

The Jaatha HowTo

Lisha Mathew, Paul R. Staab and Dirk Metzler

Version 2.2

1 Introduction

Jaatha is a fast composite likelihood method to estimate model parameters of the evolutionary history of (at the moment) two related species or populations. To do so, it uses SNP data from multiple individuals from both species and – optionally but highly recommended – one or more outgroup sequences. This HowTo describes the method and gives an example of using its implementation as an R package `jaatha`.

The package itself can be obtained from CRAN using

```
install.packages("jaatha")
```

or downloaded from http://evol.bio.lmu.de/_statgen/software/jaatha. Jaatha runs on R under Windows, OS X (Mac) and Linux with the following restrictions: On Windows the parallelization and finite sites simulation using Seq-Gen is currently not supported.

A more detailed description of the algorithm can be found in Mathew et al. [2013]. Further information about the R functions used in this document can be obtained by calling `help()` with the functions name as argument.

Please cite the above paper when using Jaatha in a publication.

2 A demographic model

Before we can apply Jaatha to estimate parameters, we first need to create a model of the evolutionary history of the two species. Jaatha cannot account for the effects of selection, hence we assume a neutral evolution. It can estimate the effects that either “demographic” events – like an expansion of the size of one population or migration between the two populations – or molecular events – like mutation or recombination – have on the genome of the populations. To emphasize that we can not include selection, we refer to a scenario of the evolution of the two species under such events as a “demographic model”.

For now, assume that we know that our two species are closely related; hence they must have separated at a certain time in the past. There may still be gene flow ongoing between them to which we will refer to as migration from one population into the other. Hence, we could propose the simple demographic model described in Figure 1.

To specify that model in R, we first need to load `jaatha`.

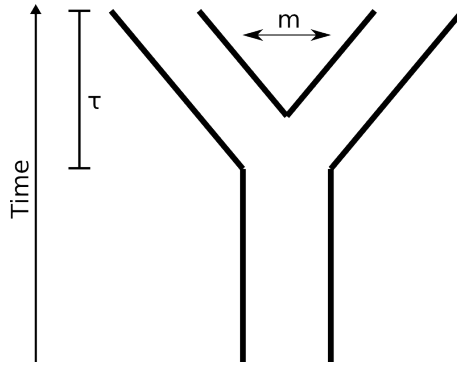


Figure 1: A simple demographic model: The ancestral population splits into two populations τ time units ago, and afterwards individuals migrated from one population to the other with a migration rate M . Mutations are occurring with rate θ and recombination with a known rate ρ .

```
library(jaatha)

## Loading required package: Rcpp
## Loading required package: methods
```

We can now create an 'empty' demographic model `dm` using the `dm.createDemographicModel()` function:

```
dm <- dm.createDemographicModel(sample.sizes = c(12, 14), loci.num = 50,
                                seq.length = 1000)
```

The parameter `sample.sizes` here corresponds to the number of individuals we have sampled from the first population and second population respectively. The second argument `loci.num` states that we are using data from 70 loci while `seq.length` gives the (average) length of each loci¹.

We can now successively add the other assumptions of our model:

```
dm <- dm.addSpeciationEvent(dm, 0.1, 5, new.time.point.name = "tau")
dm <- dm.addSymmetricMigration(dm, 0.01, 5, new.par.name = "M")
dm <- dm.addMutation(dm, 1, 20, new.par.name = "theta")
dm <- dm.addRecombination(dm, fixed = 20)
```

The first parameter is the demographic model to which we want to add an assumption/feature. The two following numbers represent the range for the corresponding parameter. The lower border has to be strictly greater the zero, as we are using a logarithmic transformation of the parameter space. The parameters are scaled as in the popular simulation program `ms` [Hudson, 2002] that we use for simulations:

¹This is only used when a finite sites model is assumed or if intra-locus recombination is included.

- The parameter for the *speciation* event is the split time τ , which states how many generations ago the split of the population has occurred. As usual in population genetics, it is measured in units of $4N_1$ generations ago, where N_1 is the (diploid) effective population size of the first population.
- The parameter for the (symmetric) *migration* is the scaled migration rate M , which is given by $M = 4N_1m$, where m is the fraction of individuals of each population which are replaced by immigrants from the other population each generation.
- The *mutation* parameter θ is $4N_1$ times the neutral mutation rate per locus.
- Finally, the *recombination* parameter ρ is $4N_1$ times the probability of recombination between the ends of the locus per generation.

Keep in mind that a ‘good’ model – which is one that approximates the real demographic history but is also as simple as possible – is crucial for obtaining meaningful estimates in the end. Jaatha will always try to find the parameters that make the model fit best to your data. If the model does not fit to the data at all, Jaatha will still return estimates, but they will not be meaningful.

3 Theoretical Background

It is important to understand the key concepts behind Jaatha before we can apply it. Like many estimation methods that rely on simulations, Jaatha tries to find the parameters that best fits ² to your data by simulating artificial data for many different parameter combinations. It uses a learning algorithm to determine how the different parameter values influence the simulated data and uses that knowledge to find the best parameter combination for your data.

You can imagine Jaatha as a method that runs through the parameter space – the space of all possible parameter combinations, in our example a cube with borders from 0.1 to 5, 0.01 to 5 and 1 to 20 – simulating in a small part of the parameter space around the current position (we call this area a *block*). It then searches the new maximum of the current blocks and moves to it, builds a new block around it and so on. The search finally stops when the likelihood cannot be improved anymore or a maximal number of steps has been reached.

To compare the simulated data to the real one, Jaatha uses *summary statistics* of the data. As default, it calculates the Joint Site Frequency Spectrum (JSFS) of the data and further summarizes it by evaluating different sums over the JSFS. Please refer to Naduvilezhath et al. [2011] for a detailed description.

4 Importing Your Data

To run Jaatha, you need to calculate the JSFS of your data. To do so, you can use the function `calculateJsfs()`. This function accepts data imported with the `read.dna()` function from the package *ape*. It assumes that you provide a joined, aligned data set with multiple samples from two populations and –

²for Jaatha, the ‘best’ parameter combination is the one with the highest composite likelihood

optional but highly recommended – one or more outgroup sequences. Additionally, you must provide the numbers of the sequences in the dataset that belong to population one, population two, and the outgroup, respectively.

Please consult the documentation of *ape* in order to get more information about `read.dna()`. For example, the import of the data could look like this:

```
library(ape)
# The path to the data
sample.file <- system.file("example_fasta_files/sample.fasta",
  package = "jaatha")

# Reading the data
sample.data <- read.dna(sample.file, format = "fasta", as.character = TRUE)

# Calculating the JSFS
sample.jsfs <- calculateJsfs(sample.data, pop1.rows = 3:7, pop2.rows = 8:12,
  outgroup.row = 1:2)
sample.jsfs

##      [,1] [,2] [,3] [,4] [,5] [,6]
## [1,]    0    1    0    0    0    0
## [2,]    1    1    1    0    0    0
## [3,]    2    0    1    0    0    0
## [4,]    1    1    0    2    0    1
## [5,]    1    0    0    0    3    0
## [6,]    1    0    0    0    0    0
```

For the purpose of this HowTo, we will use a simulated JSFS, for which we know the real parameters:

```
# Real parameters: M = 1, tau = 1 and theta = 10
real.pars <- c(1, 1, 10)

# Simulate a JSFS with this parameters
sum.stats <- dm.simSumStats(dm, real.pars)
jsfs <- sum.stats$jsfs

# Print the upper left part of the JSFS
jsfs[1:10, 1:10]

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## [1,]    0 473 211 132  79  72  54  33  32  27
## [2,] 481   35  29  19  22  18  17  25  14   9
## [3,] 195   24  21   9   9  12   8  10   9   8
## [4,] 139   16  14  11  15   8   6   6   3   1
## [5,]  93   21  12   3   6   7   2   4   5   7
## [6,]  52   13   7   6   4   2   3   1   1   8
## [7,]  49   10  12  13   5   4   3   3   3   4
## [8,]  27   12   5  10   1   5   0   6   0   4
```

##	[9,]	14	6	6	7	7	5	2	7	2	2
##	[10,]	22	12	9	4	1	6	2	6	4	3

5 Running Jaatha

Jaatha is divided into two parts. First we find good starting positions by simulating very coarsely across the entire parameter space. We call this part *initial search*. Afterwards a more thorough *refined search* is performed starting from the best positions of the first step. Before starting the search, we need to set some options like our demographic model, the summary statistics of the real data, and a seed to ensure reproducibility:

```
jaatha <- Jaatha.initialize(dm, jsfs = jsfs, seed = 12345)
```

For more options refer to `?Jaatha.initialize` or the Jaatha manual.

5.1 The Initial Search

For the initial search, we divide the parameter space into equally-sized blocks by dividing each of the n parameters ranges into `blocks.per.par` intervals such that we obtain $(\text{blocks.per.par})^n$ blocks. Within each block we simulate `sim` data sets with – on a logarithmic scale – uniformly drawn parameter values within each block. To ensure a better sampling of the edges, we simulate in addition data sets for all corner points of each parameter block.

For these data sets we then fit the GLMs and estimate the parameter combination with the maximal score³. Each of the blocks provides a single best parameter combination.

In R, the initial search is performed with the command

```
jaatha <- Jaatha.initialSearch(jaatha, sim = 100, blocks.per.par = 2)

## *** Searching starting positions ***
## Creating initial blocks ...
## *** Block 1 : 0.1-0.707 x 0.01-0.224 x 1-4.472
## Best parameters 0.195 0.224 4.472 with estimated log-likelihood -1617
##
## *** Block 2 : 0.1-0.707 x 0.01-0.224 x 4.472-20
## Best parameters 0.169 0.224 13.57 with estimated log-likelihood -199.2
##
## *** Block 3 : 0.1-0.707 x 0.224-5 x 1-4.472
## Best parameters 0.707 0.471 4.472 with estimated log-likelihood -1064
##
## *** Block 4 : 0.1-0.707 x 0.224-5 x 4.472-20
## Best parameters 0.707 0.985 10.4 with estimated log-likelihood -127.7
##
## *** Block 5 : 0.707-5 x 0.01-0.224 x 1-4.472
```

³In this phase, Jaatha uses a score instead of the likelihood for computational reasons. The likelihood is proportional to $\exp(\text{score})$. The higher the score, the higher the likelihood.

```
## Best parameters 1.631 0.224 4.472 with estimated log-likelihood -1335
##
## *** Block 6 : 0.707-5 x 0.01-0.224 x 4.472-20
## Best parameters 0.707 0.224 9.368 with estimated log-likelihood -739.5
##
## *** Block 7 : 0.707-5 x 0.224-5 x 1-4.472
## Best parameters 5 0.781 4.472 with estimated log-likelihood -637.3
##
## *** Block 8 : 0.707-5 x 0.224-5 x 4.472-20
## Best parameters 1.571 0.925 9.641 with estimated log-likelihood -105
##
##      log.likelihood    tau      M theta
## [1,]          -105.0 0.704 0.729 0.756
## [2,]          -127.7 0.500 0.739 0.782
## [3,]          -199.2 0.135 0.500 0.870
## [4,]          -637.3 1.000 0.701 0.500
## [5,]          -739.5 0.500 0.500 0.747
## [6,]         -1064.0 0.500 0.620 0.500
## [7,]         -1334.7 0.714 0.500 0.500
## [8,]         -1617.0 0.170 0.500 0.500
```

To visualise the estimates for good starting positions sorted by score, type:

```
Jaatha.getStartingPoints(jaatha)

##      log.likelihood    tau      M theta
## [1,]          -105.0 0.704 0.729 0.756
## [2,]          -127.7 0.500 0.739 0.782
## [3,]          -199.2 0.135 0.500 0.870
## [4,]          -637.3 1.000 0.701 0.500
## [5,]          -739.5 0.500 0.500 0.747
## [6,]         -1064.0 0.500 0.620 0.500
## [7,]         -1334.7 0.714 0.500 0.500
## [8,]         -1617.0 0.170 0.500 0.500
```

Here, there is a big reduction in the scores after the first seven blocks, and a smaller one after the first two. This is suggesting that we either use the first two or the first seven blocks as starting positions for the refined search, depending on how much time we want to spend. For now, we will just use the first two points.

5.2 The Refined Search

Now we can conduct the more thorough refined search described above to improve the likelihood approximations.

```
jaatha <- Jaatha.refinedSearch(jaatha, best.start.pos = 2, sim = 100)

## *** Search with starting Point in Block 1 of 2 ***
## -----
```

```

## Step No 1
## Using 108 Simulations
## Best parameters 1.292 0.922 9.115 with estimated log-likelihood -90.8
##
## -----
## Step No 2
## Using 154 Simulations
## Best parameters 1.109 0.935 9.357 with estimated log-likelihood -92
##
## -----
## Step No 3
## Using 177 Simulations
## Best parameters 1.166 0.947 9.324 with estimated log-likelihood -88.3
## No significant score changes in the last 1 Step(s)
##
## -----
## Step No 4
## Using 256 Simulations
## Best parameters 1.186 0.955 9.323 with estimated log-likelihood -90.89
## No significant score changes in the last 2 Step(s)
##
## -----
## Step No 5
## Using 344 Simulations
## Best parameters 1.085 0.957 9.463 with estimated log-likelihood -88.93
## No significant score changes in the last 3 Step(s)
##
## -----
## Step No 6
## Using 380 Simulations
## Best parameters 1.085 0.939 9.449 with estimated log-likelihood -90.81
## No significant score changes in the last 4 Step(s)
##
## -----
## Step No 7
## Using 473 Simulations
## Best parameters 1.211 0.949 9.221 with estimated log-likelihood -90.58
##
## -----
## Step No 8
## Using 477 Simulations
## Best parameters 1.047 0.961 9.532 with estimated log-likelihood -90.59
##
## -----
## Step No 9
## Using 470 Simulations
## Best parameters 1.171 0.938 9.327 with estimated log-likelihood -90.18
## No significant score changes in the last 1 Step(s)
##

```

```

## -----
## Step No 10
## Using 537 Simulations
## Best parameters 1.008 0.978 9.587 with estimated log-likelihood -89.89
##
## -----
## Step No 11
## Using 517 Simulations
## Best parameters 1.182 0.954 9.302 with estimated log-likelihood -88.83
##
## -----
## Step No 12
## Using 568 Simulations
## Best parameters 1.118 0.947 9.422 with estimated log-likelihood -90.97
## No significant score changes in the last 1 Step(s)
##
## -----
## Step No 13
## Using 647 Simulations
## Best parameters 1.111 0.947 9.444 with estimated log-likelihood -92.75
## No significant score changes in the last 2 Step(s)
##
## -----
## Step No 14
## Using 738 Simulations
## Best parameters 1.196 0.935 9.34 with estimated log-likelihood -91.69
## No significant score changes in the last 3 Step(s)
##
## -----
## Step No 15
## Using 790 Simulations
## Best parameters 1.097 0.957 9.53 with estimated log-likelihood -89.36
## No significant score changes in the last 4 Step(s)
##
## -----
## Step No 16
## Using 841 Simulations
## Best parameters 1.211 0.939 9.166 with estimated log-likelihood -91.11
##
## -----
## Step No 17
## Using 849 Simulations
## Best parameters 1.119 0.958 9.406 with estimated log-likelihood -91.39
## No significant score changes in the last 1 Step(s)
##
## -----
## Step No 18
## Using 924 Simulations
## Best parameters 1.06 0.945 9.472 with estimated log-likelihood -90.34

```



```

## No significant score changes in the last 2 Step(s)
##
## -----
## Step No 19
## Using 952 Simulations
## Best parameters 1.066 0.964 9.506 with estimated log-likelihood -89.23
## No significant score changes in the last 3 Step(s)
##
## -----
## Step No 20
## Using 1041 Simulations
## Best parameters 1.151 0.969 9.378 with estimated log-likelihood -90.03
## No significant score changes in the last 4 Step(s)
##
## -----
## Step No 21
## Using 1080 Simulations
## Best parameters 1.095 0.966 9.438 with estimated log-likelihood -90.06
## No significant score changes in the last 5 Step(s)
##
## *** Finished search ***
## Score has not change much in the last 5 steps.
## Seems we have converged.
##
## Calculating log-composite-likelihoods for best estimates:
## * Parameter combination 1 of 10
## * Parameter combination 2 of 10
## * Parameter combination 3 of 10
## * Parameter combination 4 of 10
## * Parameter combination 5 of 10
## * Parameter combination 6 of 10
## * Parameter combination 7 of 10
## * Parameter combination 8 of 10
## * Parameter combination 9 of 10
## * Parameter combination 10 of 10
##
## *** Search with starting Point in Block 2 of 2 ****
## -----
## Step No 1
## Using 108 Simulations
## Best parameters 0.86 1.006 9.963 with estimated log-likelihood -92.33
##
## -----
## Step No 2
## Using 153 Simulations
## Best parameters 1.046 0.947 9.631 with estimated log-likelihood -89.09
##
## -----
## Step No 3

```

```

## Using 149 Simulations
## Best parameters 1.23 0.962 9.211 with estimated log-likelihood -89.12
##
## -----
## Step No 4
## Using 167 Simulations
## Best parameters 1.172 0.962 9.308 with estimated log-likelihood -89.98
## No significant score changes in the last 1 Step(s)
##
## -----
## Step No 5
## Using 258 Simulations
## Best parameters 1.076 0.947 9.492 with estimated log-likelihood -91.93
## No significant score changes in the last 2 Step(s)
##
## -----
## Step No 6
## Using 306 Simulations
## Best parameters 1.201 0.948 9.297 with estimated log-likelihood -88.42
## No significant score changes in the last 3 Step(s)
##
## -----
## Step No 7
## Using 372 Simulations
## Best parameters 1.141 0.957 9.363 with estimated log-likelihood -89.88
## No significant score changes in the last 4 Step(s)
##
## -----
## Step No 8
## Using 456 Simulations
## Best parameters 1.221 0.939 9.19 with estimated log-likelihood -89.7
## No significant score changes in the last 5 Step(s)
##
## *** Finished search ***
## Score has not change much in the last 5 steps.
## Seems we have converged.
##
## Calculating log-composite-likelihoods for best estimates:
## * Parameter combination 1 of 9
## * Parameter combination 2 of 9
## * Parameter combination 3 of 9
## * Parameter combination 4 of 9
## * Parameter combination 5 of 9
## * Parameter combination 6 of 9
## * Parameter combination 7 of 9
## * Parameter combination 8 of 9
## * Parameter combination 9 of 9
##
##

```

```
## Best log-composite-likelihood values are:
##   log.cl block   tau      M theta
## 2 -88.74      2 1.046 0.9470 9.631
## 9 -89.02      2 1.221 0.9390 9.190
## 7 -89.03      1 1.151 0.9686 9.378
## 8 -89.98      2 1.221 0.9390 9.190
## 7 -90.37      2 1.141 0.9575 9.363
```

Hence we perform two independent searches, starting in the `best.start.pos` best starting positions we choose before and according to the general options we choose during initialization, which are stored in the `jaatha` object. In each step, we build a block with `half.block.size` in each direction (on a logarithmic scale) and perform simulations for `sim` random parameter combinations within this block (plus one for very corner). We use this information to estimate the composite maximum likelihood parameters within this block and take this value as new starting position for the next step.

The algorithm stops when the score has not changed more than `epsilon` for five consecutive steps or step `max.steps` is reached. To avoid getting stuck in local maxima, the `weight` option decreases the weight of simulations of previous blocks.

Finally the log composite likelihoods for the best ten parameter combinations are approximated using `sim.final` simulations. These are values printed at the end of the search. This matrix can also be accessed via

```
likelihoods <- Jaatha.getLikelihoods(jaatha)
print(likelihoods[1:3, ])

##   log.cl block   tau      M theta
## 2 -88.74      2 1.046 0.9470 9.631
## 9 -89.02      2 1.221 0.9390 9.190
## 7 -89.03      1 1.151 0.9686 9.378
```

6 Parallelization

On Linux and OS X, Jaatha can distribute the simulations on multiple CPU cores. To use this feature, set the `cores` option during initialization. The value of the option specifies the number of cores you want to use. As default, Jaatha will execute 'packages' of 10 simulation on a core in a row to reduce the overhead created by inter-process communication. This is fine for demographic models that can be quickly simulated. However, if you have simulations that takes multiple seconds and/or many cores available, you may want smaller packages. This can be achieved by setting `sim.package.size = 1`.

The use of the `cores` option does affect the seeding system. A run without `cores` will give different results as the identical run with the option set at a value greater than one, even when using the same seed. However, if you use `cores`, the actual number of cores you use does not change the results.

7 Finite sites models

As described in ?, you can use finite sites mutation models in Jaatha. However, this currently only works on Linux and OS X and requires that Seq-Gen [Rambaut and Grassly, 1997] is installed on your system. On Debian GNU/Linux and it's derivatives (Ubuntu, Mint) you can install it running `apt-get install seq-gen` as root. Otherwise download the current version from <http://tree.bio.ed.ac.uk/software/seqgen> and compile it according to the instruction within the package. Jaatha will search for the Seq-Gen executable using your PATH variable. If you get an error that it is not able to find it, you can specify the path to the executable using the `Jaatha.setSeqgenExecutable` function.

To create a finite sites model, you must specify a mutation model using the `dm.setMutationModel` function and add an outgroup to your model using `dm.addOutgroup`. Optionally, you can then add a mutation rate heterogeneity to your model using `dm.addMutationRateHeterogeneity`.

References

- Richard R. Hudson. Generating samples under a wright-fisher neutral model of genetic variation. *Bioinformatics*, 18(2):337–338, February 2002. doi: 10.1093/bioinformatics/18.2.337.
- Lisha A. Mathew, Paul R. Staab, Laura E. Rose, and Dirk Metzler. Why to account for finite sites in population genetic studies and how to do this with jaatha 2.0. *Ecology and Evolution*, 2013. doi: 10.1002/ece3.722. URL <http://onlinelibrary.wiley.com/doi/10.1002/ece3.722/abstract>.
- Lisha Naduvilezhath, Laura E Rose, and Dirk Metzler. Jaatha: a fast composite likelihood approach to estimate demographic parameters. *Molecular Ecology*, 20(13):2709–2723, July 2011. doi: 10.1111/j.1365-294X.2011.05131.x.
- Andrew Rambaut and Nicholas C Grassly. Seq-gen: An application for the monte carlo simulation of DNA sequence evolution along phylogenetic trees. *Computer applications in the biosciences : CABIOS*, 13(3):235–238, January 1997. doi: 10.1093/bioinformatics/13.3.235. URL <http://bioinformatics.oxfordjournals.org/content/13/3/235>.