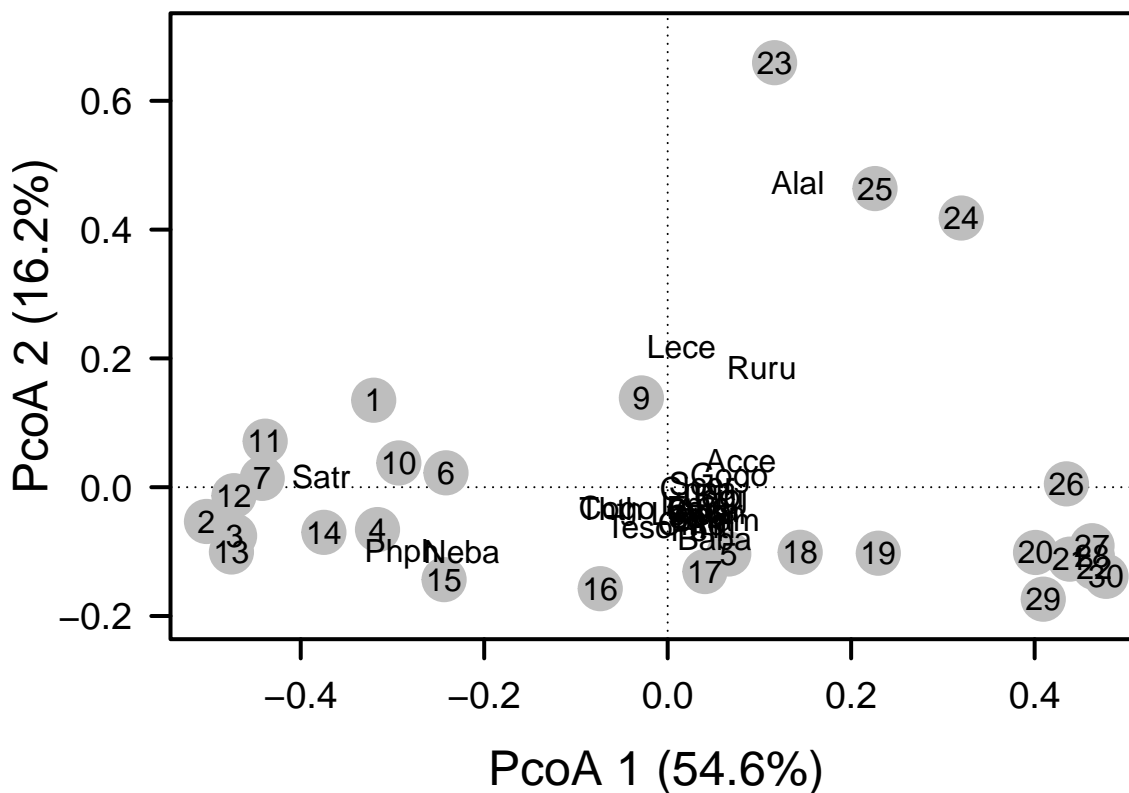


```

par(mar = c(5, 5, 1, 2) + 0.1)
plot(fish.pcoa$points[,1], fish.pcoa$points[,2], ylim = c(-0.2, 0.7),
     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(fish.pcoa$points[,1], fish.pcoa$points[,2],
       pch = 19, cex = 3, bg = "gray", col = "gray")
text(fish.pcoa$points[,1], fish.pcoa$points[,2],
     labels = row.names(fish.pcoa$points))

fishREL <- fish
for(i in 1:nrow(fish)){
  fishREL[i, ] = fish[i, ] / sum(fish[i, ])
}
fish.pcoa <- add.spec.scores.class(fish.pcoa, fishREL, method = "pcoa.scores")
text(fish.pcoa$cproj[,1], fish.pcoa$cproj[,2],
     labels = row.names(fish.pcoa$cproj), col = "black")

```



In the R code chunk below, do the following:

1. identify influential species based on correlations along each PCoA axis (use a cutoff of 0.70), and
2. use a permutation test (999 permutations) to test the correlations of each species along each axis.

```
spe.corr <- add.spec.scores.class(fish.pcoa, fishREL, method = "cor.scores")$cproj
corrcut <- 0.7
imp.spp <- spe.corr[abs(spe.corr[, 1]) >= corrcut | abs(spe.corr[, 2]) >= corrcut, ]
print(imp.spp)
```

```
##           Dim1      Dim2      Dim3
## Phph -0.8674640 -0.1699316 -0.12463098
## Neba -0.7674114 -0.1855678 -0.36963830
## Rham  0.8088751 -0.4192567  0.14136301
## Legi  0.8201759 -0.1701803  0.12423941
## Cyca  0.7595122 -0.4442926  0.17313658
## Abbr  0.7704744 -0.3452714  0.29277803
## Acce  0.7635195  0.2155765  0.10288179
## Blbj  0.8118483 -0.1324698  0.25581178
## Alal  0.4471283  0.8119843 -0.05167131
## Anan  0.7974122 -0.3918972  0.20944968
```

```
fit <- envfit(fish.pcoa, fishREL, perm = 999)
print(fit)
```

```
##
## ***VECTORS
##
##           Dim1      Dim2      r2 Pr(>r)
## Cogo -0.83884 -0.54438 0.2982 0.012 *
## Satr -0.99904  0.04371 0.4326 0.003 **
## Phph -0.94110 -0.33813 0.7814 0.001 ***
## Neba -0.91413 -0.40543 0.6234 0.001 ***
## Thth -0.87692 -0.48063 0.2634 0.019 *
## Teso -0.44704 -0.89452 0.1700 0.086 .
## Chna  0.99707 -0.07644 0.4612 0.001 ***
## Chto  0.42032 -0.90738 0.2579 0.018 *
## Lele  0.33041 -0.94384 0.0495 0.550
## Lece  0.06856  0.99765 0.3399 0.009 **
## Baba  0.54118 -0.84091 0.6752 0.001 ***
## Spbi  0.57341 -0.81927 0.4138 0.003 **
## Gogo  0.97507  0.22188 0.3753 0.003 **
## Eslu  0.72044 -0.69352 0.1673 0.096 .
## Pefl  0.43762 -0.89916 0.3048 0.014 *
## Rham  0.72476 -0.68901 0.8301 0.001 ***
## Legi  0.93461 -0.35568 0.7016 0.001 ***
## Scer  0.98569  0.16858 0.3533 0.008 **
## Cyca  0.68181 -0.73153 0.7743 0.001 ***
## Titi  0.64378 -0.76521 0.4586 0.001 ***
## Abbr  0.77254 -0.63497 0.7128 0.001 ***
## Icme  0.75626 -0.65427 0.5270 0.001 ***
## Acce  0.88799  0.45986 0.6294 0.001 ***
## Ruru  0.48379  0.87518 0.5177 0.001 ***
## Blbj  0.95802 -0.28671 0.6766 0.001 ***
## Alal  0.28755  0.95777 0.8592 0.001 ***
## Anan  0.74277 -0.66954 0.7894 0.001 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## Permutation: free
```

```
## Number of permutations: 999
```

**Question 7:** Address the following questions about the ordination results of the `doubs` data set:

- Describe the grouping of sites in the Doubs River based on fish community composition.
- Generate a hypothesis about which fish species are potential indicators of river quality.

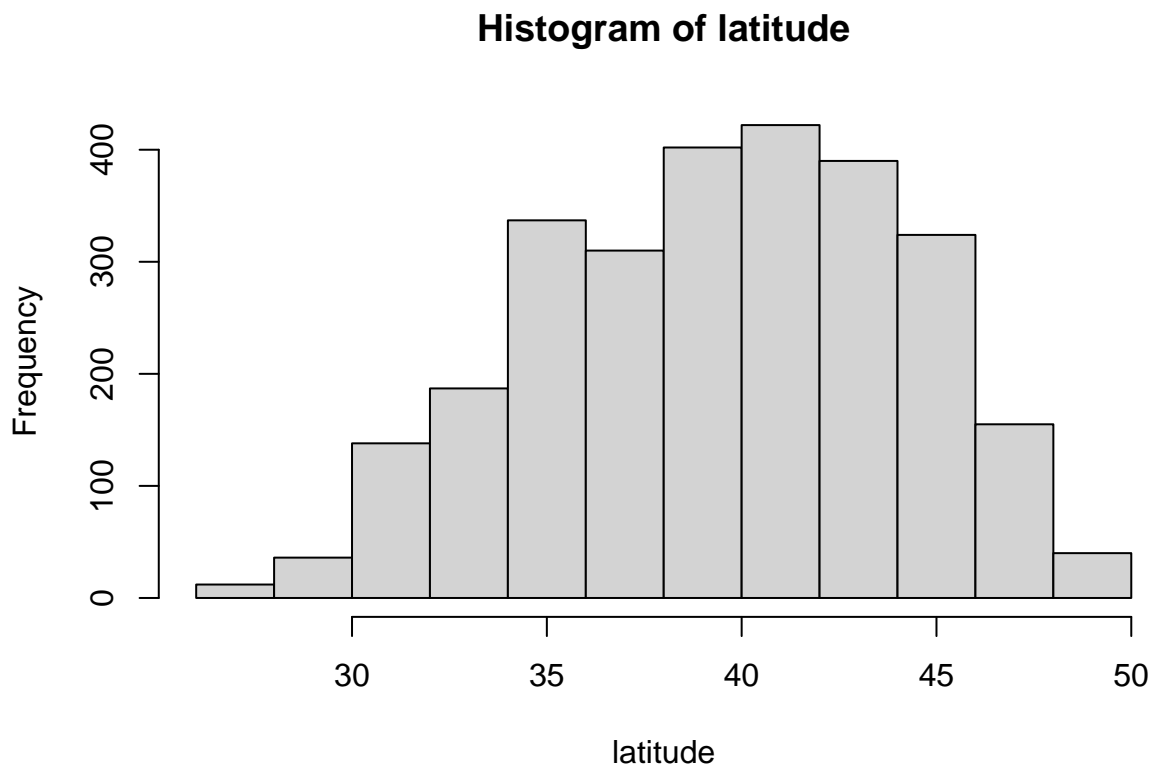
**Answer 7a:** X-axis explains 54.6% of the community composition in the Doubs river, and y-axis explains 16.4%. We can see most of the sites are located around 0 of y axis, but widely spreaded according to x-axis. Some of sites have a high y values, indicating thses sites' comosition is very different from others. **Answer 7b:** According to these information, we can know 10 significant species which is highly contributed to diversity, So, I hypothesized these 10 species are environmentally senetive and these species can represent the quality of each site.

## SYNTHESIS

Load the dataset from that you and your partner are using for the team project. Use one of the tools introduced in the beta diversity module to visualize your data. Describe any interesting patterns and identify a hypothesis is relevant to the principles of biodiversity.

**Answer** The graph groups 2,700+ sites into five latitude. Site 1 (25-30°), Site 2 (30-35°), Site 3 (35-40°), Site 4 (40-45°), and Site 5 (45-50°) represent different latitude ranges. Fish communities are most similar between latitudes 35-45° but differ significantly in the north and south. Thus, I hypothesize that biodiversity is lower in northern and southern latitudes due to harsh climate conditions.

```
fish.team <- read.csv("/cloud/project/QB2025_Choi/Fish_Dataset.csv")
latitude <- fish.team$Latitude
hist(latitude)
```



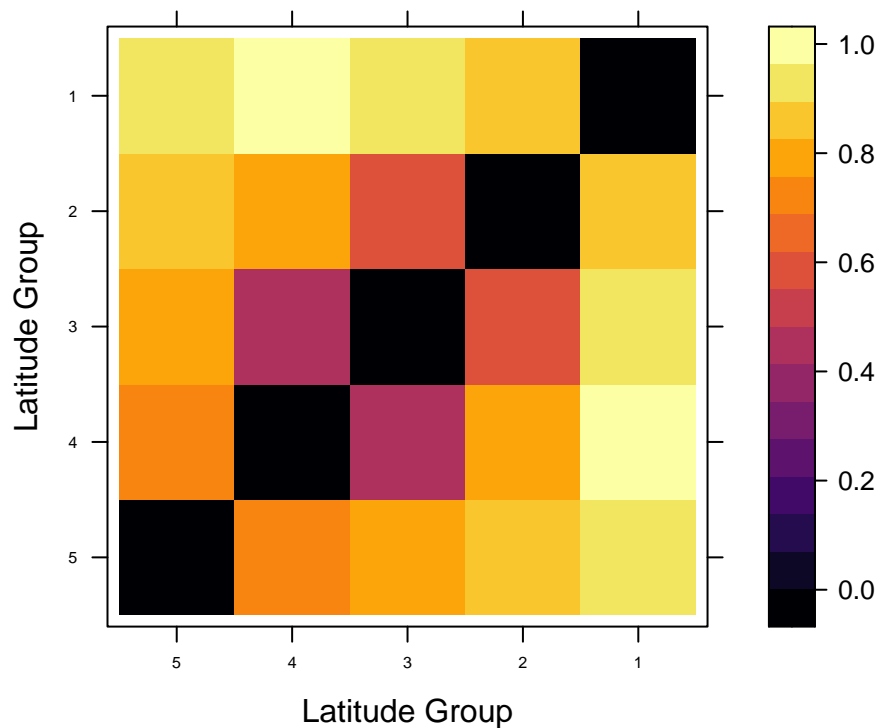
```
fish.team$Latitude <- cut(fish.team$Latitude, breaks=c(20, 25, 30, 35, 40, 45, 50),
  labels=c("20-25", "25-30", "30-35", "35-40", "40-45", "45-50"),
  include.lowest=TRUE)
```

```

fish_latitude <- cbind(fish.team[,2], fish.team[,23:658])
fish_latitude <- data.frame(Latitude = fish.team$Latitude, fish_latitude)
result <- aggregate(. ~ Latitude, data=fish_latitude, sum, na.rm=TRUE)
result <- cbind(result[,1], result[,3:638])
colnames(result)[1] <- "Latitude_Group"
new_result <- result[,2:637]
result.db <- vegdist(new_result, method = "bray", upper = TRUE, diag = TRUE)
result.db.mat <- as.matrix(result.db)
order1 <- rev(rownames(result.db.mat))
levelplot(result.db.mat[order1, order1],
           aspect = "iso",
           col.regions = viridis::inferno(256),
           xlab = "Latitude_Group",
           ylab = "Latitude_Group",
           scales = list(cex = 0.5),
           main = "Bray-Curtis Distance")

```

## Bray-Curtis Distance



```

result.ward <- hclust(result.db, method = "ward.D2")
par(mar = c(1, 5, 2, 2) + 0.1)
plot(result.ward, main = "Fish by Latitude: Ward's Clustering",
     ylab = "Squared Bray-Curtis Distance")

```

