

11. Worksheet: Phylogenetic Diversity - Traits

Student Name; Z620: Quantitative Biodiversity, Indiana University

30 April, 2021

OVERVIEW

Up to this point, we have been focusing on patterns taxonomic diversity in Quantitative Biodiversity. Although taxonomic diversity is an important dimension of biodiversity, it is often necessary to consider the evolutionary history or relatedness of species. The goal of this exercise is to introduce basic concepts of phylogenetic diversity.

After completing this exercise you will be able to:

1. create phylogenetic trees to view evolutionary relationships from sequence data
2. map functional traits onto phylogenetic trees to visualize the distribution of traits with respect to evolutionary history
3. test for phylogenetic signal within trait distributions and trait-based 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 (**11.PhyloTraits_Worksheet.Rmd**) with all code blocks filled out and questions answered) and the PDF output of Knitr (**11.PhyloTraits_Worksheet.pdf**).

The completed exercise is due on **Wednesday, April 28th, 2021 before 12:00 PM (noon)**.

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 “/11.PhyloTraits” folder, and
4. load all of the required R packages (be sure to install if needed).

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

```
## [1] "/home/patgwall/Classwork/Current/QB/2.Worksheets/11.PhyloTraits"
```

```
setwd("~/Classwork/Current/QB/2.Worksheets/11.PhyloTraits")
```

2) DESCRIPTION OF DATA

The maintenance of biodiversity is thought to be influenced by **trade-offs** among species in certain functional traits. One such trade-off involves the ability of a highly specialized species to perform exceptionally well on a particular resource compared to the performance of a generalist. In this exercise, we will take a phylogenetic approach to mapping phosphorus resource use onto a phylogenetic tree while testing for specialist-generalist trade-offs.

3) SEQUENCE ALIGNMENT

Question 1: Using your favorite text editor, compare the `p.isolates.fasta` file and the `p.isolates.afa` file. Describe the differences that you observe between the two files.

Answer 1: The `.afa` file is aligned meaning that homologous sequence regions have been placed into the same position along the sequence through the addition of gaps.

In the R code chunk below, do the following: 1. read your alignment file, 2. convert the alignment to a DNAbin object, 3. select a region of the gene to visualize (try various regions), and 4. plot the alignment using a grid to visualize rows of sequences.

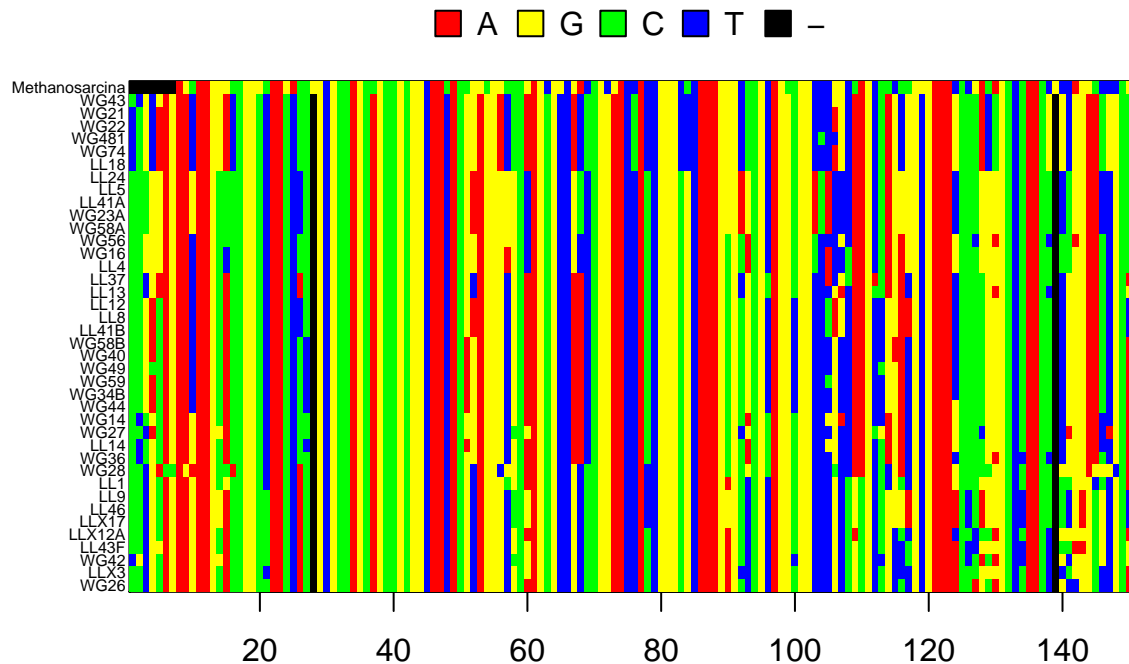
```
package.list <- c('ape', 'seqinr', 'phylobase', 'geiger', 'picante', 'stats', 'RColorBrewer', 'caper',
for (package in package.list) {
  if (!require(package, character.only=TRUE, quietly=TRUE)) {
    install.packages(package)
    library(package, character.only=TRUE)
  }
}
```

```
##
## Attaching package: 'seqinr'
##
## The following objects are masked from 'package:ape':
##
##   as.alignment, consensus
##
## Attaching package: 'phylobase'
##
## The following object is masked from 'package:ape':
##
##   edges
##
## Attaching package: 'permute'
##
## The following object is masked from 'package:seqinr':
##
##   getType
## This is vegan 2.5-7
```

```

##
## Attaching package: 'nlme'
## The following object is masked from 'package:seqinr':
##
##     gls
##
## Attaching package: 'dplyr'
## The following object is masked from 'package:MASS':
##
##     select
## The following object is masked from 'package:nlme':
##
##     collapse
## The following object is masked from 'package:seqinr':
##
##     count
## The following objects are masked from 'package:stats':
##
##     filter, lag
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
##
## Attaching package: 'phangorn'
## The following objects are masked from 'package:vegan':
##
##     diversity, treedist
# read alignment and convert/viz
read_aln <- read.alignment(file = "./data/p.isolates.afa", format = 'fasta')
p.DNAbin <- as.DNAbin(read_aln)
window <- p.DNAbin[,500:650]
image.DNAbin(window, cex.lab = 0.50)

```



Question 2: Make some observations about the muscle alignment of the 16S rRNA gene sequences for our bacterial isolates and the outgroup, *Methanosarcina*, a member of the domain Archaea. Move along the alignment by changing the values in the `window` object.

- Approximately how long are our sequence reads?
- What regions do you think would be appropriate for phylogenetic inference and why?

Answer 2a: the majority of the sequence reads are about 700-900 base pairs. The outgroup is much longer and some are much shorter. **Answer 2b:** I would think we would want to use regions without large sections missing for any taxa so maybe 500-650 would be a good choice although it is short.

4) MAKING A PHYLOGENETIC TREE

Once you have aligned your sequences, the next step is to construct a phylogenetic tree. Not only is a phylogenetic tree effective for visualizing the evolutionary relationship among taxa, but as you will see later, the information that goes into a phylogenetic tree is needed for downstream analysis.

A. Neighbor Joining Trees

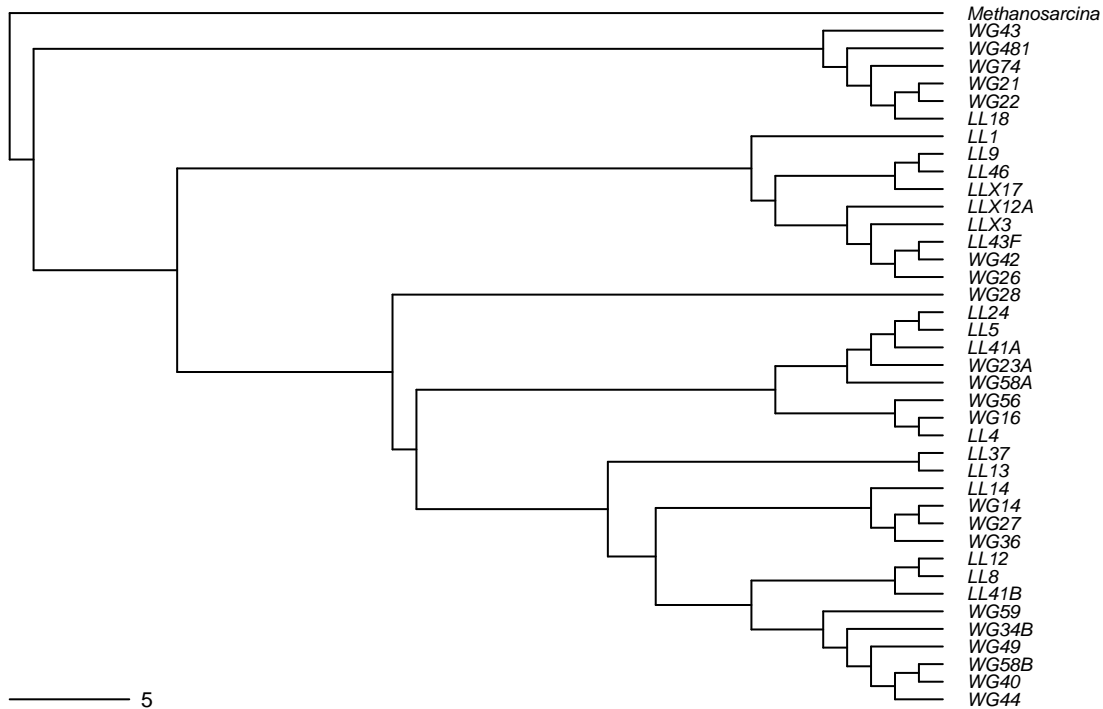
In the R code chunk below, do the following:

- calculate the distance matrix using `model = "raw"`,
- create a Neighbor Joining tree based on these distances,
- define “*Methanosarcina*” as the outgroup and root the tree, and
- plot the rooted tree.

```
seq_dist_raw <- dist.dna(p.DNAbin, model = "raw", pairwise.deletion = FALSE)
nj_tree <- bionj(seq_dist_raw)
outgroup <- match("Methanosarcina", nj_tree$tip.label)
nj_rooted <- root(nj_tree, outgroup, resolve.root = TRUE)

par(mar = c(1,1,2,1) + 0.1)
plot.phylo(nj_rooted, main = "Neighbor Joining Tree", "phylogram", use.edge.length = FALSE, direction =
add.scale.bar(cex = 0.7)
```

Neighbor Joining Tree



Question 3: What are the advantages and disadvantages of making a neighbor joining tree?

Answer 3: Its easy and fast! But also its not the most accurate method.

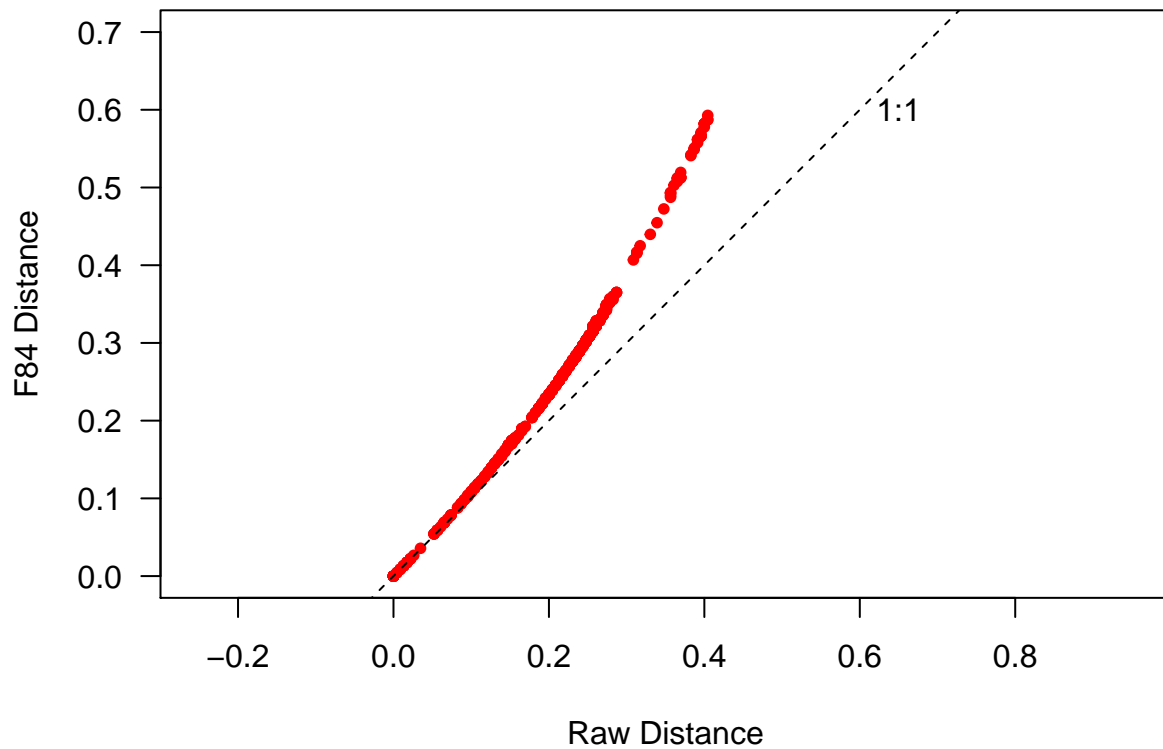
B) SUBSTITUTION MODELS OF DNA EVOLUTION

In the R code chunk below, do the following:

1. make a second distance matrix based on the Felsenstein 84 substitution model,
2. create a saturation plot to compare the *raw* and *Felsenstein (F84)* substitution models,
3. make Neighbor Joining trees for both, and
4. create a cophylogenetic plot to compare the topologies of the trees.

```
# fancy substitution model
seq_dist_F84 <- dist.dna(p.DNABin, model = "F84", pairwise.deletion = FALSE)

# Plot against raw distance
par(mar = c(5,5,2,1) + 0.1)
plot(seq_dist_raw, seq_dist_F84, pch = 20, col = "red", las = 1, asp = 1, xlim = c(0, 0.7), ylim = c(0,
  xlab = "Raw Distance", ylab = "F84 Distance"))
abline(b = 1, a = 0, lty = 2)
text(0.65, 0.6, "1:1")
```

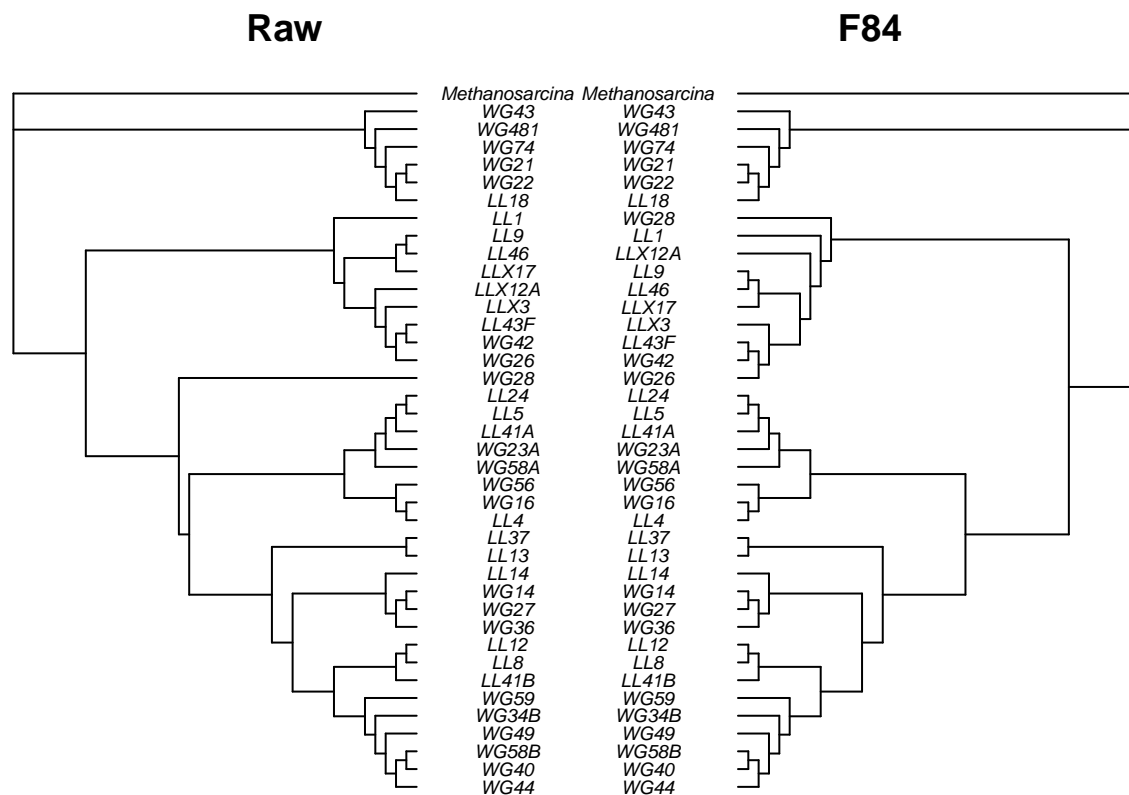


```
# now we'll make the new trees
raw_tree <- bionj(seq_dist_raw)
F84_tree <- bionj(seq_dist_F84)

# Define outgroup
raw_outgroup <- match("Methanosarcina", raw_tree$tip.label)
F84_outgroup <- match("Methanosarcina", F84_tree$tip.label)

# root the trees
raw_rooted <- root(raw_tree, raw_outgroup, resolve_root = TRUE)
F84_rooted <- root(F84_tree, F84_outgroup, resolve_root = TRUE)

# make the two tree plot
layout(matrix(c(1,2), 1, 2), width = c(1, 1))
par(mar = c(1,1,2,0))
plot.phylo(raw_rooted, type = "phylogram", direction = "right", show.tip.label = TRUE,
           use.edge.length = FALSE, adj = 0.5, cex = 0.6, label.offset = 2, main = "Raw")
par(mar = c(1,0,2,1))
plot.phylo(F84_rooted, type = "phylogram", direction = "left", show.tip.label = TRUE,
           use.edge.length = FALSE, adj = 0.5, cex = 0.6, label.offset = 2, main = "F84")
```

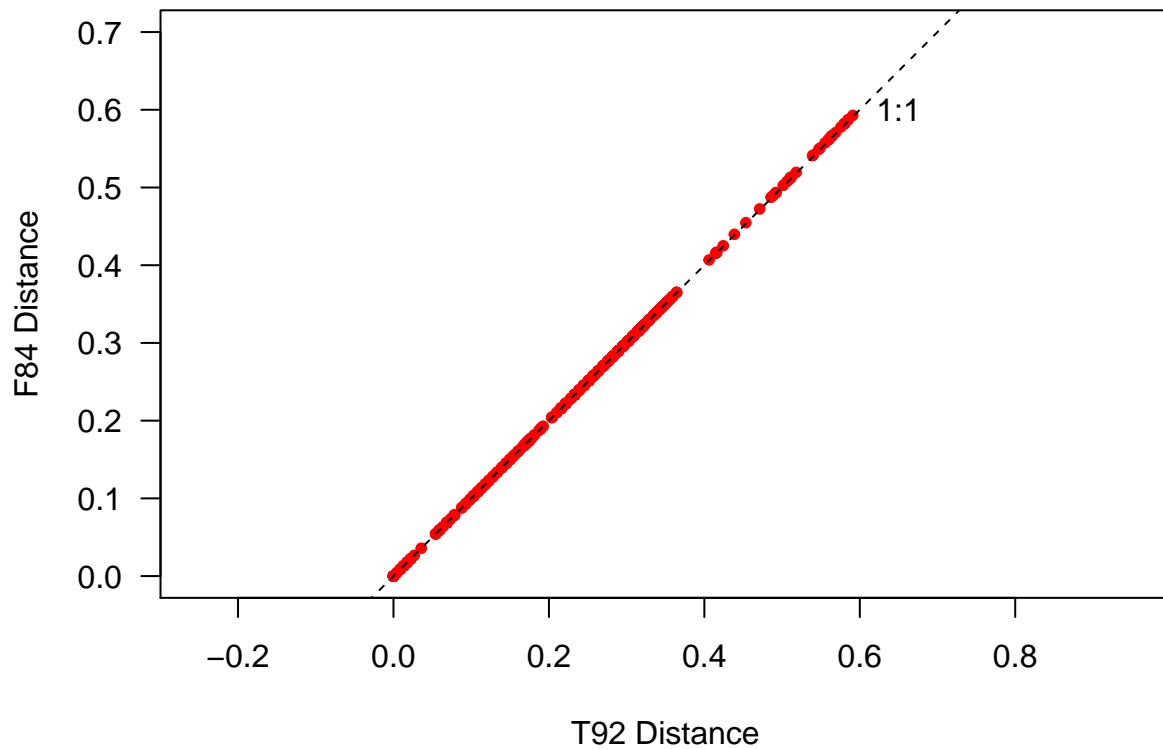


In the R code chunk below, do the following:

1. pick another substitution model,
2. create a distance matrix and tree for this model,
3. make a saturation plot that compares that model to the *Felsenstein (F84)* model,
4. make a cophylogenetic plot that compares the topologies of both models, and
5. be sure to format, add appropriate labels, and customize each plot.

```
# fancy substitution model
seq_dist_F84 <- dist.dna(p.DNABin, model = "F84", pairwise.deletion = FALSE)
seq_dist_t92 <- dist.dna(p.DNABin, model = "T92", pairwise.deletion = FALSE)

# Plot against raw distance
par(mar = c(5,5,2,1) + 0.1)
plot(seq_dist_t92, seq_dist_F84, pch = 20, col = "red", las = 1, asp = 1, xlim = c(0, 0.7), ylim = c(0,
      xlab = "T92 Distance", ylab = "F84 Distance"))
abline(b = 1, a = 0, lty = 2)
text(0.65, 0.6, "1:1")
```

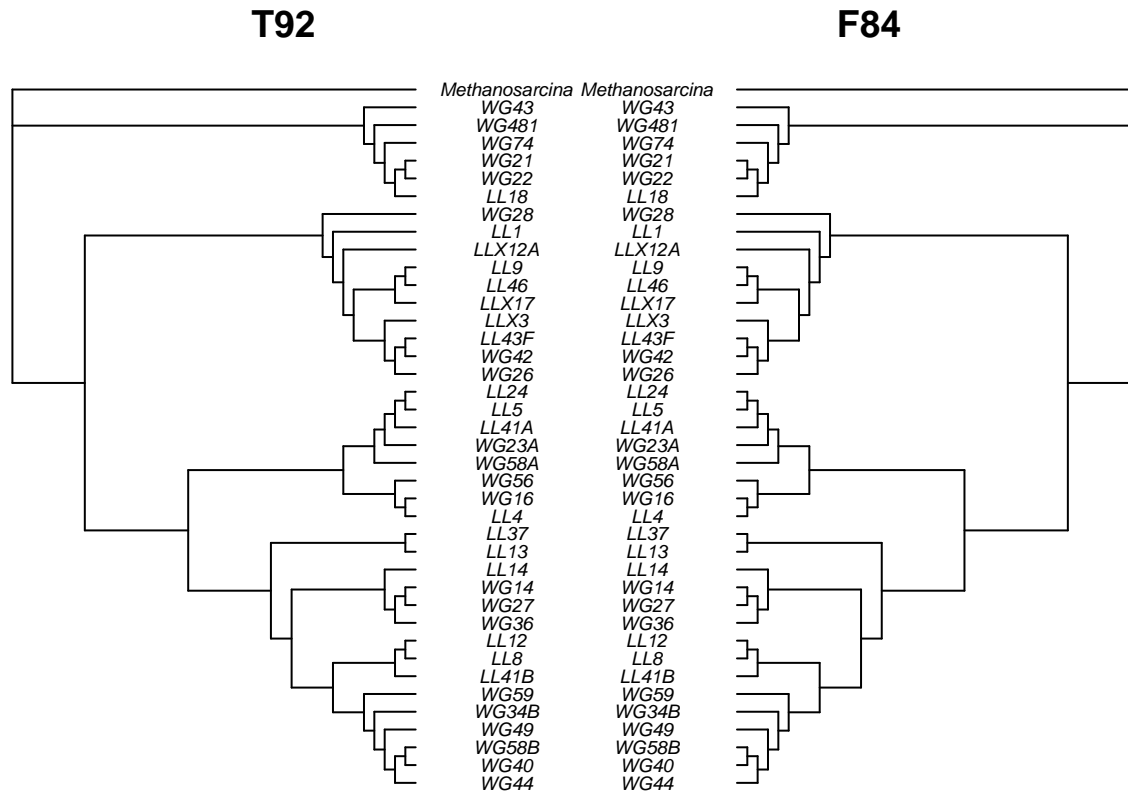


```
# now we'll make the new trees
t92_tree <- bionj(seq_dist_t92)
F84_tree <- bionj(seq_dist_F84)

# Define outgroup
t92_outgroup <- match("Methanosarcina", t92_tree$tip.label)
F84_outgroup <- match("Methanosarcina", F84_tree$tip.label)

# root the trees
t92_rooted <- root(t92_tree, t92_outgroup, resolve_root = TRUE)
F84_rooted <- root(F84_tree, F84_outgroup, resolve_root = TRUE)

# make the two tree plot
layout(matrix(c(1,2), 1, 2), width = c(1, 1))
par(mar = c(1,1,2,0))
plot.phylo(t92_rooted, type = "phylogram", direction = "right", show.tip.label = TRUE,
           use.edge.length = FALSE, adj = 0.5, cex = 0.6, label.offset = 2, main = "T92")
par(mar = c(1,0,2,1))
plot.phylo(F84_rooted, type = "phylogram", direction = "left", show.tip.label = TRUE,
           use.edge.length = FALSE, adj = 0.5, cex = 0.6, label.offset = 2, main = "F84")
```

Question 4:

- Describe the substitution model that you chose. What assumptions does it make and how does it compare to the F84 model?
- Using the saturation plot and cophylogenetic plots from above, describe how your choice of substitution model affects your phylogenetic reconstruction. If the plots are inconsistent with one another, explain why.
- How does your model compare to the *F84* model and what does this tell you about the substitution rates of nucleotide transitions?

Answer 4a:

Answer 4b: Saturation plot indicates that the distances are identical between the two methods. This is also reflected in the tree. **Answer 4c:** We can infer that GC content is not important for substitution rates of these sequences.

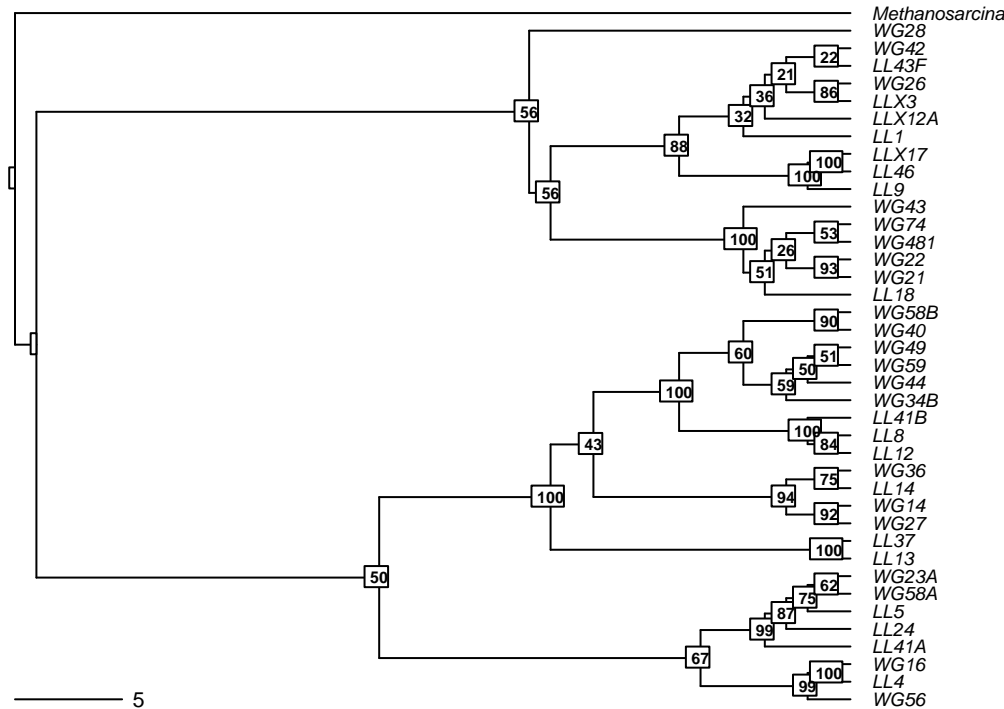
C) ANALYZING A MAXIMUM LIKELIHOOD TREE

In the R code chunk below, do the following:

- Read in the maximum likelihood phylogenetic tree used in the handout.
- Plot bootstrap support values onto the tree

```
ml_bootstrap <- read.tree("./data/ml_tree/RAxML_bipartitions.T1")
par(mar = c(1, 1, 2, 1) + 0.1)
plot.phylo(ml_bootstrap, type = "phylogram", direction = "right",
  show.tip.label = TRUE, use.edge.length = FALSE, cex = 0.6, label.offset = 1,
  main = "Maximum Likelihood with Support Values")
add.scale.bar(cex = 0.7)
nodelabels(ml_bootstrap$node.label, font = 2, bg = "white", frame = "r", cex = 0.5)
```

Maximum Likelihood with Support Values



Question 5:

- How does the maximum likelihood tree compare the to the neighbor-joining tree in the handout? If the plots seem to be inconsistent with one another, explain what gives rise to the differences.
- Why do we bootstrap our tree?
- What do the bootstrap values tell you?
- Which branches have very low support?
- Should we trust these branches?

Answer 5a: The ML tree is pretty different. The reason **Answer 5b:** Bootstrapping lets us look beyond our specific data to consider what slightly different data might look like this. This let's us make better, less overfit, predictions in general. Here the prediction is the structure of the tree. **Answer 5c:** The bootstrap values give us a confidence or significance score which lets us judge the aspects of the tree that we can safely use for futher inferences. **Answer 5d:** Many branches have low support, the lowest are all clustered together in the group near the top on the visualization above. **Answer 5e:** We probably should not trust these branches as they only appear in the minority of bootstrapped trees.

5) INTEGRATING TRAITS AND PHYLOGENY

A. Loading Trait Database

In the R code chunk below, do the following:

- import the raw phosphorus growth data, and
- standardize the data for each strain by the sum of growth rates.

```
p_growth <- read.table("./data/p.isolates.raw.growth.txt", sep = "\t", header = TRUE, row.names = 1)
p_growth_std <- p_growth / (apply(p_growth, 1, sum))
```

B. Trait Manipulations

In the R code chunk below, do the following:

1. calculate the maximum growth rate (μ_{max}) of each isolate across all phosphorus types,
2. create a function that calculates niche breadth (nb), and
3. use this function to calculate nb for each isolate.

```
umax <- apply(p_growth, 1, max)
levins <- function(p_xi) {
  p = 0
  for (i in p_xi) {
    p = p + i^2
  }
  nb = 1 / (length(p_xi) * p)
  return(nb)
}

# actually get the niche breadths
nb <- as.matrix(levins(p_growth_std))
rownames(nb) <- row.names(p_growth)
colnames(nb) <- c("NB")
```

C. Visualizing Traits on Trees

In the R code chunk below, do the following:

1. pick your favorite substitution model and make a Neighbor Joining tree,
2. define your outgroup and root the tree, and
3. remove the outgroup branch.

I have to give up on this part. I cannot get adephylo installed

```
# get the tree
nj_tree <- bionj(seq_dist_F84)
outgroup <- match("Methanosarcina", nj_tree$tip.label)
nj_rooted <- root(nj_tree, outgroup, resolve.root = TRUE)
nj_rooted <- drop.tip(nj_rooted, "Methanosarcina")
```

In the R code chunk below, do the following:

1. define a color palette (use something other than “YlOrRd”),
2. map the phosphorus traits onto your phylogeny,
3. map the nb trait on to your phylogeny, and
4. customize the plots as desired (use `help(table.phylo4d)` to learn about the options).

Question 6:

- a) Make a hypothesis that would support a generalist-specialist trade-off.
- b) What kind of patterns would you expect to see from growth rate and niche breadth values that would support this hypothesis?

Answer 6a: A generalist-specialist trade-off would likely mean that generalists have moderate growth rates on many phosphorous substrates. **Answer 6b:** Specifically this would look like a negative correlation between maximum growth rate and niche breadth.

6) HYPOTHESIS TESTING

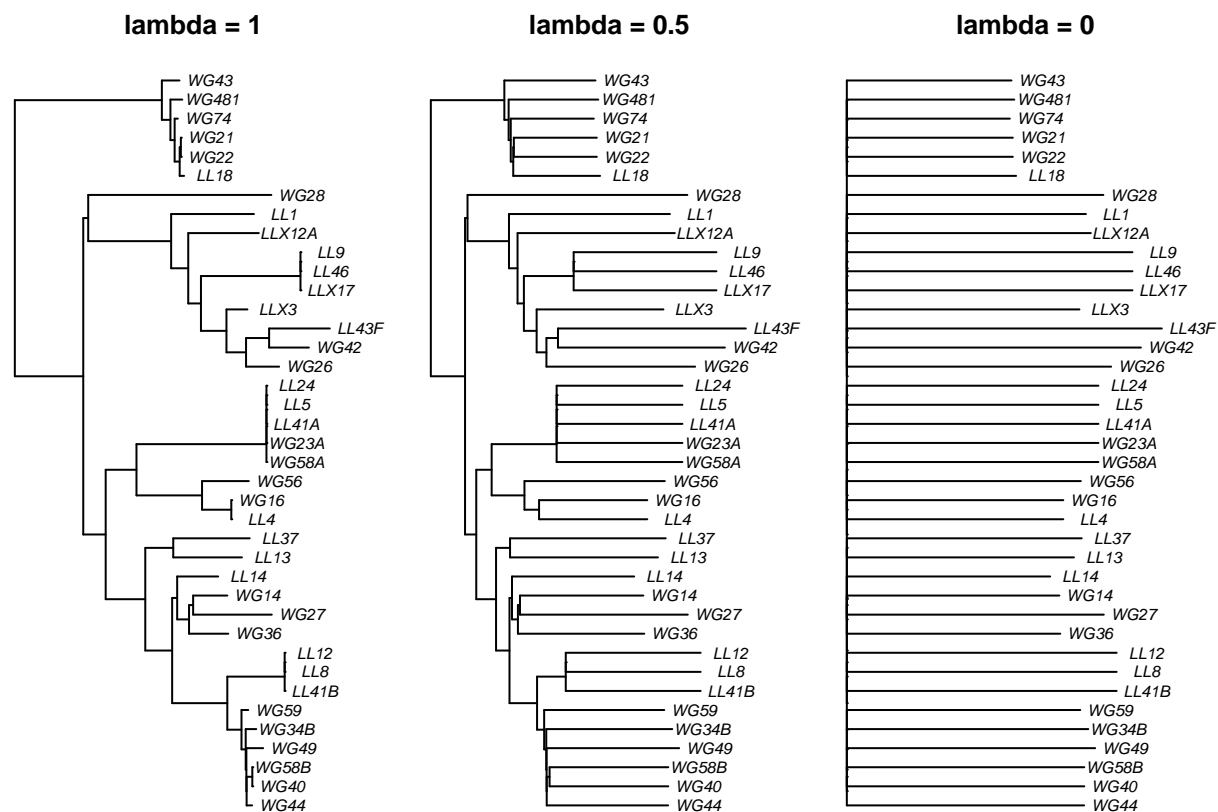
A) Phylogenetic Signal: Pagel's Lambda

In the R code chunk below, do the following:

1. create two rescaled phylogenetic trees using lambda values of 0.5 and 0,
2. plot your original tree and the two scaled trees, and
3. label and customize the trees as desired.

```
nj_lambda5 <- rescale(nj_rooted, "lambda", 0.5)
nj_lambda0 <- rescale(nj_rooted, "lambda", 0)

layout(matrix(c(1,2,3), 1, 3), width = c(1,1,1))
par(mar = c(1, 0.5, 2, 0.5) + 0.1)
plot(nj_rooted, main = "lambda = 1", cex = 0.7, adj = 0.5)
plot(nj_lambda5, main = "lambda = 0.5", cex = 0.7, adj = 0.5)
plot(nj_lambda0, main = "lambda = 0", cex = 0.7, adj = 0.5)
```



In the R code chunk below, do the following:

1. use the `fitContinuous()` function to compare your original tree to the transformed trees.

```
fitContinuous(nj_rooted, nb, model = "lambda")

## GEIGER-fitted comparative model of continuous data
## fitted 'lambda' model parameters:
## lambda = 0.020848
## sigsq = 0.106492
## z0 = 0.661368
##
## model summary:
## log-likelihood = 21.661104
```

```
## AIC = -37.322208
## AICc = -36.636494
## free parameters = 3
##
## Convergence diagnostics:
## optimization iterations = 100
## failed iterations = 49
## number of iterations with same best fit = NA
## frequency of best fit = NA
##
## object summary:
## 'lik' -- likelihood function
## 'bnd' -- bounds for likelihood search
## 'res' -- optimization iteration summary
## 'opt' -- maximum likelihood parameter estimates

fitContinuous(nj_lambda0, nb, model = "lambda")

## GEIGER-fitted comparative model of continuous data
## fitted 'lambda' model parameters:
## lambda = 0.000000
## sigsq = 0.106395
## z0 = 0.657777
##
## model summary:
## log-likelihood = 21.652293
## AIC = -37.304587
## AICc = -36.618872
## free parameters = 3
##
## Convergence diagnostics:
## optimization iterations = 100
## failed iterations = 0
## number of iterations with same best fit = 86
## frequency of best fit = 0.86
##
## object summary:
## 'lik' -- likelihood function
## 'bnd' -- bounds for likelihood search
## 'res' -- optimization iteration summary
## 'opt' -- maximum likelihood parameter estimates
```

Question 7: There are two important outputs from the `fitContinuous()` function that can help you interpret the phylogenetic signal in trait data sets. a. Compare the lambda values of the untransformed tree to the transformed ($\lambda = 0$). b. Compare the Akaike information criterion (AIC) scores of the two models. Which model would you choose based off of AIC score (remember the criteria that the difference in AIC values has to be at least 2)? c. Does this result suggest that there's phylogenetic signal?

Answer 7a: The models are equivalent in likelihood. None of this seems right but I can't figure it where/if things went wrong. **Answer 7b:** The models are equivalent according to AIC. **Answer 7c:** No, the fits are equal between the transformed and untransformed tree.

B) Phylogenetic Signal: Blomberg's K

In the R code chunk below, do the following:
1. correct tree branch-lengths to fix any zeros,

2. calculate Blomberg's K for each phosphorus resource using the `phylosignal()` function,
3. use the Benjamini-Hochberg method to correct for false discovery rate, and
4. calculate Blomberg's K for niche breadth using the `phylosignal()` function.

```
options(warn=-1)
# add a small number to remove zeros
nj_rooted$edge.length <- nj_rooted$edge.length + 10^-7
p_phylosignal <- matrix(NA, 6, 18)
colnames(p_phylosignal) <- colnames(p_growth_std)
rownames(p_phylosignal) <- c("K", "PIC_var_obs", "PIC_var_mean",
                             "PIC_var_P", "PIC_var_z", "PIC_P_BH")

for (i in 1:18) {
  x <- as.matrix(p_growth_std[, i, drop = FALSE])
  out <- phylosignal(x, nj_rooted)
  p_phylosignal[1:5, i] <- round(t(out), 3)
}

p_phylosignal[6, ] <- round(p.adjust(p_phylosignal[4, ], method = "BH"), 3)
signal_nb <- phylosignal(nb, nj_rooted)
print(tibble::as_tibble(p_phylosignal))

## # A tibble: 6 x 18
##       AEP      PEP      G1P      G6P MethCP      BGP      DNA      Peth      Pchol
##       <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
## 1      0          0          0      0.          0          0          0          0          0
## 2 4373.      664.      949.    5.92e+3 351.      536.      259.     1446.    2368.
## 3 8307.     1536.     1878.    3.60e+3 486.     1784.     5222.    1849.    3146.
## 4   0.25     0.066     0.11    7.67e-1 0.386     0.025     0.002     0.33     0.433
## 5  -0.842    -1.32    -1.24    9.60e-1 -0.397    -1.66    -1.30    -0.502   -0.447
## 6   0.643     0.297     0.396    7.82e-1 0.677     0.15      0.036     0.677     0.677
## # ... with 9 more variables: B1 <dbl>, Phyt <dbl>, SRP <dbl>, cAMP <dbl>,
## #   ATP <dbl>, PhenylCP <dbl>, PolyP <dbl>, GDP <dbl>, GTP <dbl>

print(signal_nb)

##           K PIC.variance.obs PIC.variance.rnd.mean PIC.variance.P
## 1 3.427719e-06          49966.78          49484.46          0.551
##   PIC.variance.Z
## 1      0.02334906

print(rownames(p_phylosignal))

## [1] "K"          "PIC_var_obs" "PIC_var_mean" "PIC_var_P"    "PIC_var_z"
## [6] "PIC_P_BH"
```

Question 8: Using the K-values and associated p-values (i.e., "PIC.var.P") from the `phylosignal` output, answer the following questions:

- a. Is there significant phylogenetic signal for niche breadth or standardized growth on any of the phosphorus resources?
- b. If there is significant phylogenetic signal, are the results suggestive of clustering or overdispersion?

Answer 8a: Once we correct for multiple tests, the only significant values are for growth on DNA and cAMP and there is no significant phylogenetic signal. **Answer 8b:** None!

C. Calculate Dispersion of a Trait

In the R code chunk below, do the following:

1. turn the continuous growth data into categorical data,
2. add a column to the data with the isolate name,
3. combine the tree and trait data using the `comparative.data()` function in `caper`, and
4. use `phylo.d()` to calculate D on at least three phosphorus traits.

```
# set up categorical data
p_growth_pa <- as.data.frame((p_growth > 0.01) * 1)
p_growth_pa$name <- rownames(p_growth_pa)

# merge with tree and run tests
p_traits <- comparative.data(nj_rooted, p_growth_pa, "name")
phylo.d(p_traits, binvar = BGP)

##
## Calculation of D statistic for the phylogenetic structure of a binary variable
##
## Data : p_growth_pa
## Binary variable : BGP
## Counts of states: 0 = 4
##                  1 = 35
## Phylogeny : nj_rooted
## Number of permutations : 1000
##
## Estimated D : -0.3540193
## Probability of E(D) resulting from no (random) phylogenetic structure : 0
## Probability of E(D) resulting from Brownian phylogenetic structure : 0.639
phylo.d(p_traits, binvar = DNA)

##
## Calculation of D statistic for the phylogenetic structure of a binary variable
##
## Data : p_growth_pa
## Binary variable : DNA
## Counts of states: 0 = 20
##                  1 = 19
## Phylogeny : nj_rooted
## Number of permutations : 1000
##
## Estimated D : 0.606746
## Probability of E(D) resulting from no (random) phylogenetic structure : 0.034
## Probability of E(D) resulting from Brownian phylogenetic structure : 0.003
phylo.d(p_traits, binvar = cAMP)

##
## Calculation of D statistic for the phylogenetic structure of a binary variable
##
## Data : p_growth_pa
## Binary variable : cAMP
## Counts of states: 0 = 10
##                  1 = 29
## Phylogeny : nj_rooted
## Number of permutations : 1000
```

```
##
## Estimated D : 0.1566464
## Probability of E(D) resulting from no (random) phylogenetic structure : 0
## Probability of E(D) resulting from Brownian phylogenetic structure : 0.289
```

Question 9: Using the estimates for D and the probabilities of each phylogenetic model, answer the following questions:

- Choose three phosphorus growth traits and test whether they are significantly clustered or overdispersed?
- How do these results compare the results from the Blomberg's K analysis?
- Discuss what factors might give rise to differences between the metrics.

Answer 9a: These tables are a little hard to interpret but I think this is right: growth on DNA is significantly overdispersed. Growth on BGP is not overdispersed by may be Brownian and cAMP is the same. **Answer 9b:** K did not detect phylogenetic signal so the results are qualitatively different and hard to compare. **Answer 9c:** D is based around the nature of a phylogenetic relationship while K is based in detecting such a relationship. Perhaps D finds a signal because it expects a signal.

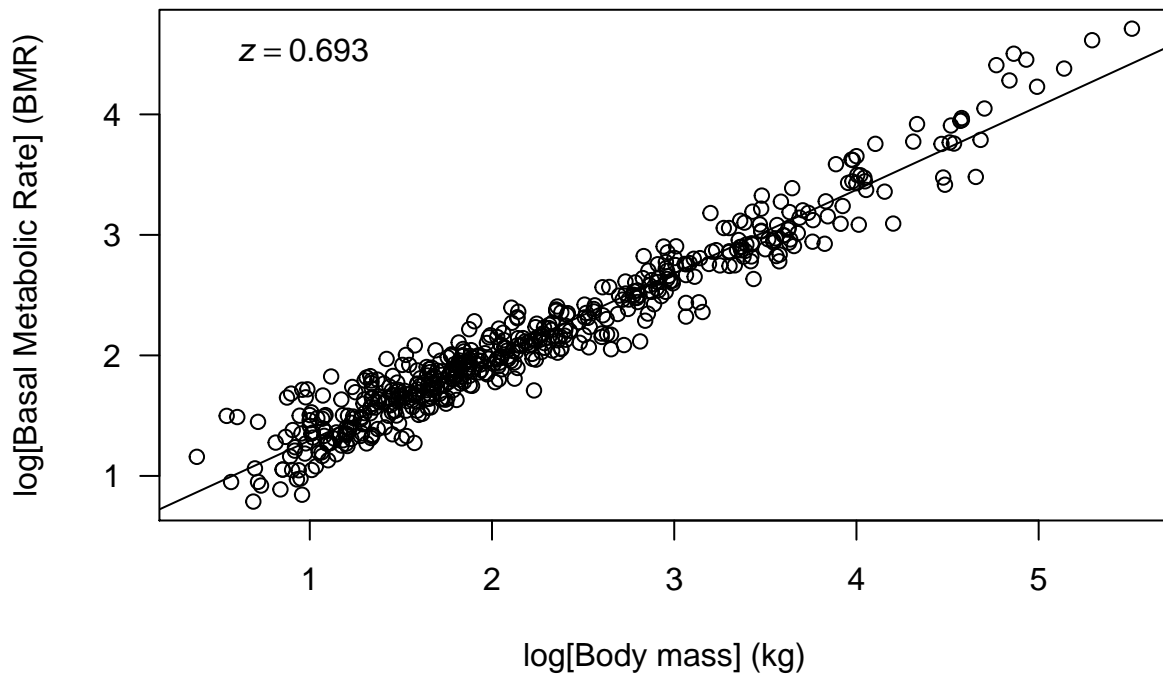
7) PHYLOGENETIC REGRESSION

In the R code chunk below, do the following:

- Load and clean the mammal phylogeny and trait dataset,
- Fit a linear model to the trait dataset, examining the relationship between mass and BMR,
- Fit a phylogenetic regression to the trait dataset, taking into account the mammal supertree

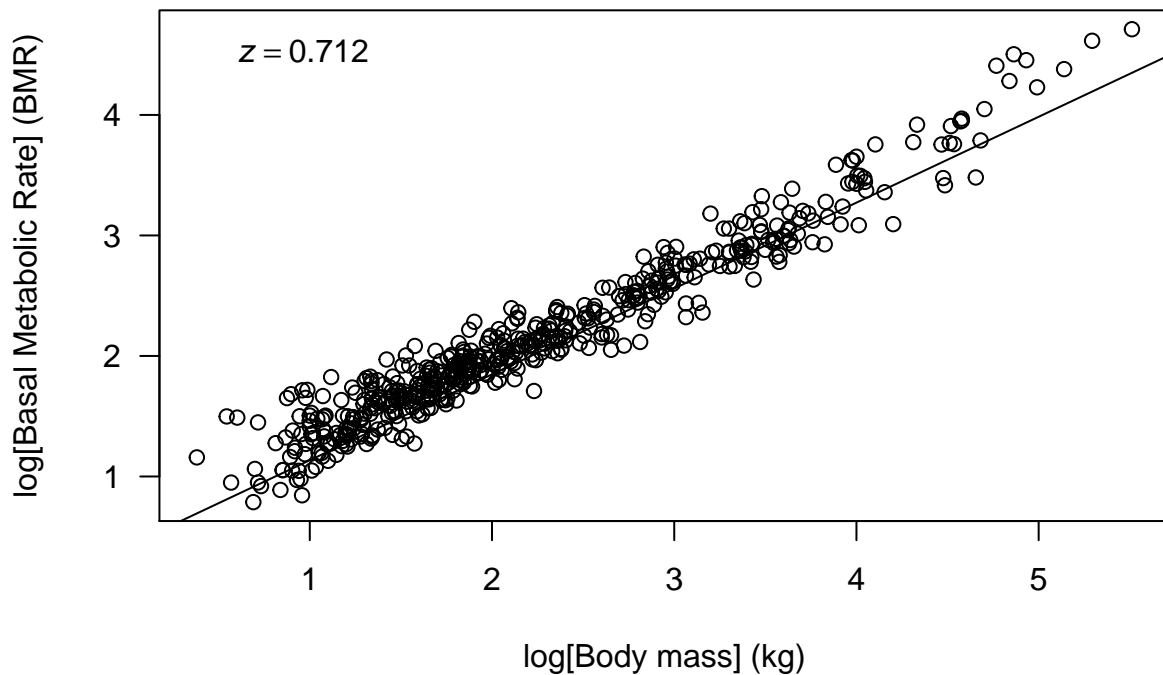
```
mammal.Tree <- read.tree("./data/mammal_best_super_tree_fritz2009.tre")
mammal.data <- read.table("./data/mammal_BMR.txt", sep = "\t", header = T)
mammal.data <- mammal.data[,c("Species", "BMR_.ml02.hour.", "Body_mass_for_BMR_.gr.")]
mammal.species <- array(mammal.data$Species)
pruned.mammal.tree <- drop.tip(mammal.Tree, mammal.Tree$tip.label[-na.omit(match(mammal.species, mammal
pruned.mammal.data <- mammal.data[mammal.data$Species %in% pruned.mammal.tree$tip.label,]
rownames(pruned.mammal.data) <- pruned.mammal.data$Species

# fit the regular OLS like
fit <- lm(log10(BMR_.ml02.hour.) ~ log10(Body_mass_for_BMR_.gr.), data=pruned.mammal.data)
# plot it
plot(log10(pruned.mammal.data$Body_mass_for_BMR_.gr.),
     log10(pruned.mammal.data$BMR_.ml02.hour.), las = 1,
     xlab= "log[Body mass] (kg)", ylab="log[Basal Metabolic Rate] (BMR)")
abline(a= fit$coefficients[1], b = fit$coefficients[2])
b1 <- round(fit$coefficients[2], 3)
eqn <- bquote(italic(z) == .(b1))
text(0.5, 4.5, eqn, pos = 4)
```

```
# fit the phylogenetic regression line
fit.phy <- phylolm(log10(BMR_.ml02.hour.) ~ log10(Body_mass_for_BMR_.gr.),
                  data = pruned.mammal.data, pruned.mammal.tree, model = 'lambda', boot = 0)

# do the plot
plot(log10(pruned.mammal.data$Body_mass_for_BMR_.gr.),
     log10(pruned.mammal.data$BMR_.ml02.hour), las = 1, xlab = "log[Body mass] (kg)", ylab = "log[Basal Metabolic Rate] (BMR)",
     abline(a = fit.phy$coefficients[1], b = fit.phy$coefficients[2])
b1.phy <- round(fit.phy$coefficients[2], 3)
eqn <- bquote(italic(z) == .(b1.phy))
text(0.5, 4.5, eqn, pos = 4)
```



```

print("\n")

## [1] "\n"
print(fit$coefficients[1])

## (Intercept)
##      0.6012236
print(fit$coefficients[2])

## log10(Body_mass_for_BMR_.gr.)
##                      0.6933002
print(fit.phy$coefficients[1])

## (Intercept)
##      0.422397
print(fit.phy$coefficients[2])

## log10(Body_mass_for_BMR_.gr.)
##                      0.7124742

```

- Why do we need to correct for shared evolutionary history?
- How does a phylogenetic regression differ from a standard linear regression?
- Interpret the slope and fit of each model. Did accounting for shared evolutionary history improve or worsen the fit?
- Try to come up with a scenario where the relationship between two variables would completely disappear when the underlying phylogeny is accounted for.

Answer 10a: Correcting for shared evolutionary history can account for many confounding features and can help isolate the features we are interested in. **Answer 10b:** Phylogenetic regression uses the branch lengths between taxa to describe the variance of the residuals which changes how the best fit line relates to the data. **Answer 10c:** The slope is steeper and the fit is better for the phylogenetic regression. **Answer 10d:** This would happen when two traits were completely clustered such that they both evolved simultaneously and once and species in that lineage always has that trait. I'm actually not even sure that is right. If your traits in question are specifically about the relatedness of the organisms that accounting for that relatedness would probably destroy the observation. An example would be if one of your traits was branch length.

7) SYNTHESIS

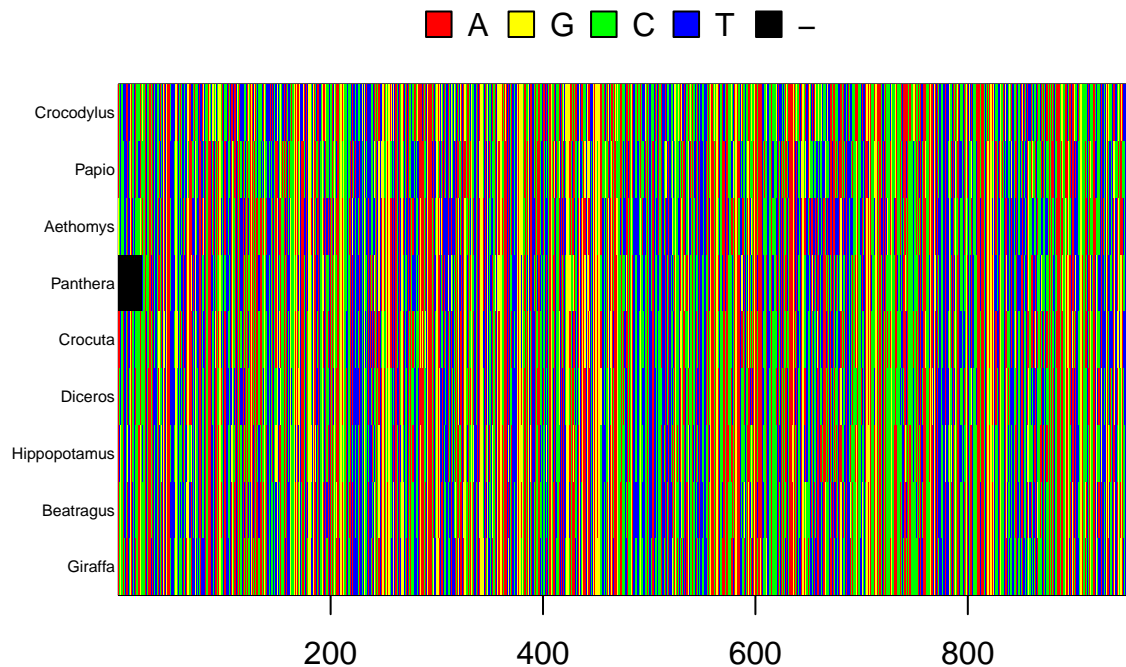
Work with members of your Team Project to obtain reference sequences for 10 or more taxa in your study. Sequences for plants, animals, and microbes can found in a number of public repositories, but perhaps the most commonly visited site is the National Center for Biotechnology Information (NCBI) <https://www.ncbi.nlm.nih.gov/>. In almost all cases, researchers must deposit their sequences in places like NCBI before a paper is published. Those sequences are checked by NCBI employees for aspects of quality and given an **accession number**. For example, here an accession number for a fungal isolate that our lab has worked with: JQ797657. You can use the NCBI program nucleotide **BLAST** to find out more about information associated with the isolate, in addition to getting its DNA sequence: <https://blast.ncbi.nlm.nih.gov/>. Alternatively, you can use the `read.GenBank()` function in the `ape` package to connect to NCBI and directly get the sequence. This is pretty cool. Give it a try.

But before your team proceeds, you need to give some thought to which gene you want to focus on. For microorganisms like the bacteria we worked with above, many people use the ribosomal gene (i.e., 16S rRNA). This has many desirable features, including it is relatively long, highly conserved, and identifies taxa with reasonable resolution. In eukaryotes, ribosomal genes (i.e., 18S) are good for distinguishing course taxonomic

resolution (i.e. class level), but it is not so good at resolving genera or species. Therefore, you may need to find another gene to work with, which might include protein-coding gene like cytochrome oxidase (COI) which is on mitochondria and is commonly used in molecular systematics. In plants, the ribulose-bisphosphate carboxylase gene (*rbcL*), which on the chloroplast, is commonly used. Also, non-protein-encoding sequences like those found in **Internal Transcribed Spacer (ITS)** regions between the small and large subunits of the ribosomal RNA are good for molecular phylogenies. With your team members, do some research and identify a good candidate gene.

After you identify an appropriate gene, download sequences and create a properly formatted fasta file. Next, align the sequences and confirm that you have a good alignment. Choose a substitution model and make a tree of your choice. Based on the decisions above and the output, does your tree jibe with what is known about the evolutionary history of your organisms? If not, why? Is there anything you could do differently that would improve your tree, especially with regard to future analyses done by your team?

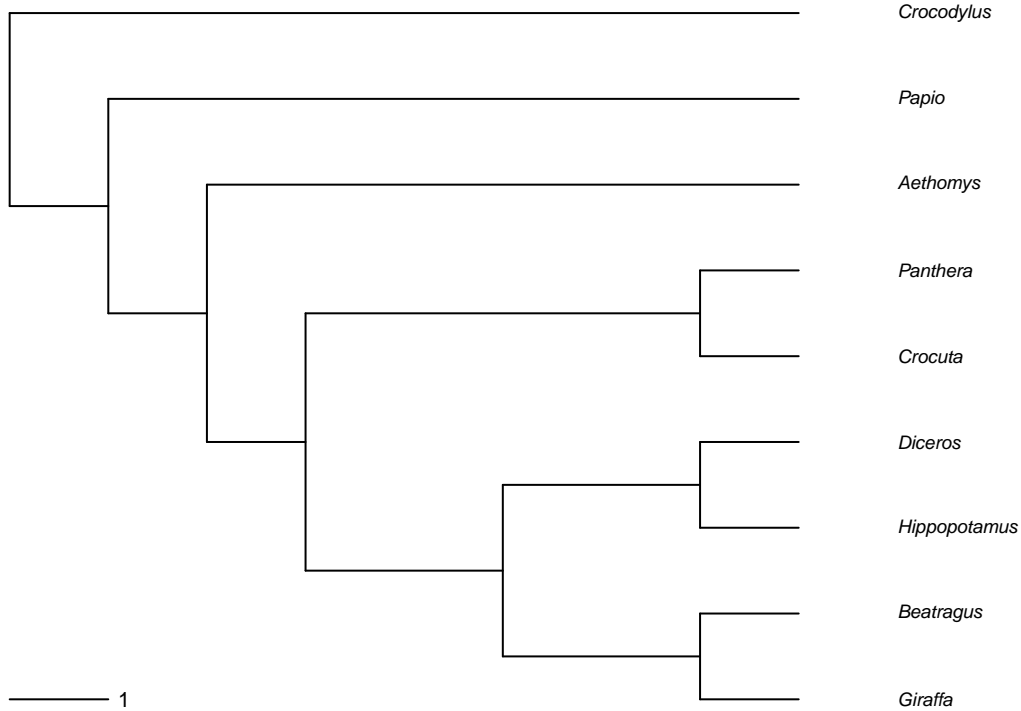
```
# read alignment and convert/viz
read_aln <- read.alignment(file = "~/projects/QB-2021-project/data/taxa.afa", format = 'fasta')
p.DNABin <- as.DNABin(read_aln)
window <- p.DNABin[-6,50:1000] # lets try building the tree on our window
image.DNABin(window, cex.lab = 0.50)
```



```
seq_dist_raw <- dist.dna(window, model = "T92", pairwise.deletion = FALSE)
nj_tree <- bionj(seq_dist_raw)
outgroup <- match("Crocodylus", nj_tree$tip.label)
nj_rooted <- root(nj_tree, outgroup, resolve.root = TRUE)

par(mar = c(1,1,2,1) + 0.1)
plot.phylo(nj_rooted, main = "Neighbor Joining Tree [T92]", "phylogram", use.edge.length = FALSE, direction = "right",
add.scale.bar(cex = 0.7))
```

Neighbor Joining Tree [T92]



This tree looks ok to me. I think that these things cluster in ways that I would, perhaps naively, expect. Hippos and rhinos seem similar. Antelope and giraffes seem similar. I think the main point of improvement would be the inclusion of more species. Our data has 77ish species which we often work with coarse-grained into about 20 families. A tree of these families would allow a lot more analysis.