

Lab 2 — Phylogenetic trees

Evolution of Language and Music

October 22, 2018

Goals Last week we implemented computer simulations of evolution in R. This week, we extend this simulation, and think about the patterns of genetic variation the evolutionary process leaves in a population. We look at methods for *phylogenetic tree reconstruction*. These methods use the genetic variation in the most recent (current) generation to reconstruct the evolutionary history of a population or a set of species.

We use a simple clustering algorithm¹ for phylogenetic tree reconstruction. The goal of this lab is to learn about the possibilities and difficulties that are involved in phylogenetic tree reconstruction.

1 Simulated evolution (continued)

In the previous lab, we simulated the evolution of strings of symbols, and looked at the effect of different fitness functions. Today we will repeat this simulation, but during the evolutionary process we will keep track of ancestry, so that we can reconstruct family trees of different individuals.

- Download the scripts for this week's lab and extract all the scripts to a folder on your computer.
- Start R-studio (or an R terminal) and set R studio's working directory to the folder where the scripts reside (if you forgot how, have a look at this webpage)

We will start with a very small simulation.

- If you are working on a university computer, you might have to reinstall the package `stringr`:
`install.packages("stringr")`

The script `lab-3.R` runs the same simulation we saw in the previous lab. This time, however, the script also generates a matrix called `parent_matrix` that specifies the parent of each member in each generation. Where the parent is the individual of the previous generation whose genetic material was inherited.

¹An algorithm is a description of a series of steps to do arrive at a certain end result or perform a calculation. Algorithms can for example be implemented by a computer program.

At the end of the simulation, a plot illustrating the development of both the average population fitness and the diversity of the population is generated.

- Change the parameters at the top of the file `lab-3.R`: Set both `population_size` and `simulation_length` to 10.
- ? What values do you expect on the y-axes of these plots?
- ? What do you think the curves of average population fitness and population diversity look like?
- ? At what point do you expect the curves to start and finish?
- Run the script by executing the following command in the console:

```
source('lab-3.R')
```

Are the results as you expected?

- Visualise the parent matrix by running

```
print_parent_matrix(parent_matrix)
```
- ? Where in this plot can you find the first generation?
- ★ Follow some paths up and down. Why do downward paths often end in dead ends, whereas upward paths always go all the way up?
- Change the parameters back to their original settings:

```
population_size <- 100
simulation_length <- 1000
```
- Run the simulation again (this may take a while).

Printing the parent matrix for such large simulations is not very helpful (you may try if you want), because the network is too dense to properly visualise. Rather than looking the parent matrix, we will use the parent matrix to reconstruct a family tree for only the last generation (that is, we only look at the members of previous generations whose offspring appears in the last generation.)

- From the data you just generated, generate a family tree with the function `reconstruct_tree` and print it with the function `print_tree`, which will generate a textual representation of the phylogenetic tree:

```
tree <- reconstruct_tree(parent_matrix)
print_tree(tree) \end{verbatim}
```

To generate a visual representation of the tree, we will use an online tree viewer.

- Copy everything between the double quotes in the output of `print_tree`.
- Go to <http://evolangmus.knownly.net/newick.html>.
- On the website, change the tree type from *Cladogram* to *Rectangular cladogram*.
- Paste the tree representation you copied into the text area.
- Click "show". Now, you can zoom in by scrolling and move the tree around by dragging the mouse.
- ★ As far as you can judge, how many generations ago did the LCA of the

current population live?

- ★ Which aspects of evolution leave traces that we can detect in the current generation and which aspects do not?

2 Phylogenetic reconstruction with R

In evolutionary research, the elaborate ancestry information represented by the parent matrix is usually not available. To reconstruct family trees we have to resort to different methods. Information about when species branched off ('speciated') can be deduced from genetic variation in the current population. For instance, horses are genetically more similar to donkeys than to, say, frogs. So, the last common ancestor of horses and frogs most likely lived much further in the past than the last common ancestor of horses and donkeys. In other words, the branch that would eventually evolve into frogs split off from the branch that would eventually evolve into horses earlier than the branch that would eventually evolve into donkeys. This type of analysis is called phylogenetic reconstruction. It's based on genetic similarity between members of the current generation, which we measure using a *distance measure*. There are R packages that can automatically perform this reconstruction. Let's start with installing these packages:

- Install the packages `ape` and `phangorn` using:

```
install.packages("ape")  
install.packages("phangorn")
```

Activate them using:

```
library(ape)  
library(phangorn)
```

The `phangorn` package comes with a dataset that contains real genetic data (i.e., RNA samples) from many different species. You can load this dataset by typing:

```
data(Laurasiatherian)
```

To show a summary of the data you can type `str(Laurasiatherian)`. The data originates from <http://www.allanwilsoncentre.ac.nz/>. If you want to find out more about this data, have a look at that website.

We will try to reconstruct a phylogenetic tree for these species. That is, we will try to reconstruct when different species branched off from each other, based only on genetic information of the current population (the last generation). The first step is to measure 'genetic distance' between the genetic samples for each species. For simplicity, we assume that all species ultimately originate from a single common ancestor (an uncontroversial assumption in evolutionary biology), and that species have diverged genetically by picking up mutations at a roughly constant rate (a more problematic assumption).

? (OPTIONAL) If you would like to understand this in more depth, try and convince yourself that the phylogenetic tree reconstruction method we described requires the second assumption, and that this assumption is problematic when considering evolution in the real world.

The distance between strings of DNA or RNA is typically measured by counting the number of mutations required to change one into the other. Because of the second assumption, the genetic distance between two species is proportional to the time that has passed since their last common ancestor.

- Select five species from the Laurasiatherian dataset (for instance three that you think are closely related and two that are more distantly related)
 - Create a subset of the data containing just these five species using:

```
mysubset <- subset(Laurasiatherian, subset=c(19,20,28,29,30))
```

The numbers correspond to the position of the species in the list printed by `str(Laurasiatherian)`, i.e. Platypus = 1, Possum = 3, etc.) You have to replace these numbers by the numbers corresponding to the species that you chose.
 - Verify that your subset contains the right species using:

```
str(mysubset)
```
 - Compute the *pairwise distance* between all elements in the set using the function `dist.ml` and print it

```
distance_matrix <- dist.ml(mysubset)
print(distance_matrix)
```
- ? The results are stored in a *distance matrix*. How can you read off the distance between two species from this matrix?
- ? Why are the numbers on the diagonal of this matrix zero?
- ? Do the computed distances correspond to your intuitions about the selected species' relatedness?
- ? Using pen and paper, or your favourite drawing software, reconstruct a phylogenetic tree that describes the evolutionary relations between your selected species. Use the principles described earlier. You shouldn't need to do any calculations.

We can use a simple method called 'hierarchical clustering' to build such phylogenetic trees automatically. Hierarchical clustering can be done with a simple algorithm that follows these steps:

1. treat each datapoint (for example, a RNA sample of a species) as a separate "cluster" containing just one datapoint;
2. compute the distances between all clusters (using some distance measure; for example, genetic distance);
3. merge the two clusters that are nearest to each other into a new cluster;
4. repeat steps 2 and 3 until only all datapoints are in cluster.

To construct a phylogenetic tree, we can think of each merging of clusters as the joining of two branches. In the simplest version of this algorithm, we define 'distance' between a cluster A and a cluster B as the average distance between

any datapoint in A and any datapoint in B (a slightly more complicated method, Ward's clustering, uses the square root of the average of the squared point-to-point distances).

- ★ Using the distances between species in `mysubset`, manually perform three cycles of the hierarchical clustering algorithm with pen and paper.

The `phangorn` package we installed earlier provides pre-defined functions implementing different hierarchical clustering methods.

- Generate a phylogenetic tree for your subset and plot it using:

```
tree <- upgma(distance_matrix, method='average')  
plot(tree)
```

Is the tree the same as the one that you created before with pen and paper?
- Create a tree for the entire dataset. Does it agree with your expectations?
- (OPTIONAL) Try different methods for computing the distance between clusters by changing the parameter `method` (options are, for instance, *ward.D*, *single* and *median*).
- ? (also OPTIONAL) Do you notice any changes in the resulting phylogenetic trees?

3 Phylogenetic reconstruction of simulated data

We will now investigate what happens if we perform phylogenetic analysis on the population resulting from our own simulated evolution. Remember that, since this is a simulation over which we have full control, we can reconstruct the *actual* phylogenetic tree using the information that we stored in the parent matrix.

- Run the script `lab-3.R` again to generate a new population and parent matrix
- Generate a distance matrix of the last generation from your simulation using the function `compute_distance_matrix`:

```
distance_matrix <- compute_distance_matrix(population)
```
- Reconstruct a phylogenetic tree with the `upgma` function (choose your own *method*) and plot it:

```
tree <- upgma(distance_matrix, method='ward.D')  
plot(tree, cex=0.3)
```

The parameter `cex` sets the font size of the plot, adjust it if the numbers are illegible.
- Now generate the *actual* family tree of the simulation by running

```
gold_standard_tree <- reconstruct_tree(parent_matrix)  
print_tree(gold_standard_tree)
```

and plot it using the online tree visualiser we have used before.

- ? How well does the reconstruction produced by the hierarchical clustering algorithm match the actual family tree?
- ★ How can you explain the differences between the reconstructed and the actual family tree?