

Running Coalescent Analyses With `coalescentMCMC`

Emmanuel Paradis

March 3, 2015

Coalescent analyses have emerged in the recent years as a powerful approach to investigate the demography of populations using genetic data. The coalescent is a random process describing the coalescent times of a genealogy with respect to population size and mutation rate. In the majority of cases, the genealogy of individuals within a population is unknown. So a coalescent analysis typically integrates over the “likely” genealogies to make inference on the dynamics of the population. This uses computer-intensive methods such as Monte Carlo simulations of Markov chains. Besides, if priors are defined on the distributions of the parameters, Bayesian inference can be done. Several methods have been proposed for such integrations, although currently there is no consensus on which method is the best, or which ones are the most appropriate in some circumstances [1].

`coalescentMCMC` aims to provide a general framework to run coalescent analyses. In its current version, the package provides a simple MCMC algorithm based on Hastings’s ratio.

`coalescentMCMC` has three main groups of functions that have different roles:

- the function `coalescentMCMC` itself which runs the chain;
- some functions doing operations on tree which are called by the previous one to move from one tree to another;
- some functions to infer demography from genealogies under various coalescent models which are typically used to analyse the output of a chain run.

The motivating idea behind `coalescentMCMC` is that the user can have full control over the analysis. The options of the main function are:

```
> library(coalescentMCMC)
> args(coalescentMCMC)

function (x, ntrees = 3000, burnin = 1000, frequency = 1, tree0 = NULL,
         model = NULL, printevery = 100)
NULL
```

where `ntrees` are the number of trees to output, `burnin` is the number of trees discarded before output of trees starts, `frequency` is the sampling frequency of trees, and `tree0` is the initial tree (if not provided, a UPGMA tree from a

JC69-based distance matrix is used). `model` is either `NULL` in which case Θ is assumed to be constant, or `"time"` in which case a model where Θ follows an exponential growth is used.¹ Finally, `printevery` is an integer controlling the print frequency of the progress of the chain.

The current implementation uses only neighborhood rearrangement as proposed in [2] calling the function `NeighborhoodRearrangement` at each step of the chain. This can be modified by using other functions described in `?treeOperators`.

Let us now consider a very simple analysis with the woodmouse data available in `ape`. For the purpose of this vignette, we run a very light analysis in order to produce small outputs in a reasonable time.

```
> data(woodmouse)
> out <- coalescentMCMC(woodmouse, ntrees = 300, burnin = 100)
```

Running the Markov chain:

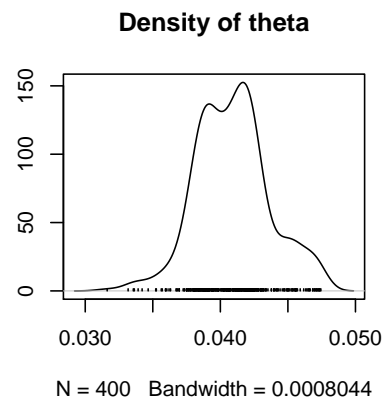
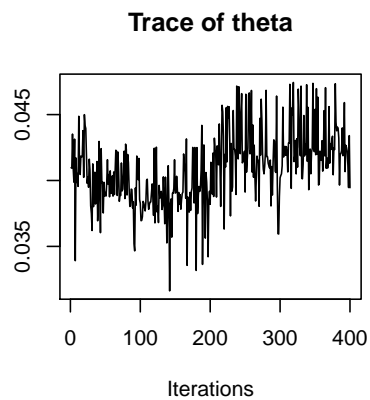
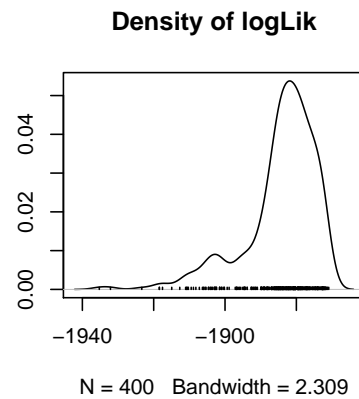
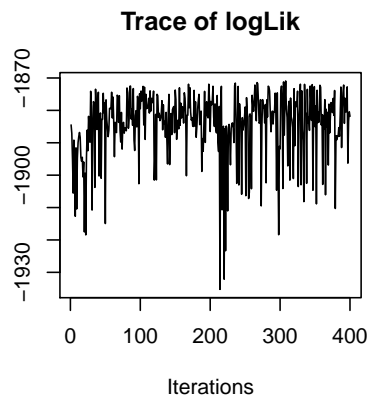
```
Number of trees to output: 300
Burn-in period: 100
Sampling frequency: 1
Number of generations to run: 400
Generation      Nb of accepted trees
    400                62
```

Done.

The output object is of class `"coda"`, so we can visualise it with the package of the same name (which has already been loaded):

```
> plot(out)
```

¹Other models will be implemented later. See the vignette "CoalescentModels".



The log-likelihood was relatively stable between -1870 and -1880 . The trees are stored in a special place of the memory (an environment in R's jargon) from where they can be retrieved with a specific function:

```
> TR <- getMCMCtrees()
> TR
```

300 phylogenetic trees

Note that the trees generated during the burn-in period are not output, but the corresponding values of log-likelihood and Θ are. Hence `out` has 400 rows.

```
> dim(out)

[1] 400  2

> colnames(out)

[1] "logLik" "theta"
```

We now run a model of time-dependent coalescent where Θ follows an exponential change through time:

```
> out2 <- coalescentMCMC(woodmouse, ntrees = 300, burnin = 100, model = "time")
```

Running the Markov chain:

Number of trees to output: 300

Burn-in period: 100

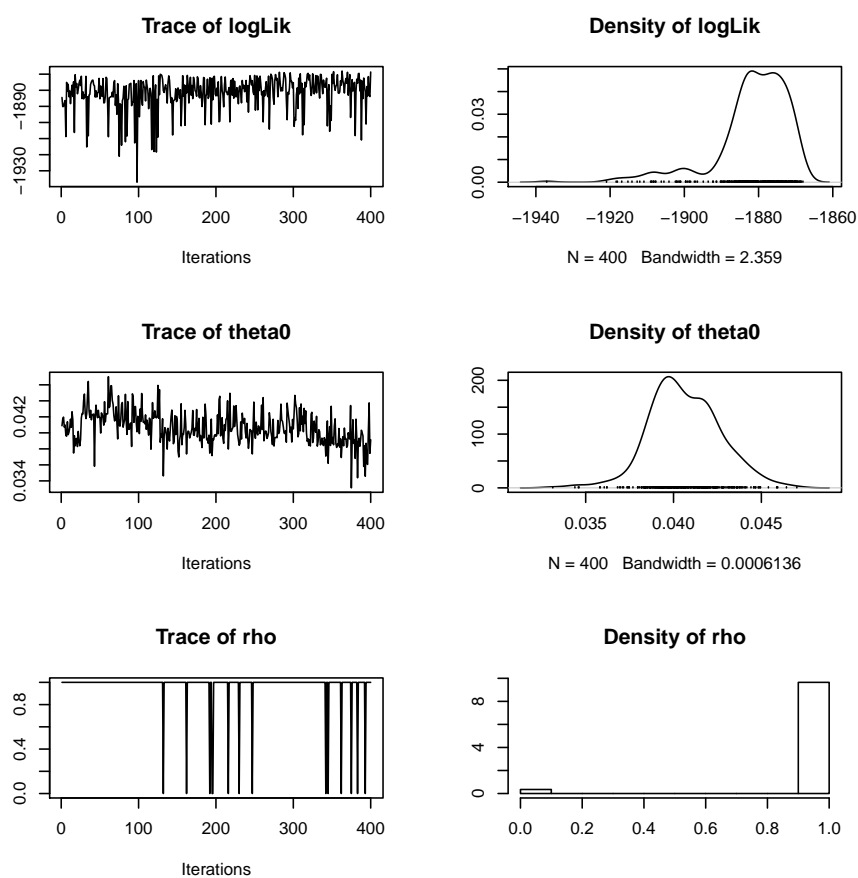
Sampling frequency: 1

Number of generations to run: 400

Generation	Nb of accepted trees
400	54

Done.

```
> plot(out2)
```



The change in log-likelihood along the chain is similar to what was observed above. The object out2 has now three columns:

```
> dim(out2)
```

```
[1] 400 3
```

```
> colnames(out2)
```

```
[1] "logLik" "theta0" "rho"
```

If we try to extract the trees as previously done and R is running in interactive mode, we will be asked which list of trees to extract:

```
> getMCMCtrees()
Several lists of MCMC trees are stored:
```

```
1 : TREES_001
2 : TREES_002
```

```
Return which number?
```

It is also possible to extract the trees of a specific chain with its number, which is useful when the R code is not run interactively (such as this vignette):

```
> TR2 <- getMCMCtrees(2)
```

The parameters of the MCMC runs are stored separately and can be extracted with:

```
> getMCMCstats()
```

```
MCMC chain summaries (chains as columns):
```

```

                001 002
Number of trees output 300 300
Burn-in period        100 100
Sampling frequency     1   1
Number of generations 400 400
Nb of accepted moves   50  60
```

We can now compare both coalescent models: the two hypotheses under consideration are:

- H_0 : Θ is constant;
- H_1 : Θ changes through time following an exponential model.

We need to calculate the likelihood under both hypotheses. This can be done with functions provided in `coalescentMCMC` (see `?dcoal`). We use the last 100 trees of each chain.² Because we are using a list of trees (which is a vector) and also a vector of estimates of Θ , we use here the function `mapply`. For clarity, we extract the trees and the values of $\hat{\Theta}$ that we need:

```
> tr <- TR[201:300]
> THETA <- out[301:400, 2]
> logLik0 <- mapply(dcoal, phy = tr, theta = THETA, log = TRUE)
> summary(logLik0)
```

```
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
72.06   73.18   73.68   73.54   73.99   74.70
```

²These calculations are quite fast, even with 1000 trees, but we use here a subset of the trees to illustrate how we select some of them which might be useful when running longer Markov chains.

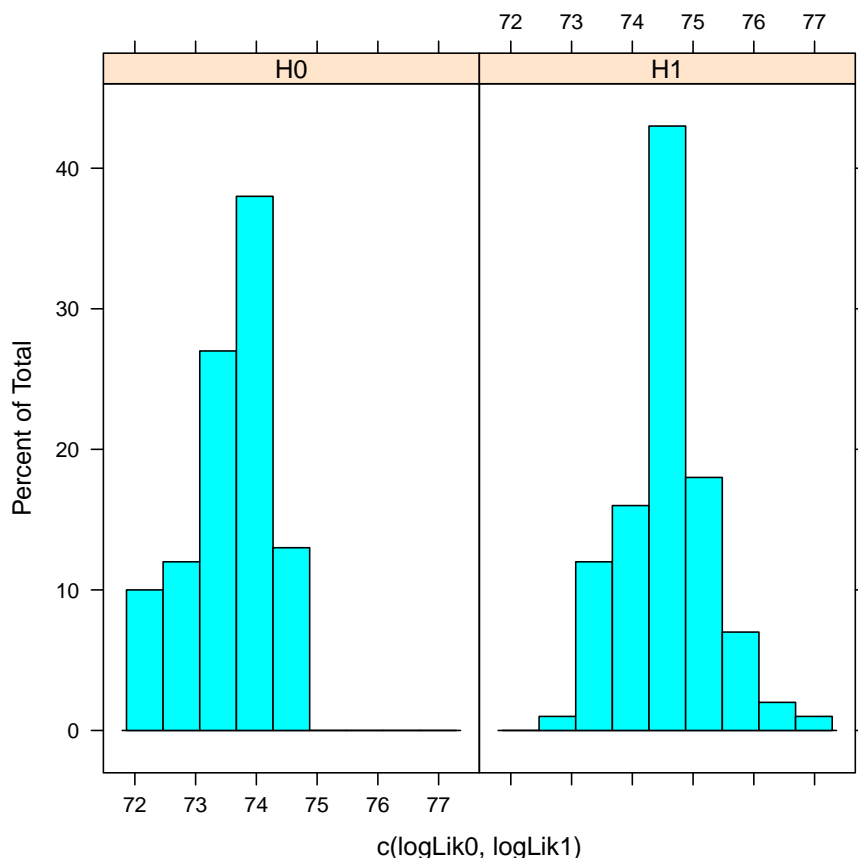
We can now repeat this operation for the second model:

```
> tr2 <- TR2[201:300]
> THETA0 <- out2[301:400, 2]
> RHO <- out2[301:400, 3]
> logLik1 <- mapply(dcoal.time, phy = tr2, theta = THETA0, rho = RHO, log = TRUE)
> summary(logLik1)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
72.84	74.17	74.63	74.55	74.93	77.09

A conditional histogram shows the two distributions:

```
> library(lattice)
> print(histogram(~c(logLik0, logLik1) | gl(2, 100, labels = c("H0", "H1"))))
```



Since the increase in log-likelihood for the second model is about 8, the LRT comparing both models is $\chi_1^2 \approx 16$ which is highly significant ($P \approx 6 \times 10^{-5}$). This suggests that the population of woodmice (*Apodemus sylvaticus*) from where these sequences have been sampled has expanded (reminding that in coalescent models the time scale is reversed so a negative value of ρ means that the population has expanded).

We can also compare both models with the deviance information criterion (DIC) developed by Spiegelhalter *et al.* [3]. These authors defined this information criterion as $DIC = \bar{D} - D(\bar{\Theta})$, where \bar{D} is the deviance³ averaged over the chain, and $D(\bar{\Theta})$ is the deviance for the mean (or estimates) of the parameters. We first produce estimates of the parameters with their means:

```
> (MLE0 <- colMeans(out[301:400, 2, drop = FALSE]))
```

```
      theta
0.04274699
```

```
> (MLE1 <- colMeans(out2[301:400, 2:3]))
```

```
      theta0      rho
0.0396134 0.9400073
```

We then calculate the second term, $D(\bar{\Theta})$, using these estimates and averaging the deviance values over the trees:

```
> D2.0 <- -2 * mean(mapply(dcoal, phy = tr, theta = MLE0, log = TRUE)) + 2
> D2.1 <- -2 * mean(mapply(dcoal.time, phy = tr2, theta = MLE1[1],
+                          rho = MLE1[2], log = TRUE)) + 4
```

The values of DIC for both models can now be calculated simply:

```
> mean(-2 * logLik0 + 2) - D2.0
```

```
[1] -0.03350816
```

```
> mean(-2 * logLik1 + 4) - D2.1
```

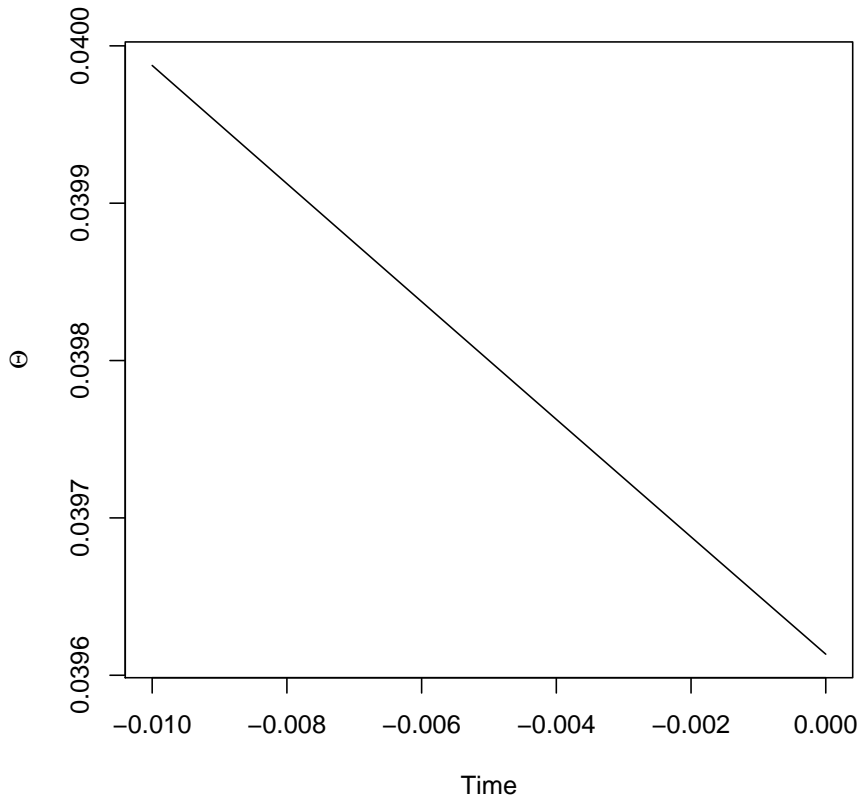
```
[1] -0.03877519
```

Like for other information criteria (AIC, BIC, ...) the model with the smallest value of DIC should be preferred: this is here the time-dependent model.

We can represent the temporal variation in Θ predicted by this model (remember that the time scale is the one of molecular change):

```
> x <- seq(0, 0.01, 0.0001)
> y <- MLE1["theta0"] * exp(MLE1["rho"] * x)
> plot(-x, y, "l", xlab = "Time", ylab = expression(Theta))
```

³The deviance is defined as $D = -2 \log L + 2k$, with L the likelihood and k the number of estimated parameters.



Other things that could be done with simple R commands include:

- Compute confidence intervals around $\hat{\Theta}_0$ and $\hat{\rho}$ (alternatively, posterior distributions of these parameters if a Bayesian sampling is done);
- Re-run the chain(s) with different initial trees, for instance to run branching chains taking a tree from TR or TR2.
- Use other other models of time-dependent coalescent (see the other vignette in this package).

References

- [1] J. Felsenstein. Trees of genes in populations. In O. Gascuel and M. Steel, editors, *Reconstructing Evolution: New Mathematical and Computational Advances*, pages 3–29. Oxford University Press, Oxford, 2007.
- [2] M. K. Kuhner, J. Yamato, and J. Felsenstein. Estimating effective population size and mutation rate from sequence data using Metropolis-Hastings sampling. *Genetics*, 140:1421–1430, 1995.
- [3] D. J. Spiegelhalter, N. G. Best, B. R. Carlin, and A. van der Linde. Bayesian measures of model complexity and fit. *64(4):583–616*, 2002.