

BINF*6210—Bioinformatics Software Tools

Karl Cottenie

October 25, 2024

Ulli Bodnar—1351935

Assignment 2

Introduction

Phenotypic variations offer organisms the ability to survive in a catalogue of environments. Pigment patterns can elicit effective camouflage; a simple change in body mass—increasing overall volume, thus reducing the surface area to volume ratio—can allow an organism to survive in a harsh winter environment (McQueen et al., 2022). Bergmann’s rule (Bergmann, C. 1847) explains that “within species and amongst closely related species of homeothermic animals a larger size is often achieved in colder climates than in warmer ones” (Salewski et al., 2017). As explained by Bodnar (2024), “Bergmann’s rule can be observed within a combination of the Pantheria and BOLD System’s databases among species containing the *Lepus* (hare) genus designation.” To continue exploring this interesting topic, this project will build on previous work (Bodnar, 2024) and determine whether related species that are genetically close (within the *Lepus* genus) share similar body mass traits, i.e., within the *Lepus* genus, does Bergmann’s rule follow evolutionary relationships or random chance?

Code Setup

```
1 # Setup
2
3 # Import the libraries so that their functions can be used
4 # Just a heads up, I find R kinda unpredictable with the cache and holding libraries loaded so if this doesn't work
  by just the ones i've uncommented, please just uncomment all of these below and run it again. It will for sure work
  with all loaded, but I am pretty sure I only need the ones that are labelled 'in use'
5 # MIGHT NEED LIBRARIES:
6 # library(stats)
7 # library(sf)
8 # library(randomForest)
9 # library(viridis)
10 # library(sequinr)
11
12 # Libraries in use
13 library(tidyverse)
14 library(rentrez)
15 library(Biostrings)
16 library(fmsb)
17 library(ape)
18 library(dendextend)
19 library(phytools)
20 library(DECIPHER)
21 library(muscle)
22
23
24 # Notes:
25 # Write confirmation checks (that the NAs were removed, etc.)
26 # Write comments that state what the returned information is and even results and interpretation
27 # How is this helping me address my main question
28
29 # Set working directory to utilize data from datasets
30 setwd("/Users/ullibodnar/Documents/School/Guelph Masters/Bioinformatics Software Tools/Assignment 2/code")
31
32
```

```

33
34 ▾ # Global variables -----
35
36 ▾ #### CHECK ONLINE IF THESE ARE THE CORRECT CHOICES FOR WHAT I'M DOING ####
37
38 missingData <- 0.01
39 lengthVar <- 50
40 chosenModel <- "K80" # K2P
41 clusteringMethod <- "ML"
42
43
44
45 ▾ # Global functions -----
46
47 # Use the species name from BOLD dataset to obtain corresponding mass (in grams) from Pantheria dataset
48 ▾ getAverageMass <- function(df, speciesName) {
49   averageMass <- df >
50     filter(MSW05_Binomial == speciesName) >
51     pull("5-1_AdultBodyMass_g")
52
53   # Substitute -999 masses
54   ifelse(averageMass >
55     length() == 0
56     , 0, averageMass)
57 ▴ }
58
59
60

```

```

61 ▾ # Importing data -----
62 # BOLD DB
63 # API CALL TO SHOW I KNOW HOW TO DO IT
64 # Lepus <- read_tsv("http://v3.boldsystems.org/index.php/API_Public/combined?taxon=Lepus&format=tsv")
65 # write_tsv(Lepus, "lepus_bold_data.txt")
66
67 # Import lepus data from BOLD database to have larger selection of COI-5P sequences
68 rawBoldLepus <- read_tsv(file = "../data/lepus_bold_data.txt")
69
70 # Import Pantheria DB for body mass in grams data
71 pantheriaData <- read_tsv(file = "../data/Pantheria.tsv")
72
73 # NCBI's nucleotide ---
74 # Determine possible database search locations
75 entrez_dbs()
76
77 # Sequence length range chosen because NCBI returned 671bp length for this sequence of interest when inspecting the
  database. Range to catch ones with missing data
78 # Query NCBI's nucleotide database to test the waters of possible data
79 lepusSearch <- entrez_search(db = "nucleotide", term = "Lepus[ORGN] AND COI AND 600:800[SLEN]", retmax = 100)
80
81 # Determine result count and change entrez_search to provide unique IDs for each possible returned result
82 lepusRetmax <- lepusSearch$count
83 lepusSearch <- entrez_search(db = "nucleotide", term = "Lepus[ORGN] AND COI AND 600:800[SLEN]", retmax = lepusRetmax)
84
85 # Fetch the data as fasta
86 lepusFetch <- entrez_fetch(db = "nucleotide", id = lepusSearch$id, rettype = "fasta")
87
88 # Write the data to fasta file; commented out because the data have already been imported
89 # write(lepusFetch, "lepus_fetch.fasta", sep = "\n") --- Done on Oct. 18
90
91 # Import already created NCBI data file
92 nucleotideLepusStringSet <- readDNAStringSet("../data/lepus_fetch.fasta")
93
94 # Convert to a dataframe object for easy manipulation
95 nucleotideLepus <- data.frame(title = names(nucleotideLepusStringSet), nucleotides = paste(nucleotideLepusStringSet))
96

```

```

99 # Formatting data
100 # Subset bold lepus for only sections needed to reduce cognitive load when viewing dataframe
101 # remove ones that don't contain nucleotide data and ones that aren't COI
102 boldLepus <- rawBoldLepus[, c("processid", "species_name", "markercode", "nucleotides")] >
103   filter(!is.na(nucleotides)) >
104   filter(markercode == "COI-5P")
105
106 # Change BOLD's processid column name to "id" to match NCBI database
107 names(boldLepus)[names(boldLepus) == "processid"] <- "id"
108
109 # Manipulate the NCBI dataframe to have correct column names for later merge with boldLepus
110 nucleotideLepus$id <- word(nucleotideLepus$title, 1L)
111 nucleotideLepus$species_name <- word(nucleotideLepus$title, 2L, 3L)
112
113 # filter out any non cytochrome sequences
114 nucleotideLepus <- filter(nucleotideLepus, grepl("cytochrome oxidase subunit", title))
115
116 # Add marker code column and rearrange columns to allow for clean merge with bold
117 nucleotideLepus$markercode <- "COI-5P"
118 nucleotideLepus <- nucleotideLepus[, c("id", "species_name", "markercode", "nucleotides")]
119
120
121 # MERGE DATAFRAMES
122 lepusSeq <- merge(nucleotideLepus, boldLepus, all = T)
123
124 # Map body masses from Pantheria to the subsetted column for downstream analysis of body masses
125 lepusSeq$mass_g <- purrr::map(lepusSeq$species_name, function(x) {getAverageMass(pantheriaData, x)}) >
126   as.numeric()
127
128 # Remove duplicates, trim Ns from the ends and remove gaps, remove entries with NA species_name, and remove entries
  with no mass data
129 lepusSeq <- lepusSeq[!duplicated(lepusSeq$nucleotides), ]
130
131 lepusSeq <- lepusSeq >
132   filter(!is.na(species_name)) >
133   filter(mass_g > 0) >
134   mutate(nucleotides2 = str_remove_all(nucleotides, "^N+|N+$|-")) >
135   filter(str_count(nucleotides2, "N") <= (missingData * str_count(nucleotides))) >
136   filter(str_count(nucleotides2) >= median(str_count(nucleotides2)) - lengthVar & str_count(nucleotides2) <= median
  (str_count(nucleotides2)) + lengthVar)
137
138 # Create a subset dataframe containing a random sequence from each species; this allows for comparison between just
  one sample per species
139 set.seed(1234) # so that we get the same result
140
141 lepusSeqSubset <- lepusSeq >
142   group_by(species_name) >
143   slice_sample(n = 1) > # Randomly selects one row per species, because Karl told me to do it :)
144   ungroup() >
145   as.data.frame()

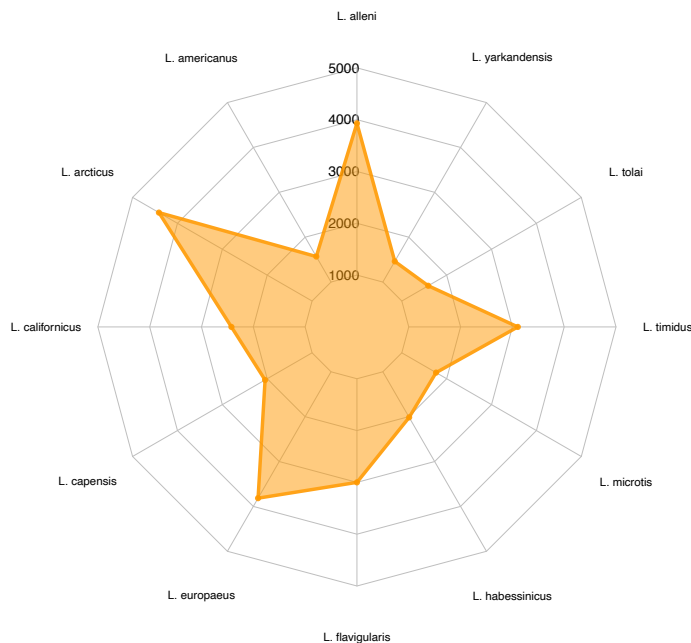
```

Figures

Figure 1

```
148
149 # View data in radar chart
150 # The following chart idea taken from https://r-graph-gallery.com/142-basic-radar-chart.html
151
152 # Map the data as names and mass to visualize in a radar chart
153 massAndNames <- as.data.frame(matrix(lepustSeqSubset$mass_g , ncol=12))
154 colnames(massAndNames) <- paste("L.", word(lepustSeqSubset$species_name, 2L), sep = " ")
155
156 # Add in the upper and lower limits to comply to the formatting requirements for fmsb library
157 massAndNames <- rbind(rep(5000,12), rep(1000,12), massAndNames)
158
159 # check data formatted properly
160 # head(massAndNames)
161
162 # Remove styling from other graph so it doesn't display weird when running
163 par(mar=c(1,1,1,1))
164
165 # Render the radar chart
166 radarchart( massAndNames, axistype=1,
167             #customize the polygon, grid, and labels
168             pcol=rgb(1,0.6,0,0.9) , pfc=rgb(1,0.6,0,0.5) , plwd=4 ,
169             cglcol="grey", cglty=1, axislabcol="black", caxislabels=seq(1000,5000, 1000), cglwd=0.8,
170             vlce=0.8
171 )
172
173
```

Body Mass (in Grams) of Each *Lepus* Species



Explanation for the radar chart:

Since body mass was the only trait being explored (absent are extra traits such as latitude), a radar chart was employed to effectively depict the body mass (in grams) of each of the *Lepus* species to better understand the distribution of weights and compare across species.

Alignment (Preliminary to Figure 2)

```
174
175 ▾ # Aligning sequences -----
176 # put entire lepus subset into DNASTringSet to work with the library
177 lepusSeqSubset$nucleotides2 <- DNASTringSet(lepuseqSubset$nucleotides2)
178
179 # Map the names, substituting L. for Lepus for readability.
180 names(lepuseqSubset$nucleotides2) <- paste("L.", word(lepuseqSubset$species_name, 2L), sep = " ")
181 # Check that it worked
182 names(lepuseqSubset$nucleotides2)
183
184 # Conduct alignment with muscle
185 lepusSeqSubsetAlignment <- DNASTringSet(muscle::muscle(lepuseqSubset$nucleotides2))
186 # Check it out in the browser to see if anything is out of place → originally, I saw that one of the names was NA
187 # and I had to go back to remove NA species_name entries
187 #BrowseSeqs(lepuseqSubsetAlignment)
188
189
190
191 ▾ # Clustering -----
192 # Convert to a dataclass used by Ape for distance clustering
193 lepusSeqBin <- as.DNABin(lepuseqSubsetAlignment)
194 # Check conversion is correct
195 class(lepuseqBin)
196
197 # Create distance matrix
198 distanceMatrix <- dist.dna(lepuseqBin, model = chosenModel, as.matrix = TRUE, pairwise.deletion = TRUE)
199 # Check out distance matrix worked
200 head(distanceMatrix)
201
202 # PHYLOGENETIC TREE
203 clustersLepusCOI <- DECIPHER::TreeLine(lepuseqSubsetAlignment,
204                                       myDistMatrix = distanceMatrix,
205                                       method = clusteringMethod,
206                                       model = chosenModel,
207                                       reconstruct = TRUE,
208                                       maxTime = 0.01)
209
```

```
33
34 ▾ # Global variables -----
35
36 ▾ #### CHECK ONLINE IF THESE ARE THE CORRECT CHOICES FOR WHAT I'M DOING ####
37
38 missingData <- 0.01
39 lengthVar <- 50
40 chosenModel <- "K80" # K2P
41 clusteringMethod <- "ML"
42
```

Explanation for the chosen method:

The maximum likelihood (ML) method was chosen for evolutionary tree generation because of its “statistical consistency, robustness,...and ability to...make full use of original data within a statistical framework” (Zou, Yue, et al., 2024).

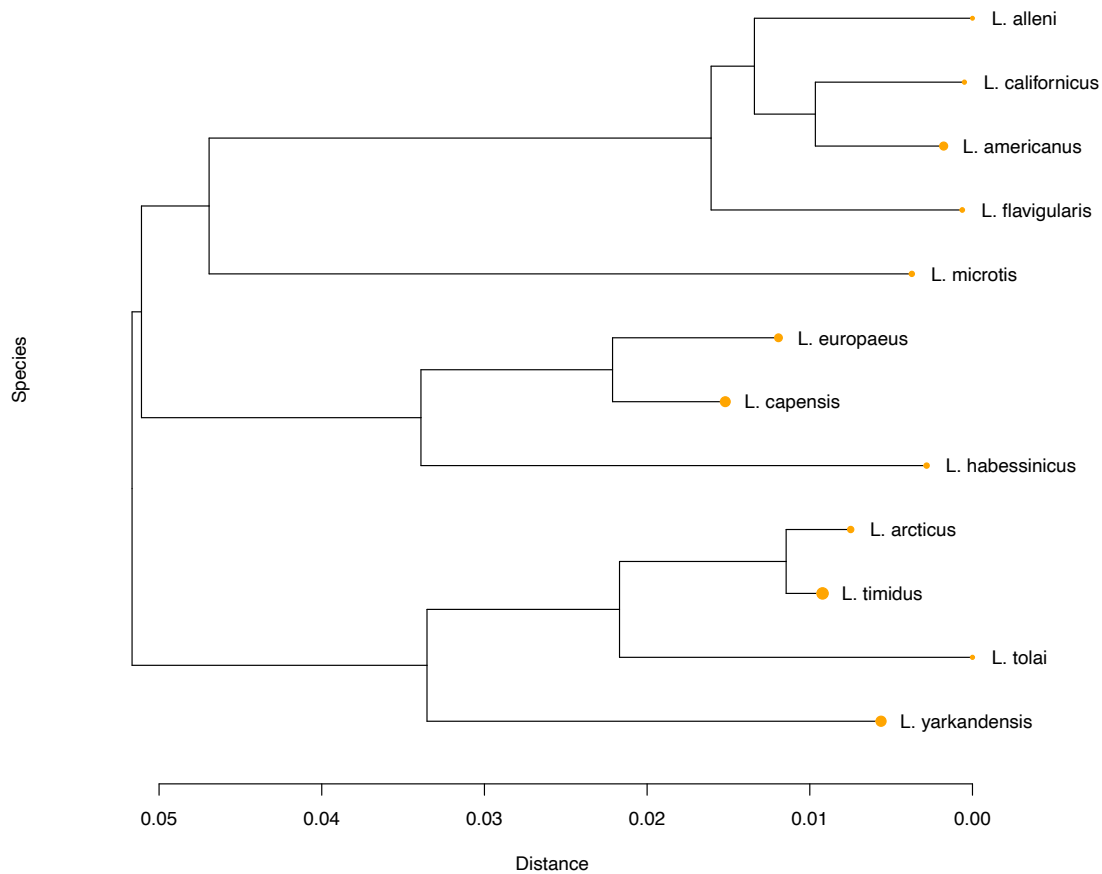
Explanation for using the K80 model:

The Kimura 2-parameter (K2P/K80) model is typically used early in phylogenetic analysis for constructing trees (Salemi et al., 2009, p. 163), which fits handsomely within the context and scope of this analysis.

Figure 2

```
210
211 # Bottom, left, top, right margins
212 par(mar=c(5,5,1,10))
213
214 # Create a vector of masses scaled relative to the max mass
215 maxMass <- max(lepuseqSubset$mass_g)
216 scaledMass <- lepuseqSubset$mass_g / maxMass * 2
217
218 # Plot the phylogenetic tree
219 clustersLepusCOI >->
220   set("leaves_pch", 20) >->
221   set("leaves_cex", scaledMass) >->
222   set("nodes_col", "orange") >->
223   set("labels_col", "black") >->
224   plot(horiz = TRUE, xlab = "Distance", ylab = "Species")
225
```

Phylogenetic Tree Depicting the Genetic Relationship Among Lepus Species Based on COI Gene Sequences, with Relative Body Mass Represented by Tip Sizes



● = Relative body mass of *Lepus* sp.

Table 1

```
227 # Test for phylogenetic conservatism
228
229 # Export tree as a format used by ape package
230 tree_phylo <- as.phylo(clustersLepusCOI)
231
232 # Lambda estimation will test the extent to which body mass shows phylogenetic signal.
233 # If it is close to 1, the trait follows the phylogeny
234 # If it is close to 0, the trait does not really follow the phylogenetic signal and does not exhibit phylogenetic conservatism
235 lambda_estimation <- phylosig(tree_phylo, lepusSeqSubset$mass_g, method = "lambda")
236
237 print(lambda_estimation)
238
```

Summary of Important Phylogenetic Signal Outputs

Phylogenetic Signal Lambda	LogL (lambda)
4.787 x 10 ⁻⁵	-99.78

Results and Discussion

A phylogenetic tree (Figure 2) and phylogenetic signal lambda (Table 1) were created and calculated with data of *Lepus* species from the BOLD, NCBI, and Pantheria databases. The data were aligned with the muscle library’s muscle function. A distance matrix was created by utilizing ape’s dist.dna function and a phylogenetic tree was created using DECIPHER’s TreeLine function. The phylogenetic signal was computed via phytool’s phylosig function. The output of the phylogenetic signal lambda test demonstrated that there is no significant evidence that the trait (body mass) follows the phylogenetic signal. Within the scope of this assessment, the trait does not exhibit phylogenetic conservatism, i.e., related species that are genetically close (within the *Lepus* genus) do not share similar body mass traits. The results of this analysis can be confirmed following visual inspection of the phylogenetic tree (Figure 2): closely related species clustered within the same clade do not consistently appear to have a similar body mass, e.g., *L. arcticus* and *L. timidus*.

There are several limitations to this work that should be considered. First, the data were filtered to remove species entries with missing data. While this step was necessary to perform an analysis, i.e., sequence alignment cannot be conducted on a species that has no sequence data, it may have removed key species that are necessary for a thorough analysis. The absence of these species may have adversely impacted the tree and, more importantly, the phylogenetic signal test. In a future study, 47S ribosomal genes may be better suited for analysis than COI genes. Even though COI genes are robust in their ability to determine phylogenetic relationships, the 47S ribosomal genes may produce a more reliable analysis when determining whether the body mass trait follows phylogeny (Law et al., 2024) among *Lepus* species.

Acknowledgements

The author wishes to thank Sophia Teng for the moral support and the conversations that prompted him to question the strength of his analysis and adjust where suitable. The author also wishes to acknowledge that he could not have created the figures without the help of The R Graph Gallery.

References

Bergmann, C. (1847). Über die verhältnisse der wärmeökonomie der thiere zu ihrer grösse. *Gött. Stud*, 1, 595-708.

Bodnar, U. (2024). *Assignment 1*. [Unpublished manuscript]. University of Guelph.

Law, Pui Pik, et al. “Ribosomal DNA Copy Number Is Associated with Body Mass in Humans and Other Mammals.” *Nature Communications*, vol. 15, no. 1, June 2024, p. 5006. *DOI.org (Crossref)*, <https://doi.org/10.1038/s41467-024-49397-5>.=

McQueen, Alexandra, et al. “Thermal Adaptation Best Explains Bergmann’s and Allen’s Rules across Ecologically Diverse Shorebirds.” *Nature Communications*, vol. 13, no. 1, Aug. 2022, p. 4727. *DOI.org (Crossref)*, <https://doi.org/10.1038/s41467-022-32108-3>.

Salemi, M., Vandamme, A.-M., & Lemey, P. (Eds.). (2009). *The phylogenetic handbook : a practical approach to phylogenetic analysis and hypothesis testing* (2nd ed.). Cambridge University Press.

Salewski, Volker, and Cortney Watt. “Bergmann’s Rule: A Biophysiological Rule Examined in Birds.” *Oikos*, vol. 126, no. 2, Feb. 2017, p. oik.03698. *DOI.org (Crossref)*, <https://doi.org/10.1111/oik.03698>.

Zou, Yue, et al. “Common Methods for Phylogenetic Tree Construction and Their Implementation in R.” *Bioengineering*, vol. 11, no. 5, May 2024, p. 480. *DOI.org (Crossref)*, <https://doi.org/10.3390/bioengineering11050480>.