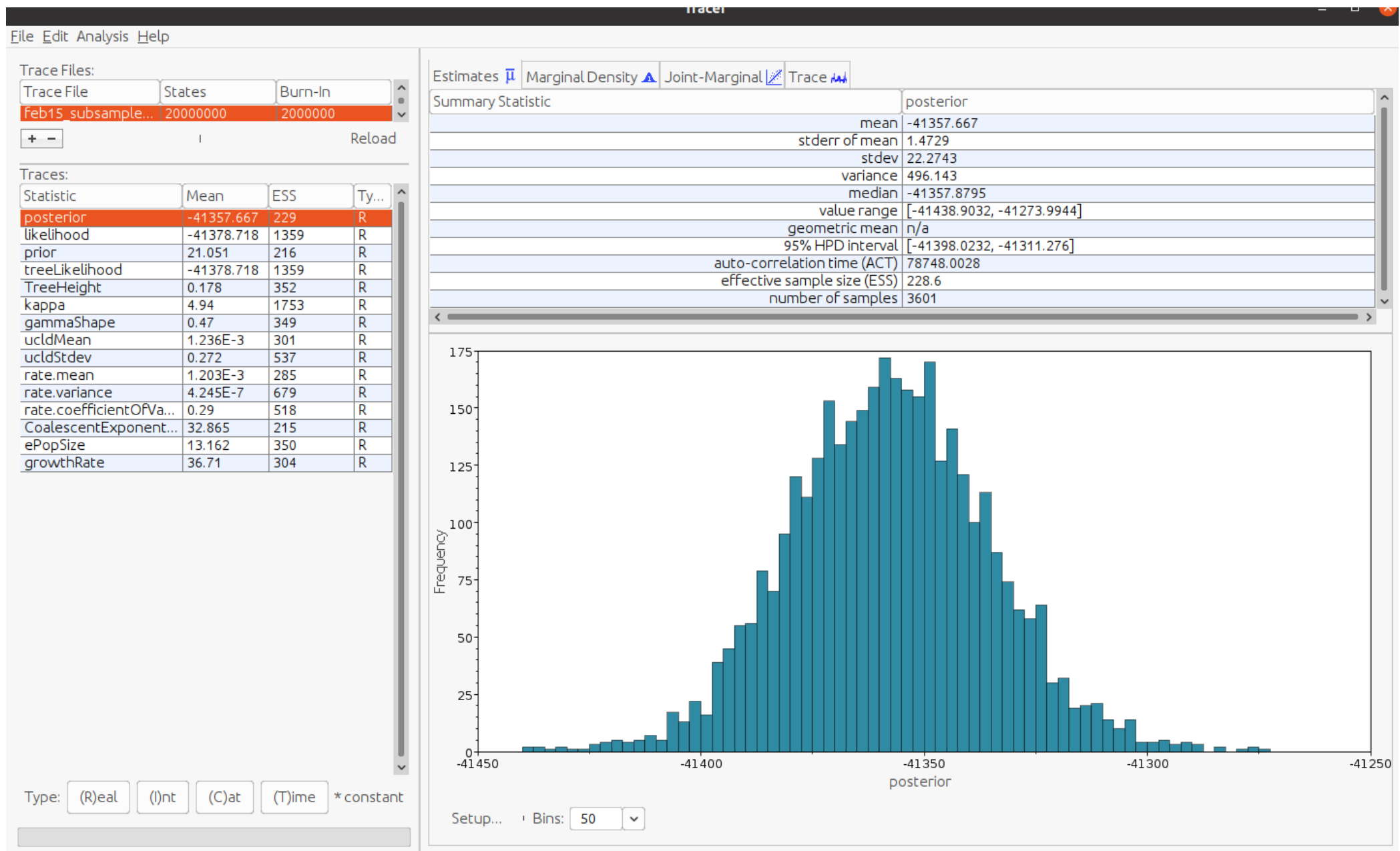


Interpreting results and
summarising trees

- Refresher of MCMC traces
- Summarising estimates
 - Means, medians, mode, credible intervals, highest posterior densities
- Hypothesis support with posterior distributions
- Key traces
- Summary trees
 - Typical summary trees
 - Summarising node heights

Refresher of MCMC traces



File Edit Analysis Help

Trace Files:

Trace File	States	Burn-In
Feb15_subsample...	20000000	2000000

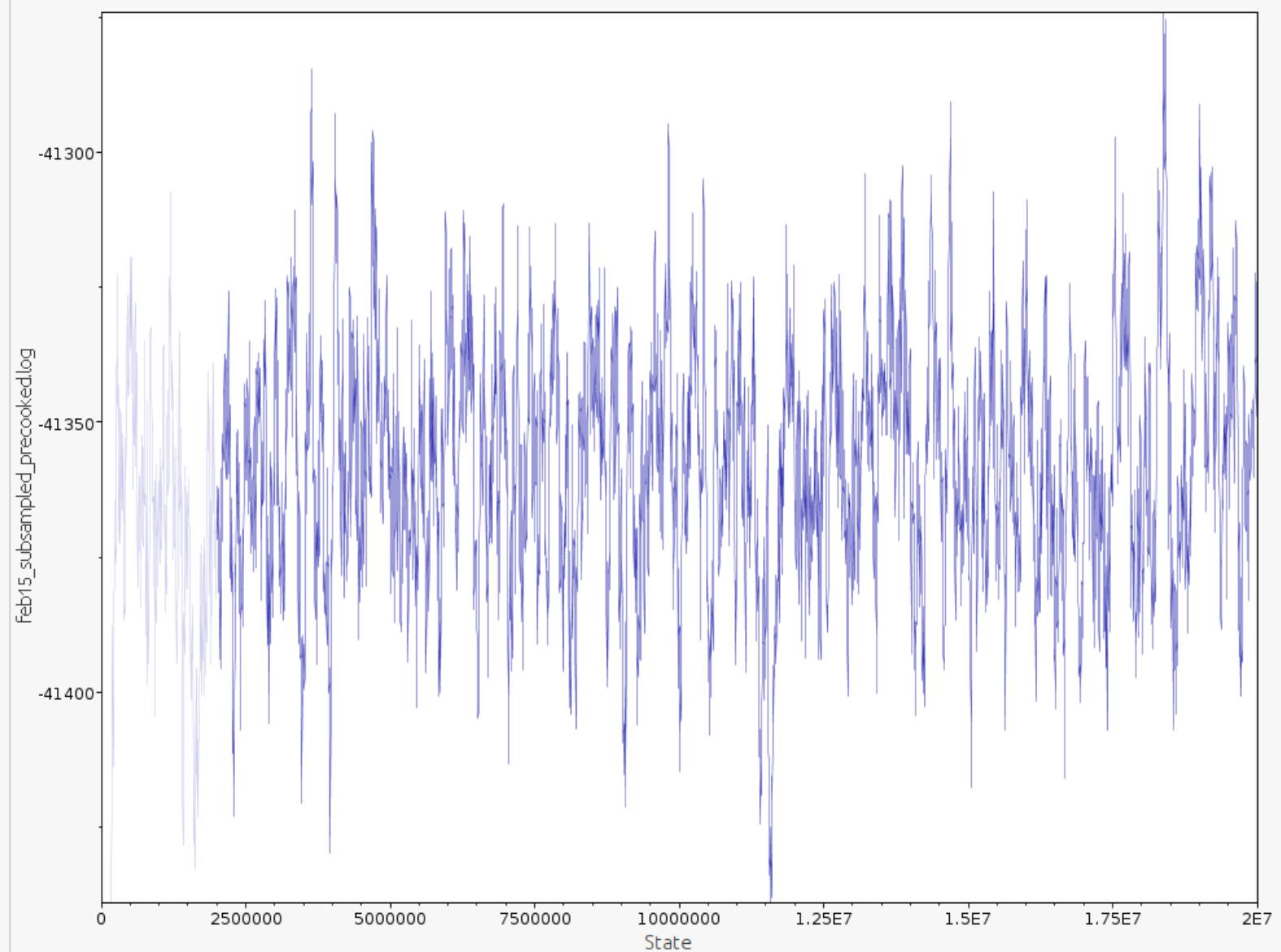
+ - | Reload

Traces:

Statistic	Mean	ESS	Ty...
posterior	-41357.667	229	R
likelihood	-41378.718	1359	R
prior	21.051	216	R
treeLikelihood	-41378.718	1359	R
TreeHeight	0.178	352	R
kappa	4.94	1753	R
gammaShape	0.47	349	R
uclMean	1.236E-3	301	R
uclStdev	0.272	537	R
rate.mean	1.203E-3	285	R
rate.variance	4.245E-7	679	R
rate.coefficientOfVa...	0.29	518	R
CoalescentExponent...	32.865	215	R
ePopSize	13.162	350	R
growthRate	36.71	304	R

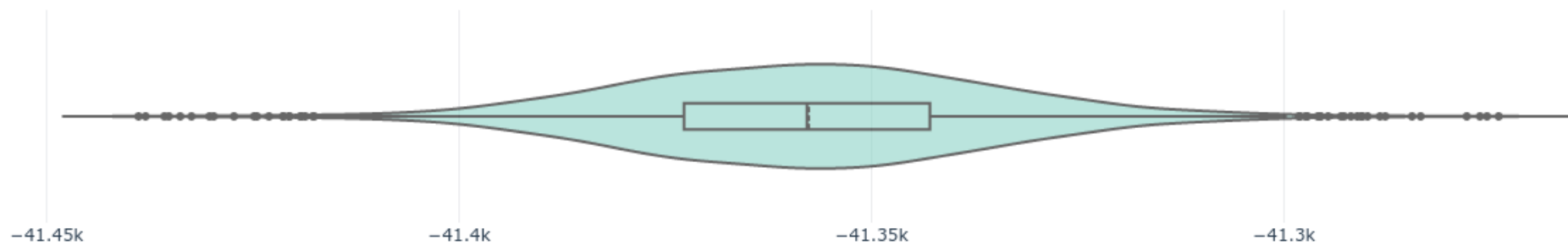
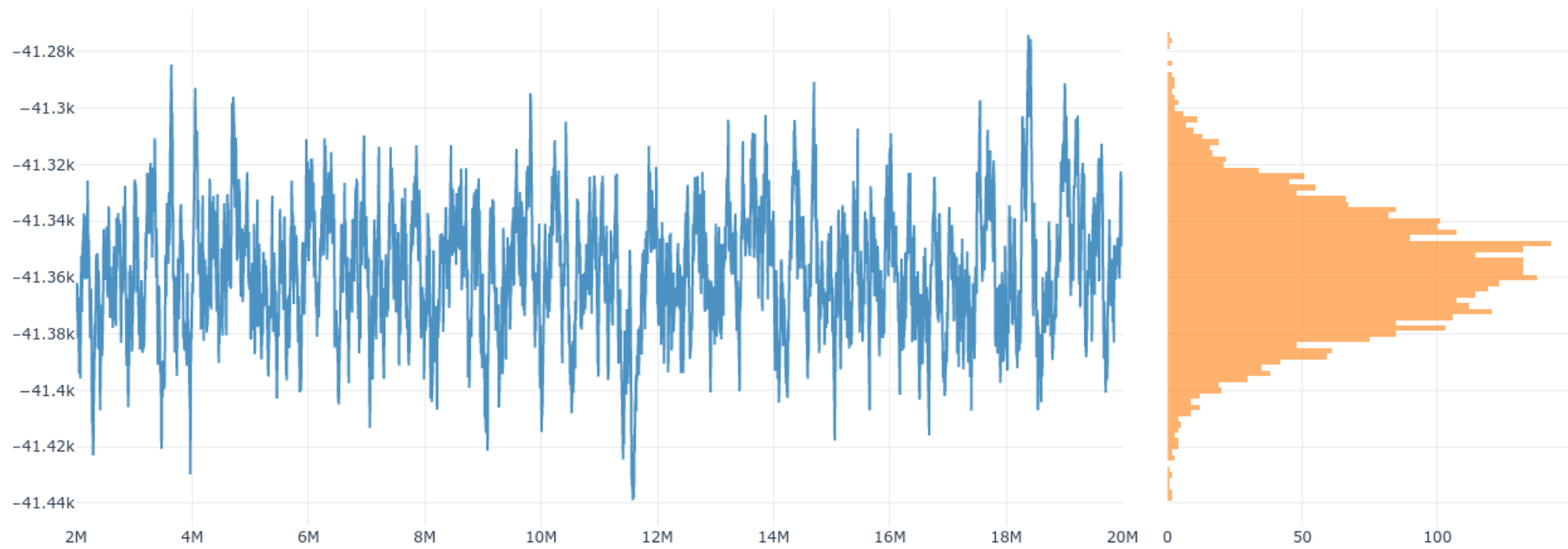
Type: (R)eal (I)nt (C)at (T)ime * constant

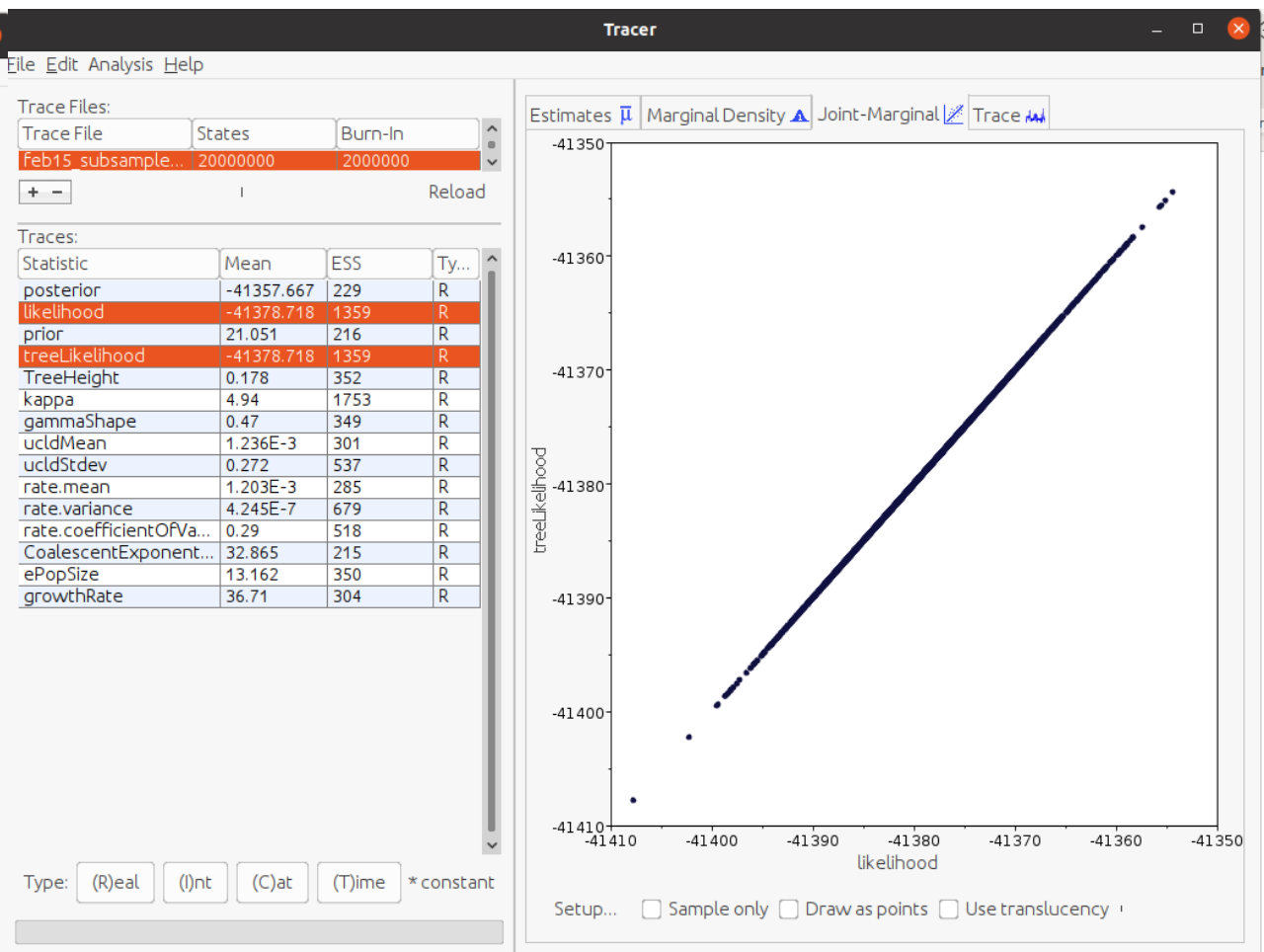
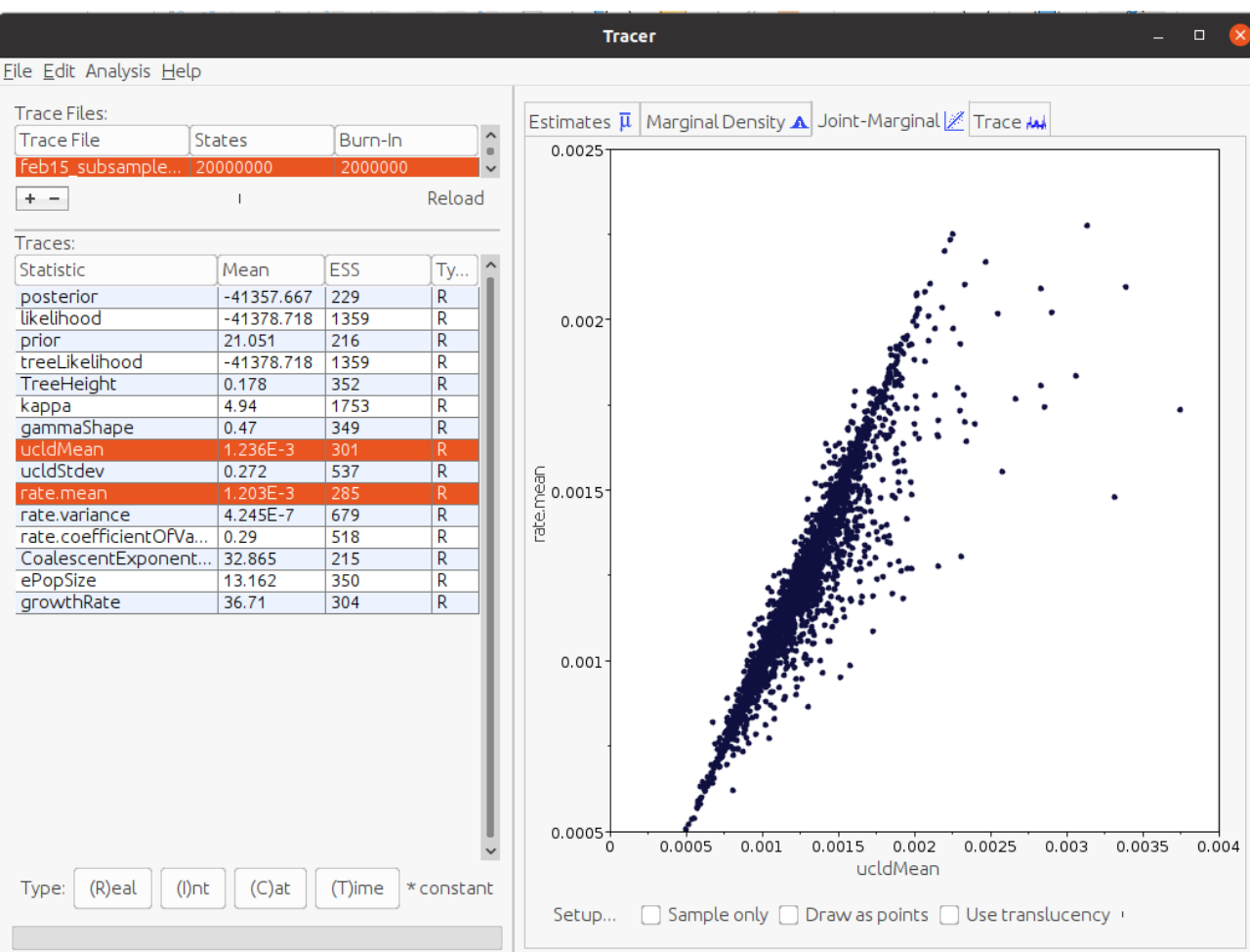
Estimates Marginal Density Joint-Marginal Trace

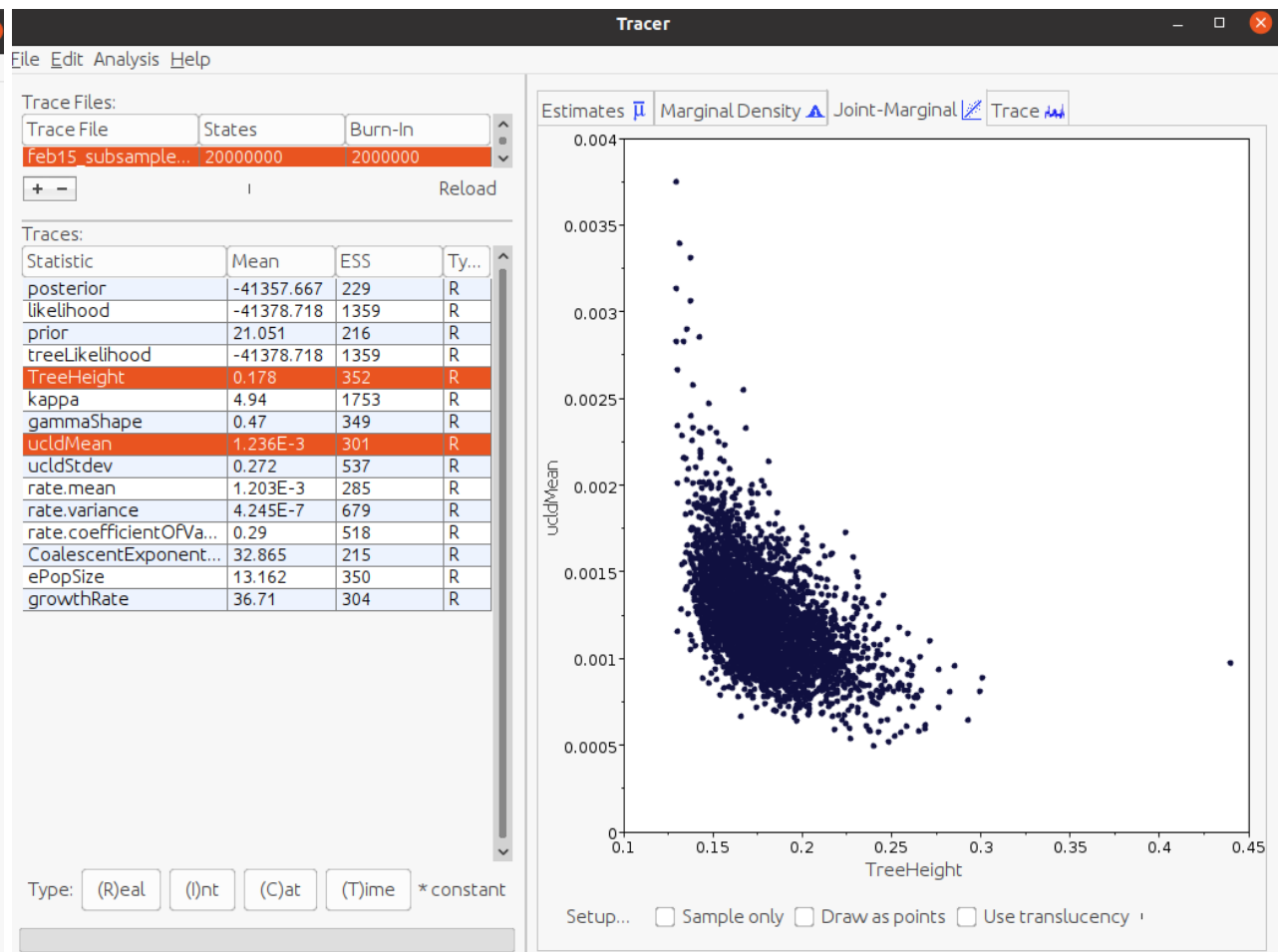
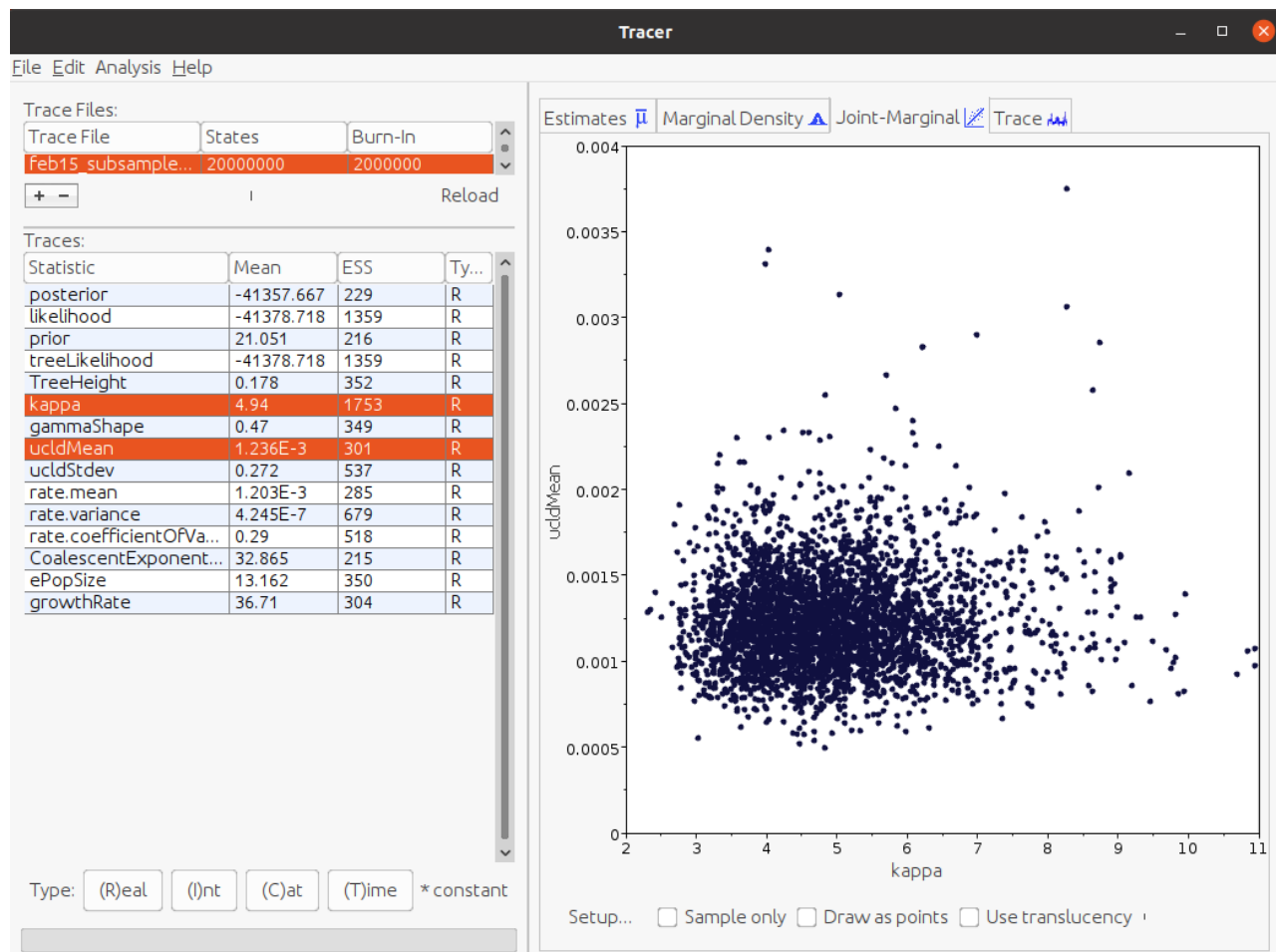
☒ Show Burn-in ☐ Sample only ☒ Draw line plot Legend: None Colour by: Trace



POSTERIOR



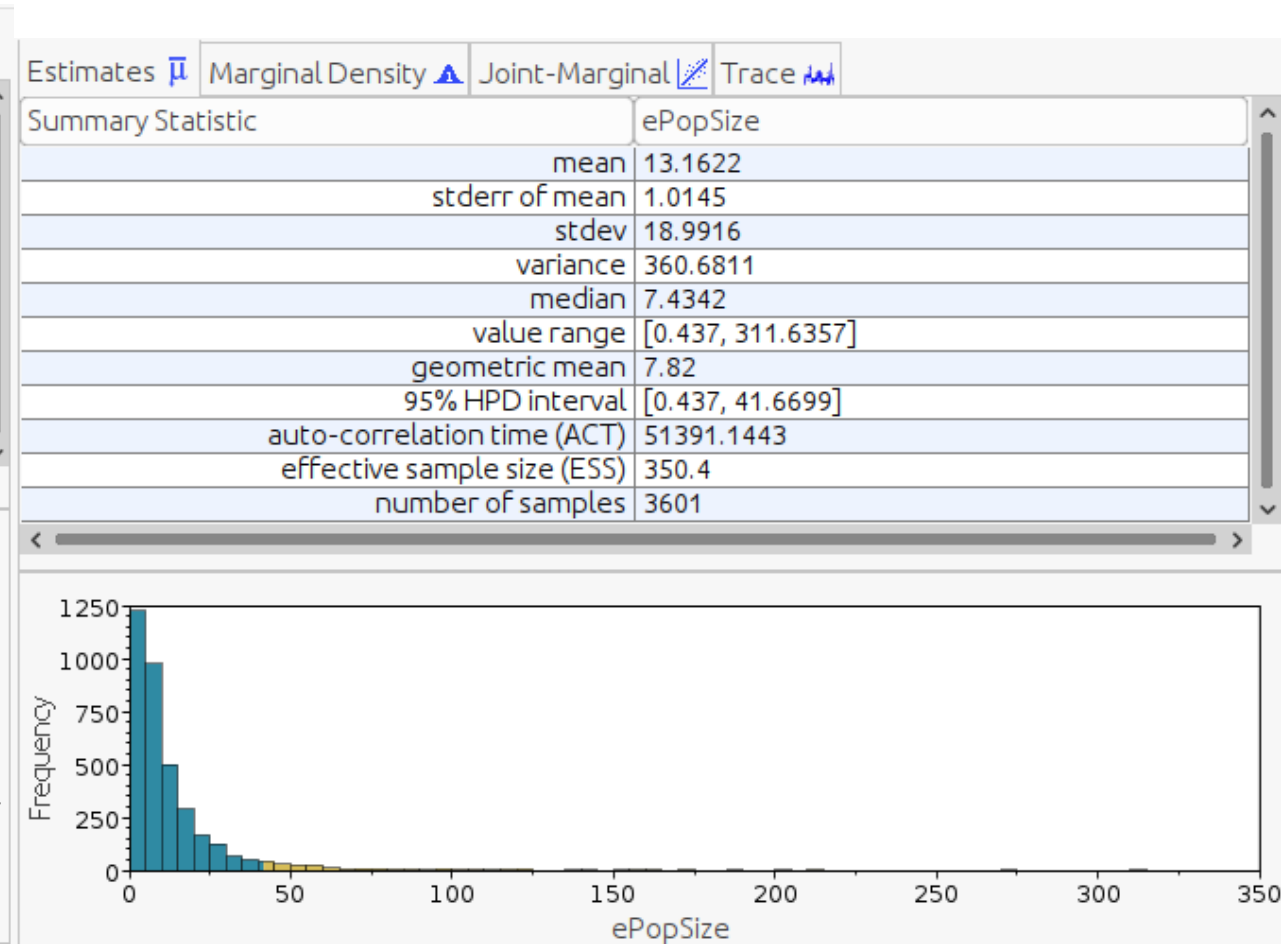
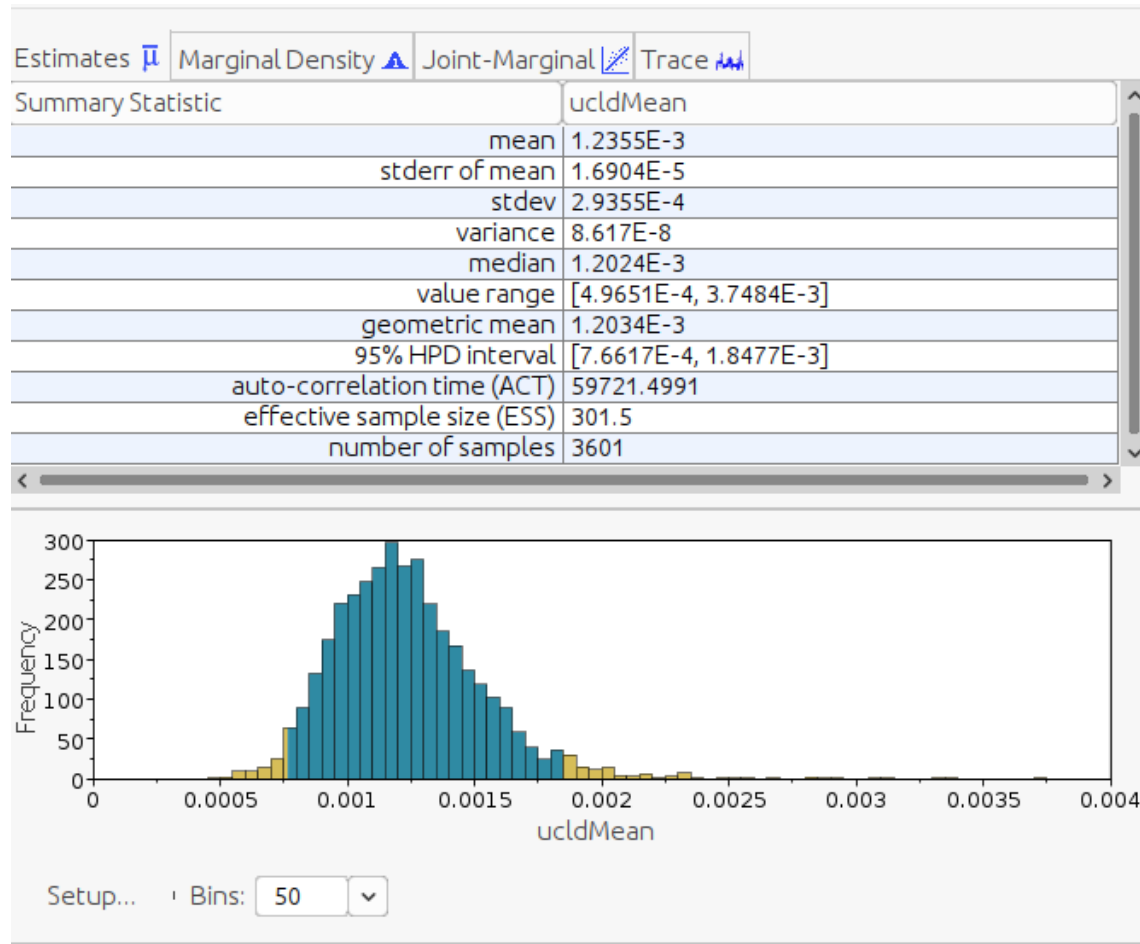




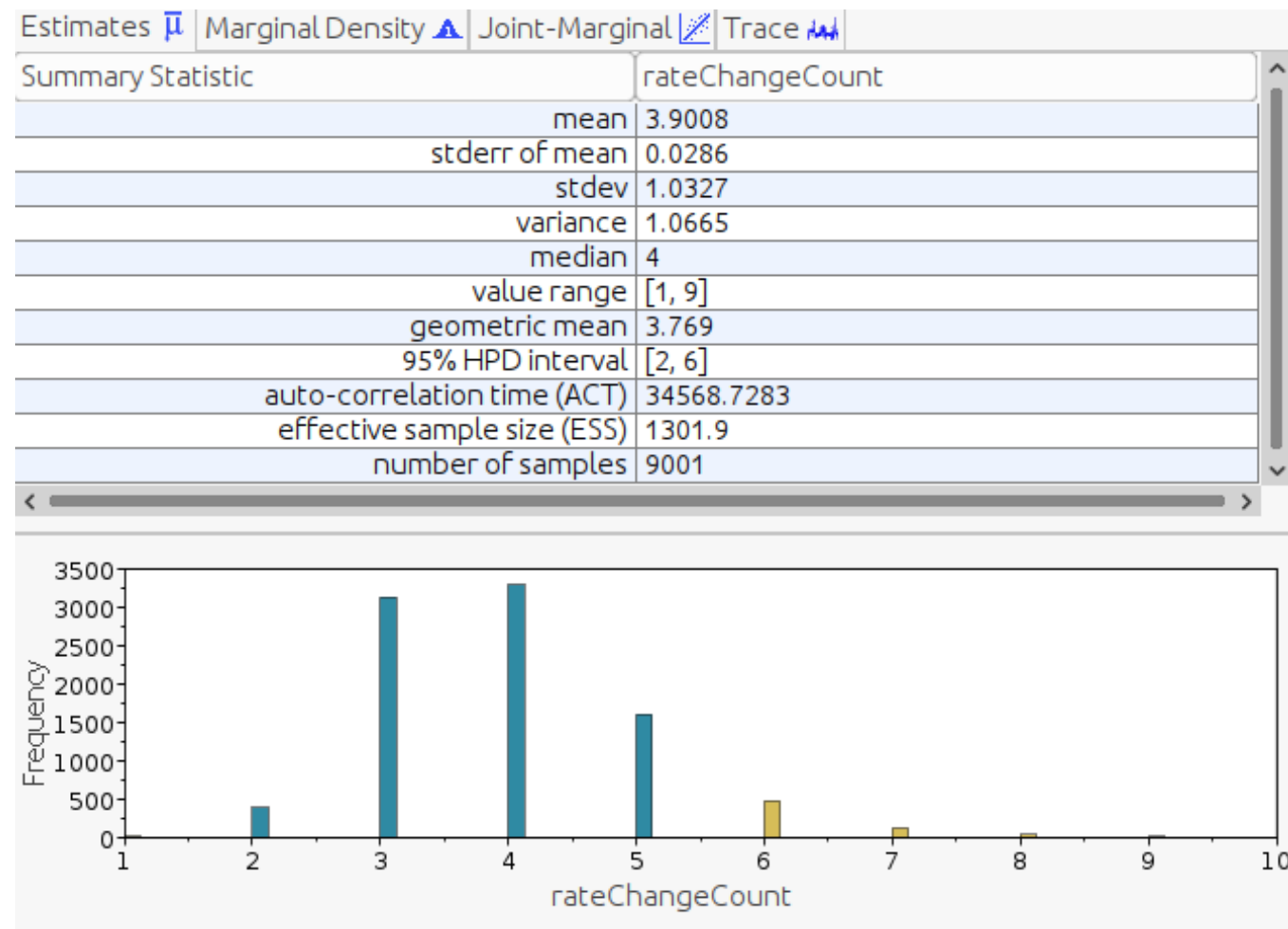
Concept summary

- MCMC sampling frequency is proportional to the posterior probability.
- Correlation between parameters is sometimes expected (but not for most parameters -> overparameterisation).

Summarising estimates



- Choice of one summary statistic (e.g. mean, median) depends on the shape of the distribution.



- For discrete statistics we can use the mode=maximum a posteriori estimate (MAP).

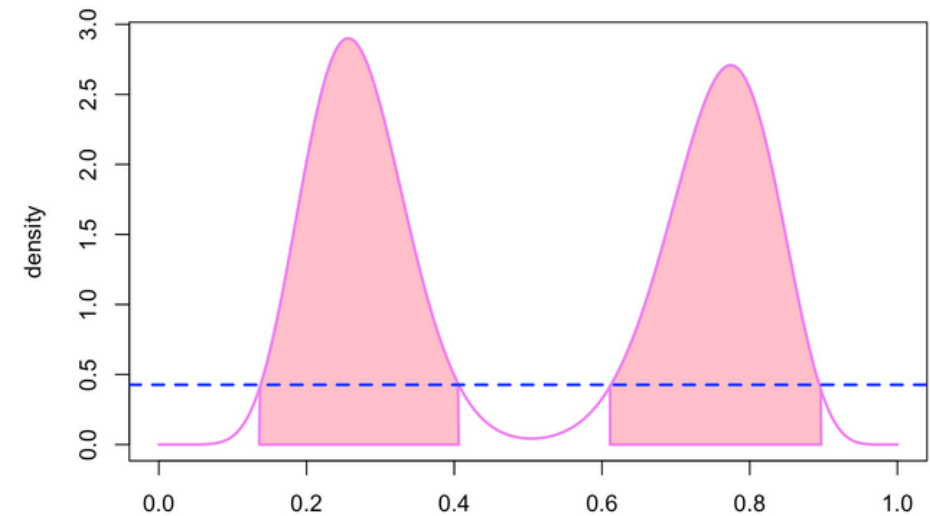
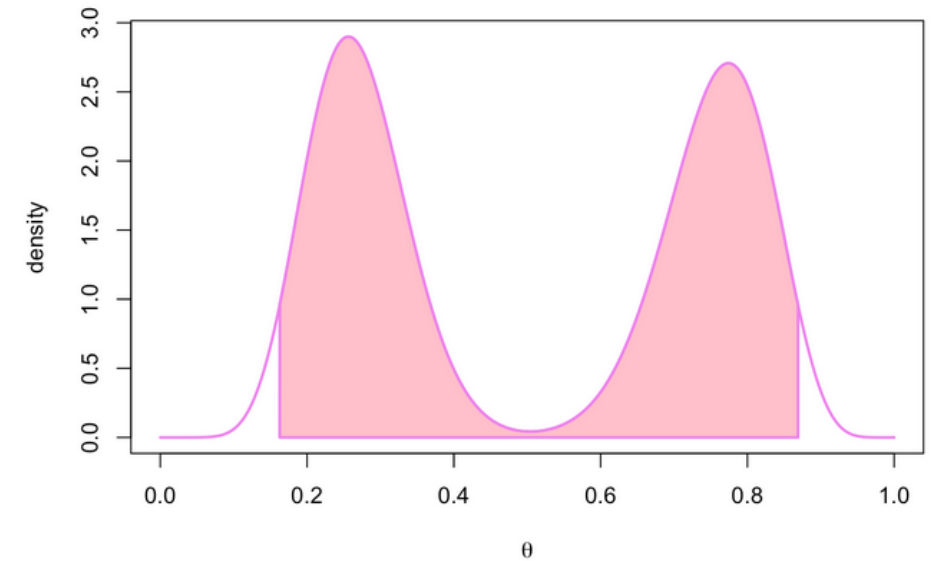
Reporting uncertainty

- Credible interval: 'given our data there is a 0.95 probability that our parameter falls in a range.'
- Highest posterior density: 'every point within the interval has a higher density than any point outside.'
- Confidence interval: 'there is a 0.95 probability that if we repeat this experiment 100 times, the confidence interval will include the true value 95 times.'

Reporting uncertainty

Credible intervals are very similar to HPDs when the distribution is unimodal. Can be calculated with quantiles or percentiles.

HPDs are very useful for bimodal distributions or those that are very skewed. They require additional calculations (in R use package HDinterval).



Concept summary

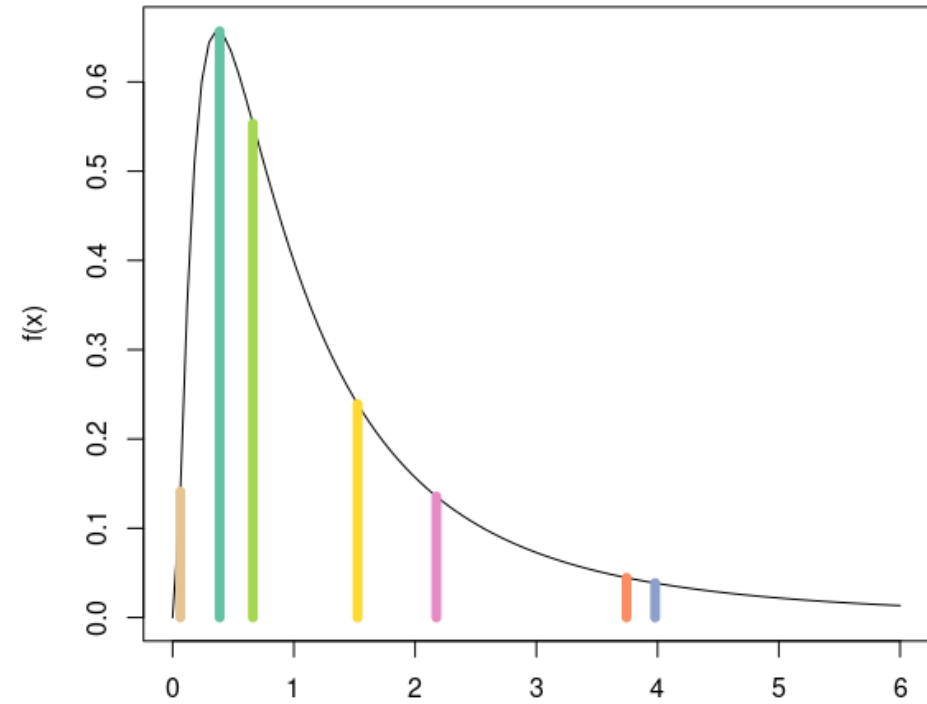
- All parameters and statistics have uncertainties.
 - Report using credible intervals or highest posterior densities NOT confidence intervals.
- The choice of summary values (e.g. mean) depends on the distribution.
- Reporting the MAP value is an alternative to the mean, particularly if the distribution is not very skewed or if it is discrete.

Hypothesis support using
posterior distributions

Bayesian hypothesis testing

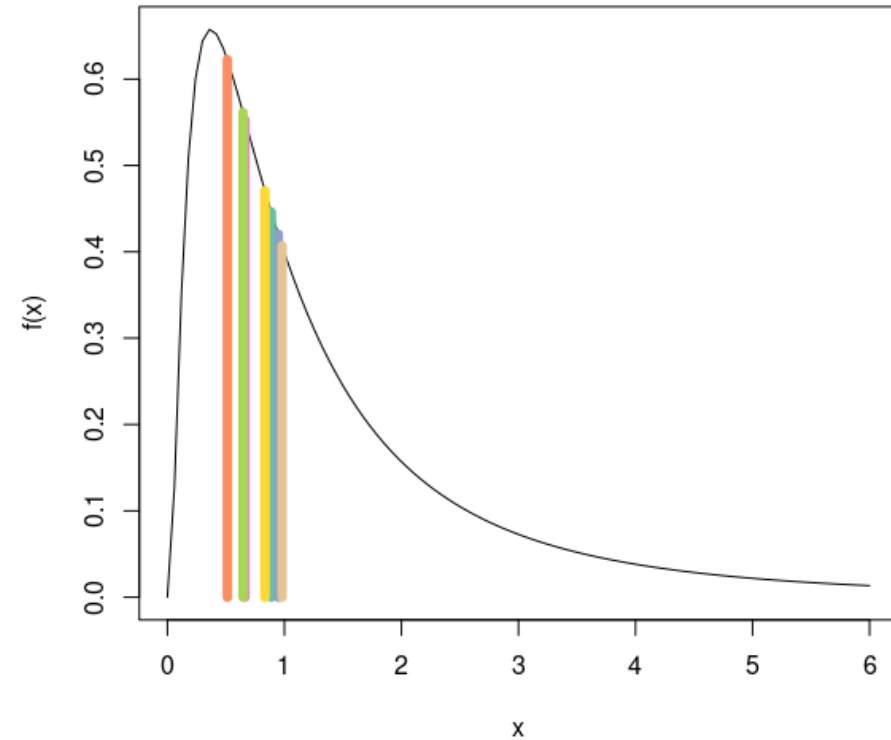
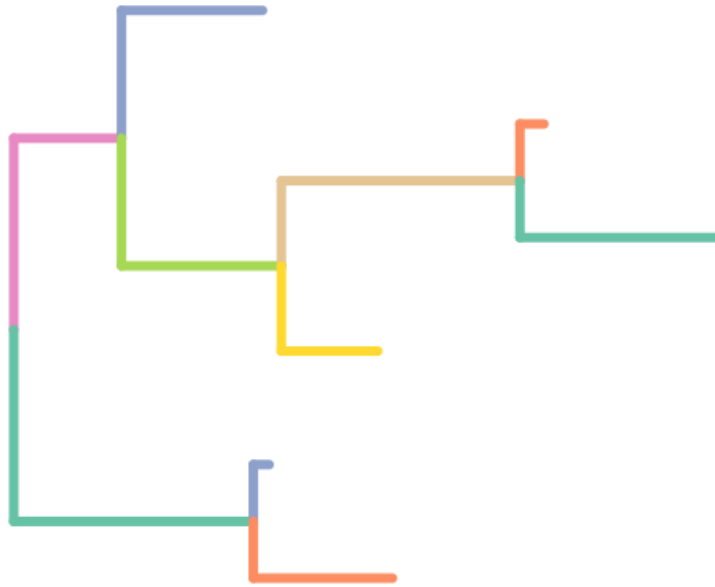
- 'Formal' methods via Bayes factors and marginal likelihoods
 - Requires additional computation, proper priors, and careful prior selection.
- Inspection of the posterior distributions
 - Straight-forward if model is correctly parameterised.
 - May require setting up complex hierarchical models.
 - Only for 'nested' models
- Posterior model checking
 - Requires simulations and data summaries.
 - Can assess the overall, not just relative, fit of models.
- Bayesian model averaging
 - Alternative to marginal likelihood calculations.
 - Traversing model space can be difficult.

We still have the
responsibility of proposing
a pool of models!



The mean of the lognormal distribution (**uclid.mean**) is 0.1 and the sd (**uclid.stdev**) is 1.^x
 The mean of the branch rates (*meanRate*) is 1.79 and the sd is 1.58

$$\begin{aligned} \text{Coefficient of rate variation} &= \text{sd of branch rates} / \text{meanRate} \\ &= 0.88 \end{aligned}$$

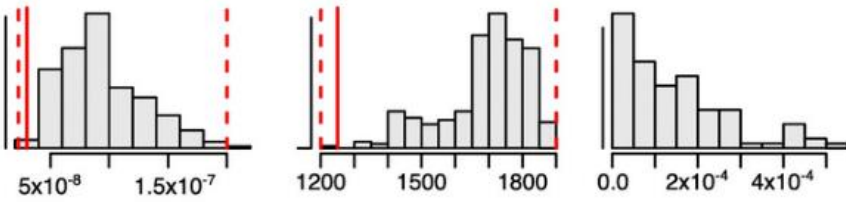


The mean of the lognormal distribution (**uclid.mean**) is 0.1 and the sd (**uclid.stdev**) is 1.
 The mean of the branch rates (*meanRate*) is 0.78 and the sd is 0.17

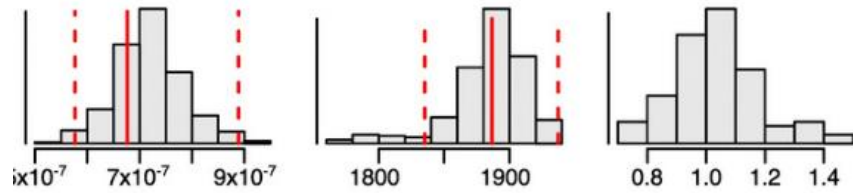
Coefficient of rate variation = *sd of branch rates / meanRate*
 = 0.22

Assume that the strict clock is a relaxed clock with very low rate variation (it is nested)!

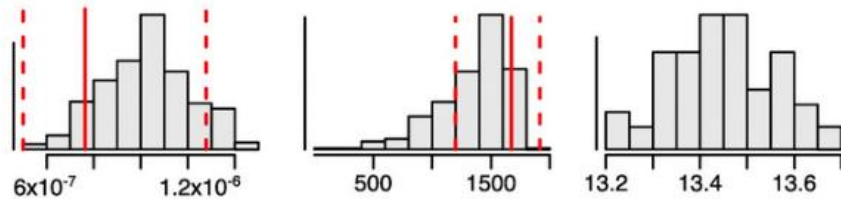
Mycobacterium tuberculosis Lineage 2



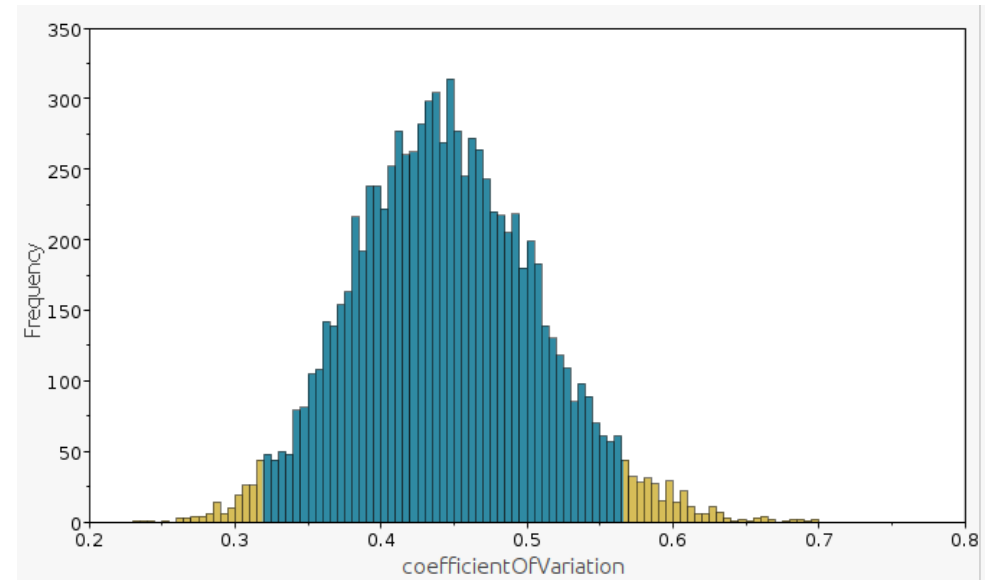
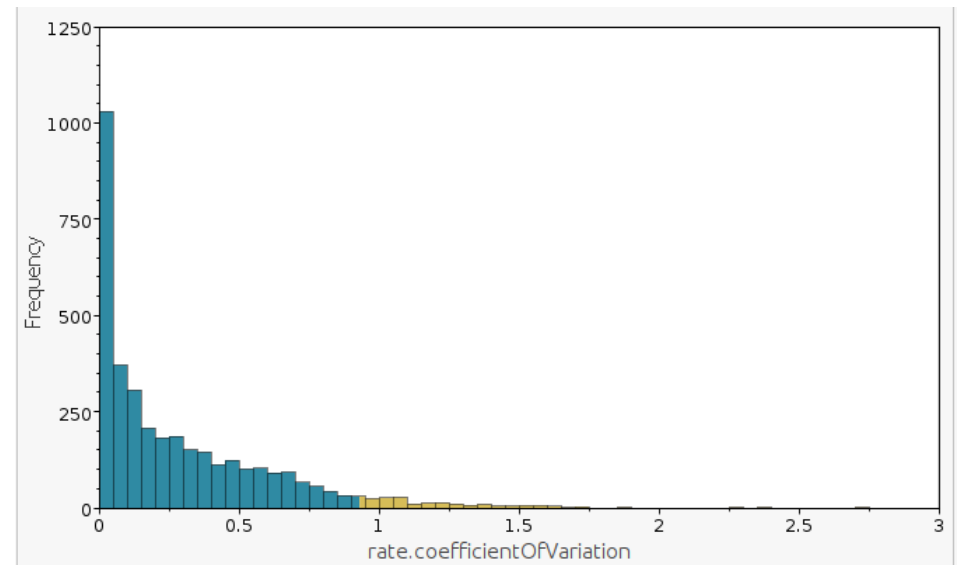
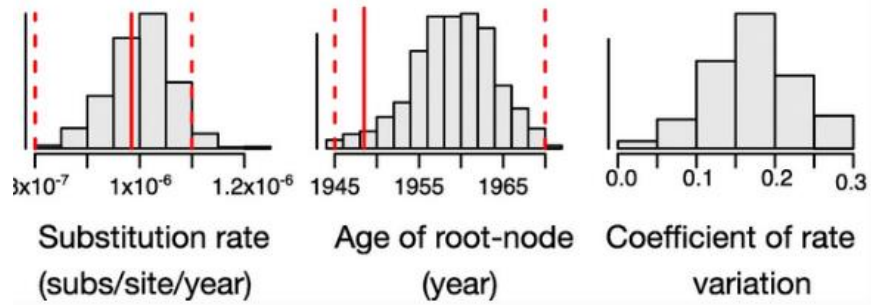
Vibrio cholerae

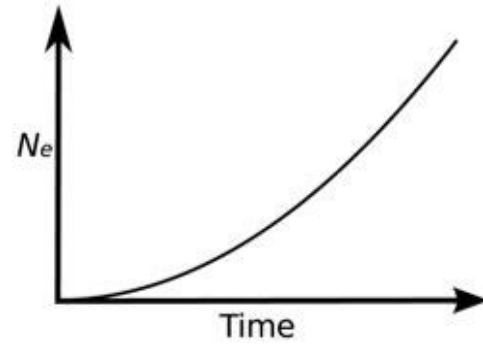
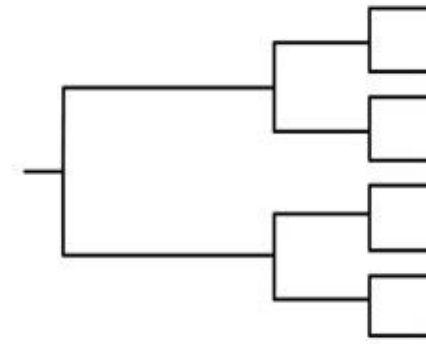
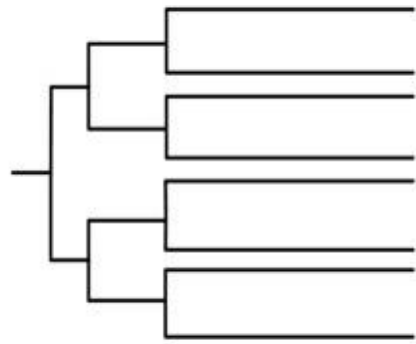


Shigella dysenteriae type I



Staphylococcus aureus ST239





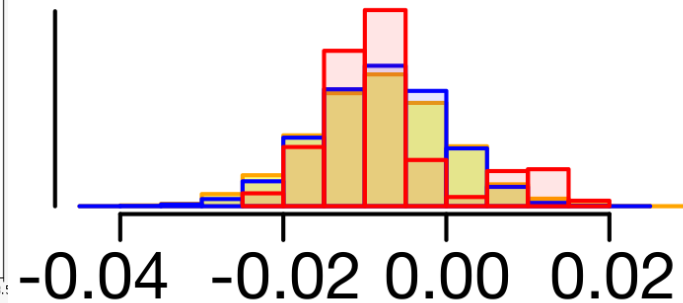
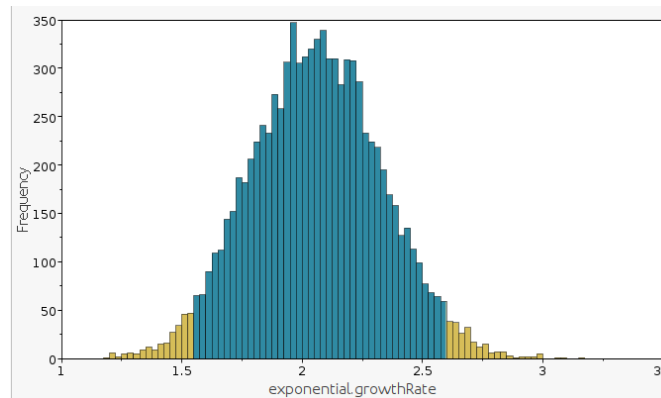
Under the exponential growth:

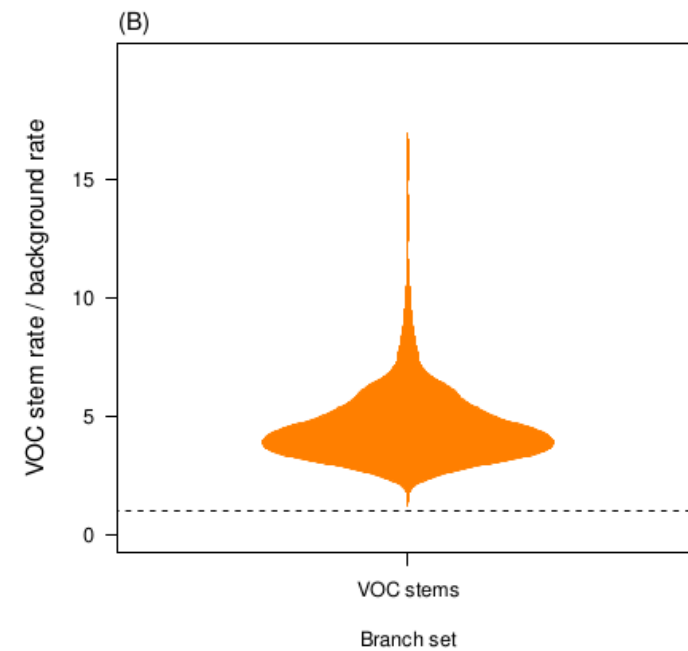
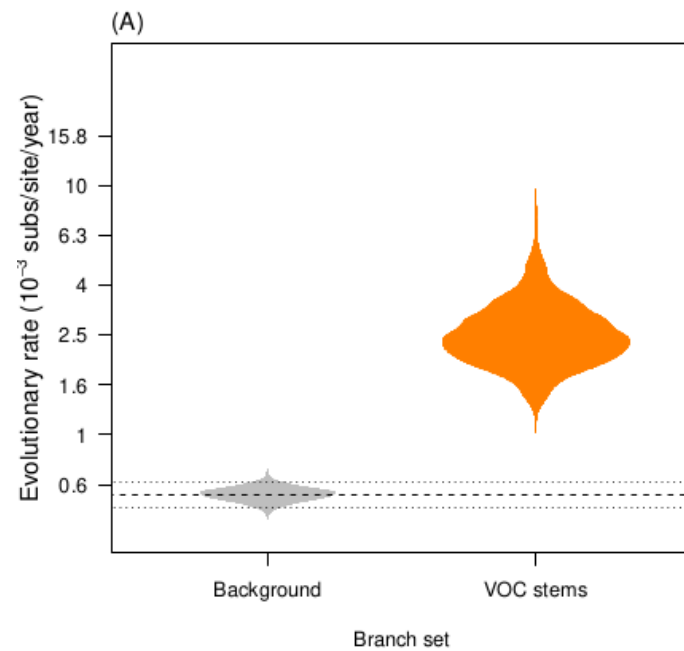
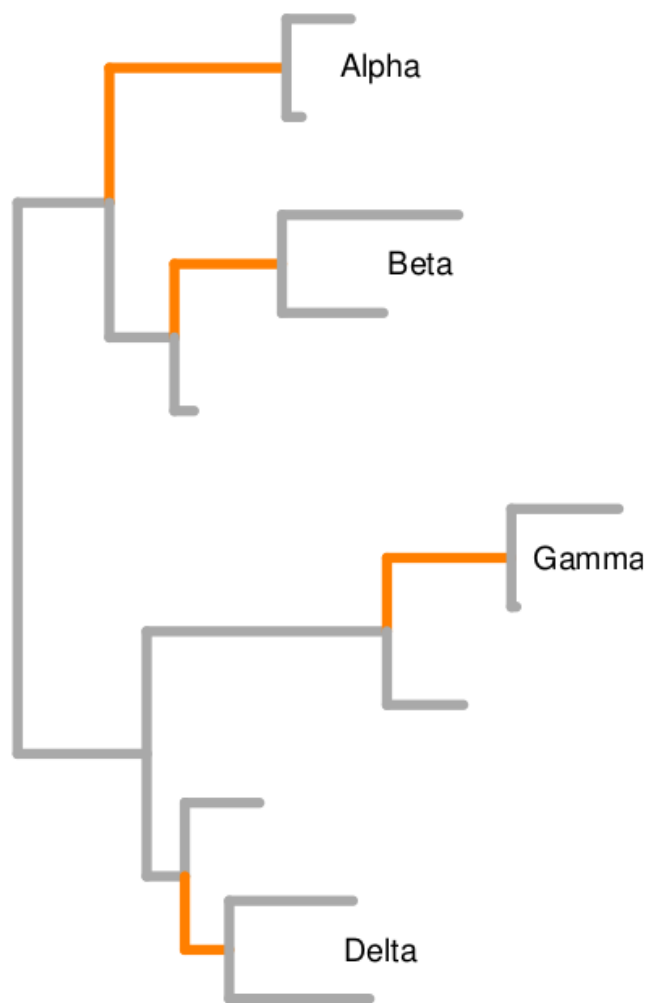
$$N_e = N_0 * e^{rt}$$

For $r > 0$

$$N_e = N_0 * e^{rt}$$

If $r = 0$
 $N_e = N_0$





Concept summary

- For some models, parameters or statistics can give an idea of how well it is supported, **relative to nested models**.

Key traces

$$P(\text{E} \text{ } \text{ } \text{ } \text{ } | \text{ }) = \frac{P(\text{ } | \text{E} \text{ } \text{ } \text{ }) P(\text{E} | \text{ }) P(\text{ }) P(\text{ }) P(\text{ })}{P(\text{ })}$$


Alignment


Chronogram


Branching model
(can be an epi model)


Substitution
model


Clock
model

- HKY+G: 2 parameters (kappa, shape)
- Coalescent exponential: 2 parameters (growthRate, ePopSize)
- UCLN clock model: 2 parameters (ucldMean, ucldStdev)*

* arguably the branch lengths are also parameters, but they are much less tractable.

Probability densities

Summary statistics

Parameters of the model *

Traces:

Statistic	Mean	ESS	Ty...
posterior	-41357.667	229	R
likelihood	-41378.718	1359	R
prior	21.051	216	R
treeLikelihood	-41378.718	1359	R
TreeHeight	0.178	352	R
kappa	4.94	1753	R
gammaShape	0.47	349	R
uclMean	1.236E-3	301	R
uclStdev	0.272	537	R
rate.mean	1.203E-3	285	R
rate.variance	4.245E-7	679	R
rate.coefficientOfVa...	0.29	518	R
CoalescentExponent...	32.865	215	R
ePopSize	13.162	350	R
growthRate	36.71	304	R

Traces:

Statistic	Mean	ESS
posterior	-41357.667	229
likelihood	-41378.718	1359
prior	21.051	216
treeLikelihood	-41378.718	1359
TreeHeight *	0.178	352
kappa	4.94	1753
gammaShape	0.47	349
uclMean	1.236E-3	301
uclStdev	0.272	537
rate.mean	1.203E-3	285
rate.variance	4.245E-7	679
rate.coefficientOfVa...	0.29	518
CoalescentExponent...	32.865	215
ePopSize	13.162	350
growthRate	36.71	304

Traces:

Statistic	Mean	ESS	Ty...
posterior	-41357.667	229	R
likelihood	-41378.718	1359	R
prior	21.051	216	R
treeLikelihood	-41378.718	1359	R
TreeHeight	0.178	352	R
kappa	4.94	1753	R
gammaShape	0.47	349	R
uclMean	1.236E-3	301	R
uclStdev	0.272	537	R
rate.mean	1.203E-3	285	R
rate.variance	4.245E-7	679	R
rate.coefficientOfVa...	0.29	518	R
CoalescentExponent...	32.865	215	R
ePopSize	13.162	350	R
growthRate	36.71	304	R

$$P(\text{Alignment} \mid \text{Chronogram} \mid \text{Branching model} \mid \text{Substitution model} \mid \text{Clock model}) = \frac{P(\text{Alignment} \mid \text{Chronogram} \mid \text{Branching model} \mid \text{Substitution model} \mid \text{Clock model}) P(\text{Branching model} \mid \text{Substitution model}) P(\text{Substitution model}) P(\text{Clock model})}{P(\text{Alignment})}$$


Alignment


Chronogram


Branching model
(can be an epi model)

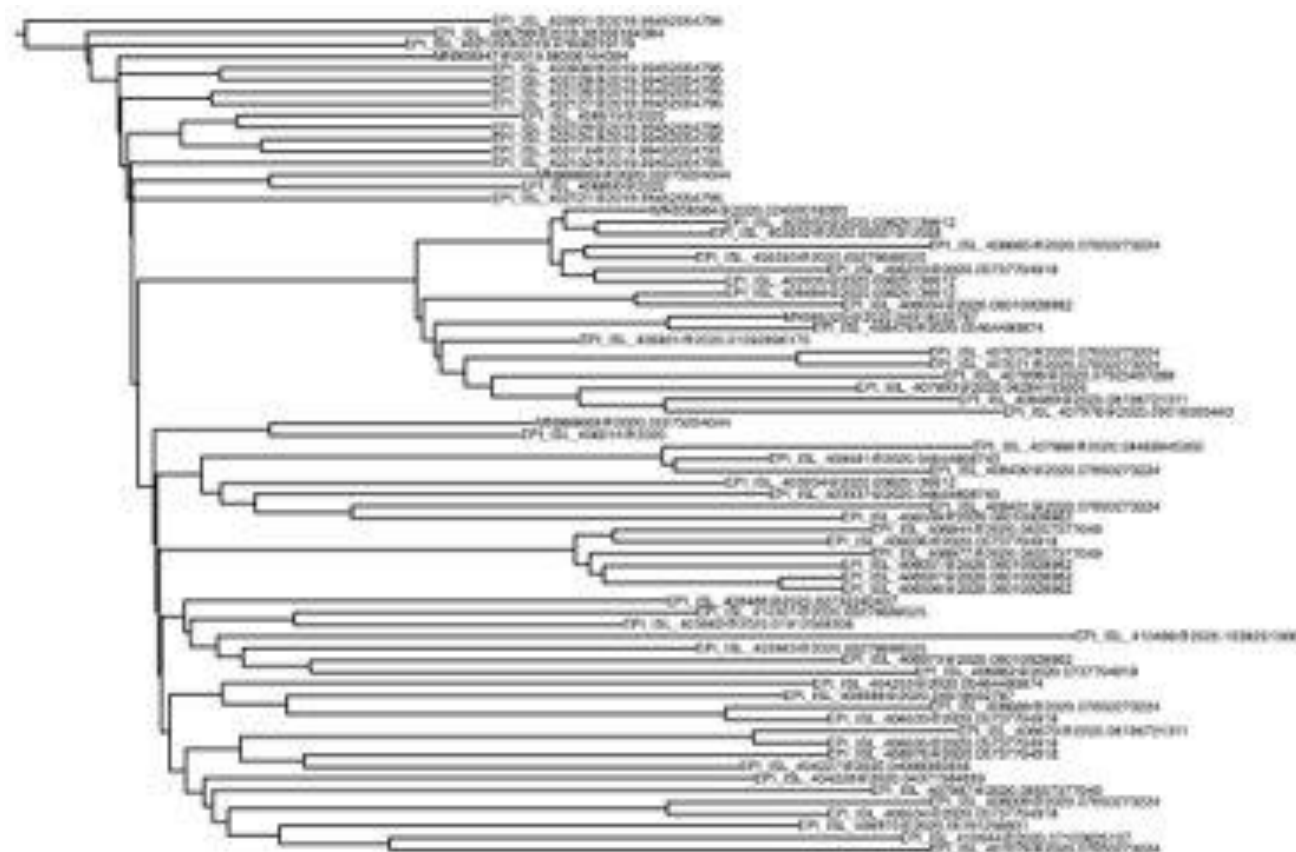

Substitution
model


Clock
model

Concept summary

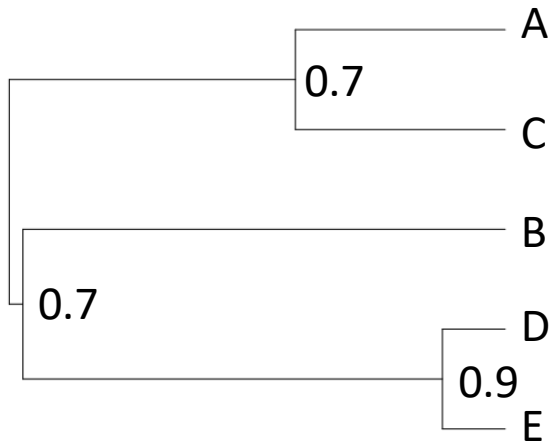
- Model parameters can be found in the posterior function, $P(M, \text{par} | \text{data})$.
- Arguably, the tree is also a parameter, but we rarely inspect its trace (use AWTY).
- Summary statistics can be calculated from the model parameters.
- Probability densities are in log units and should not be compared between data sets, or used to select models.
 - The 'posterior' trace is an unnormalised posterior and NOT the posterior probability.

Summarising trees

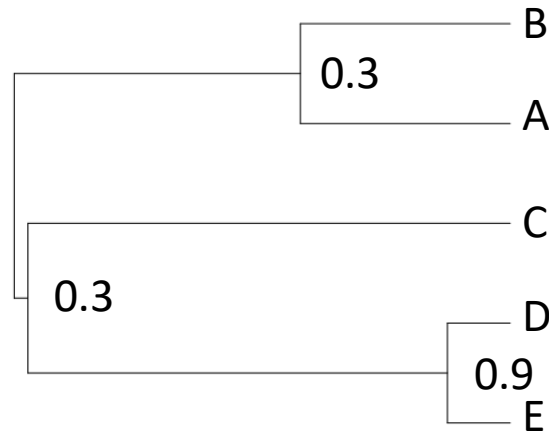


The highest clade credibility tree

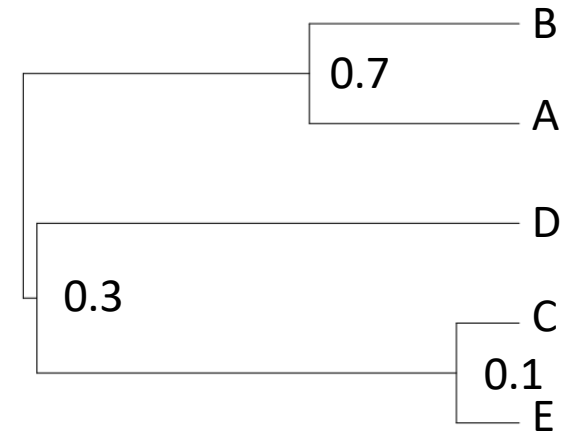
- The highest product of node posterior probabilities



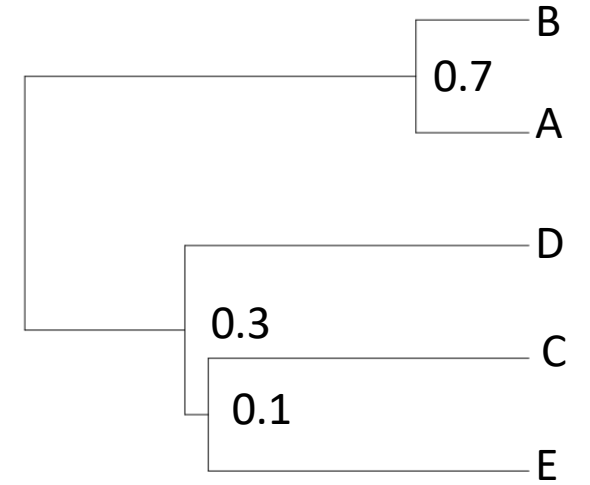
7,000 trees



2,000 trees



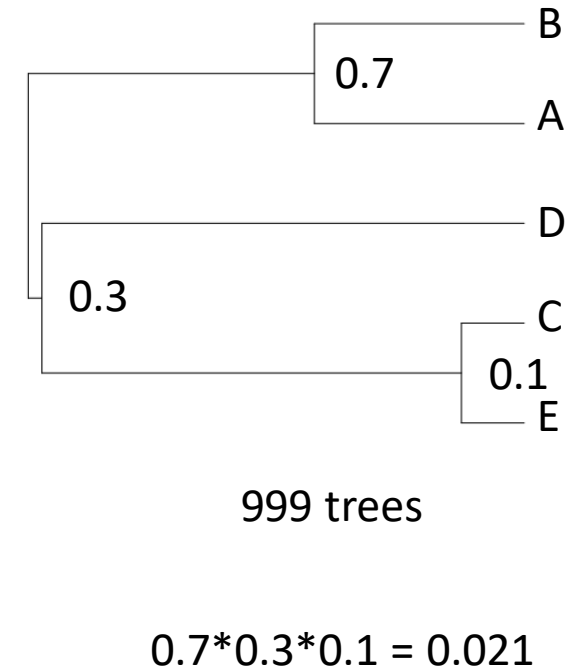
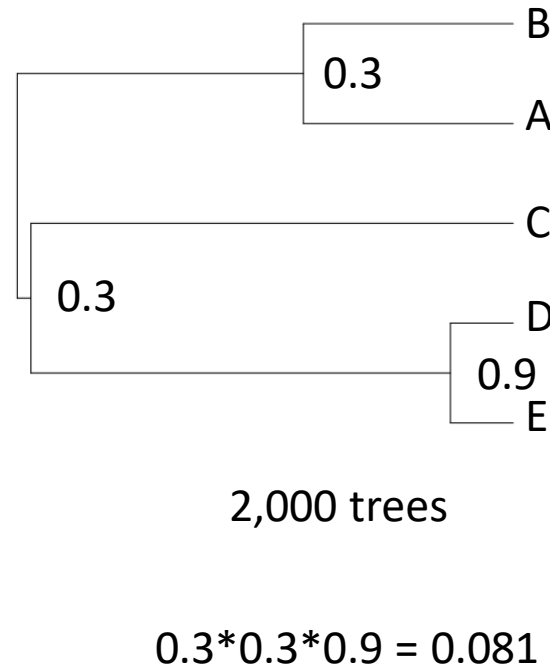
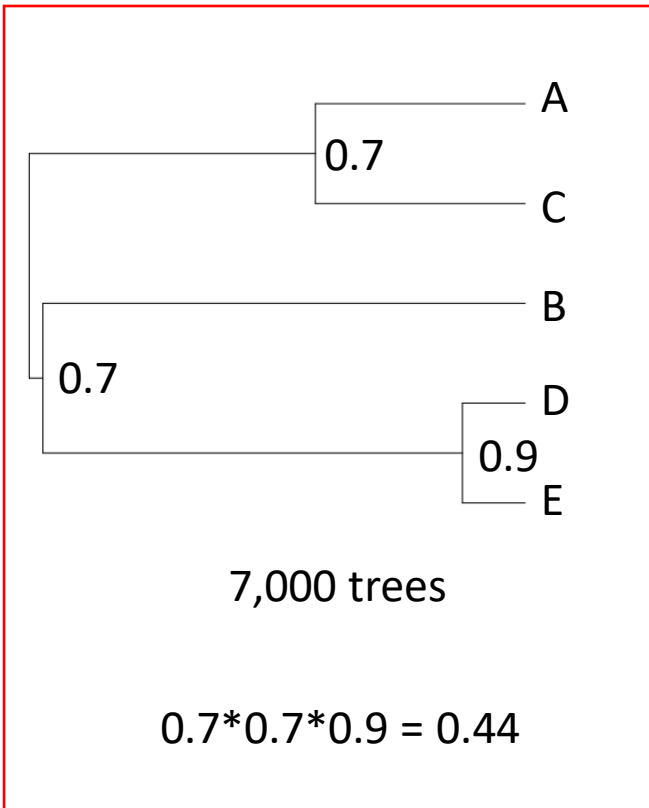
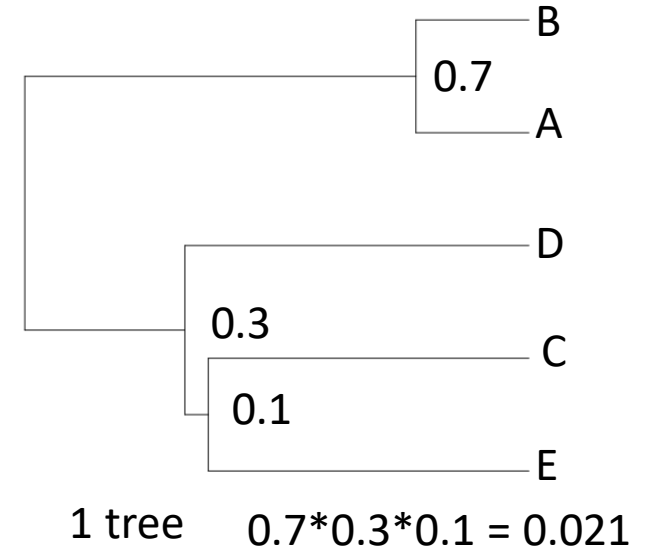
999 trees



1 tree

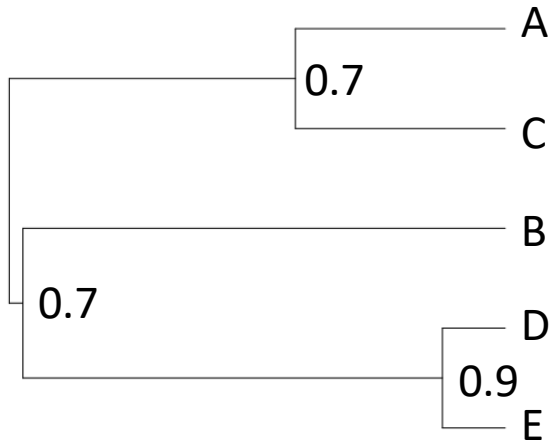
The highest clade credibility tree

- The highest product of node posterior probabilities

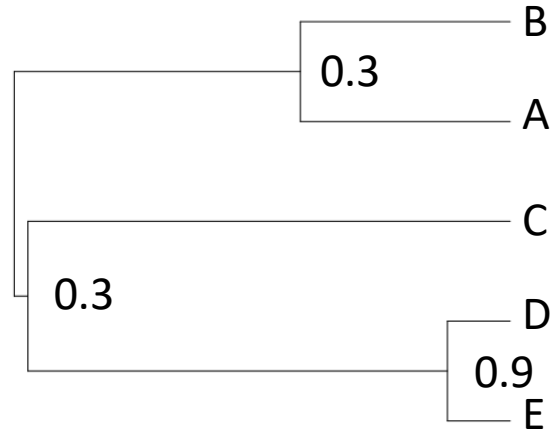


The maximum *a posteriori* tree

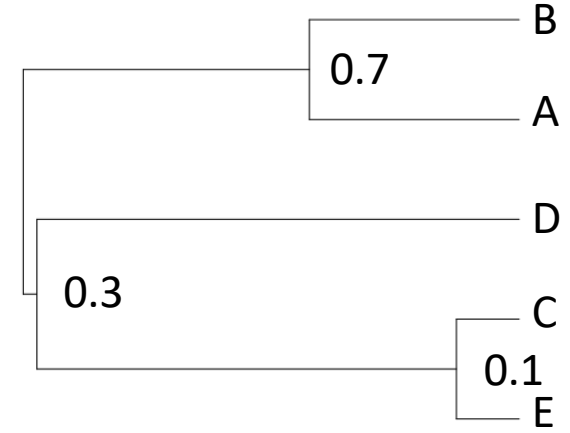
- The tree with highest posterior density.



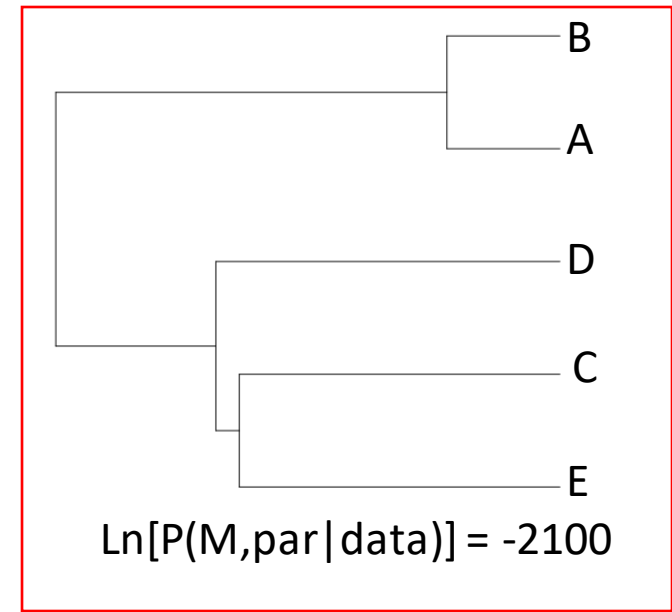
$$\text{Ln}[P(\text{M,par} | \text{data})] = -3000$$



$$\text{Ln}[P(\text{M,par} | \text{data})] = -4000$$

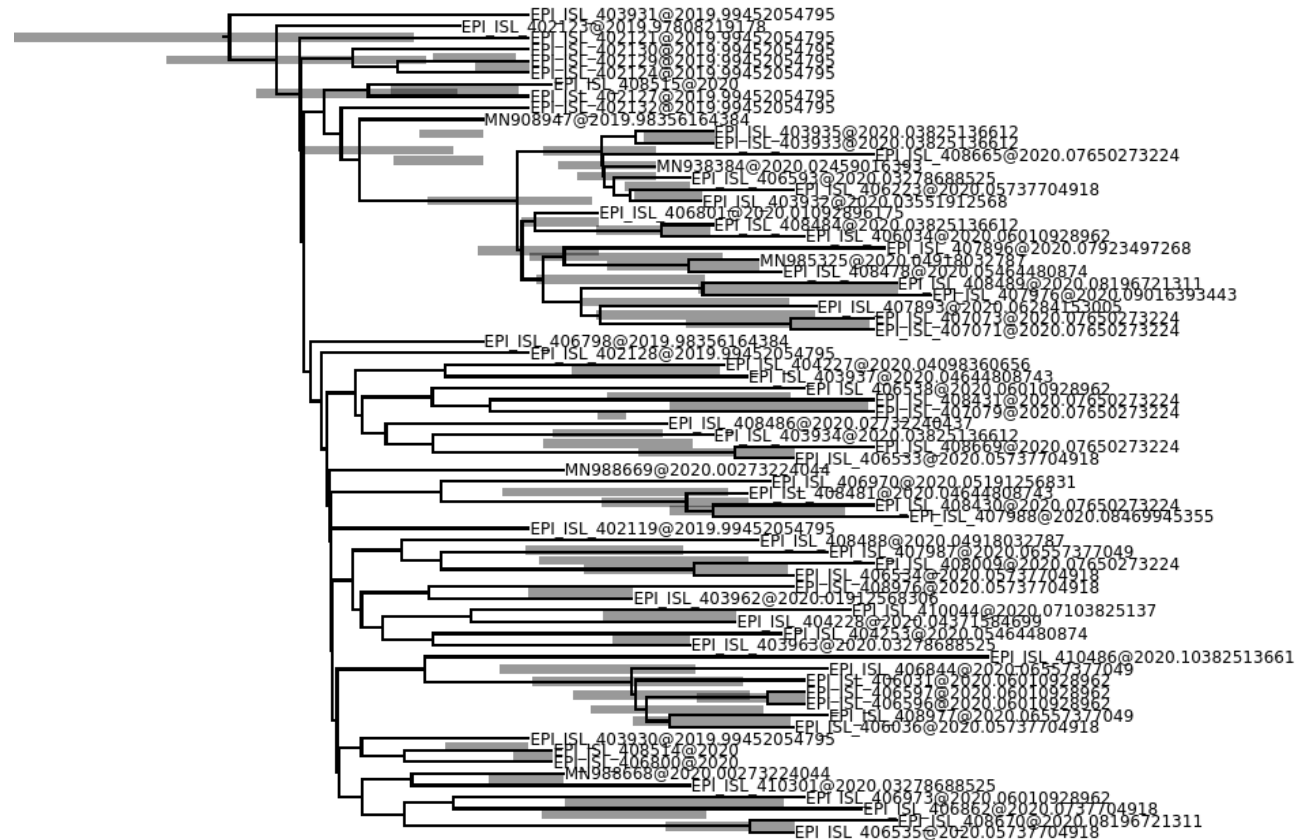


$$\text{Ln}[P(\text{M,par} | \text{data})] = -4100$$



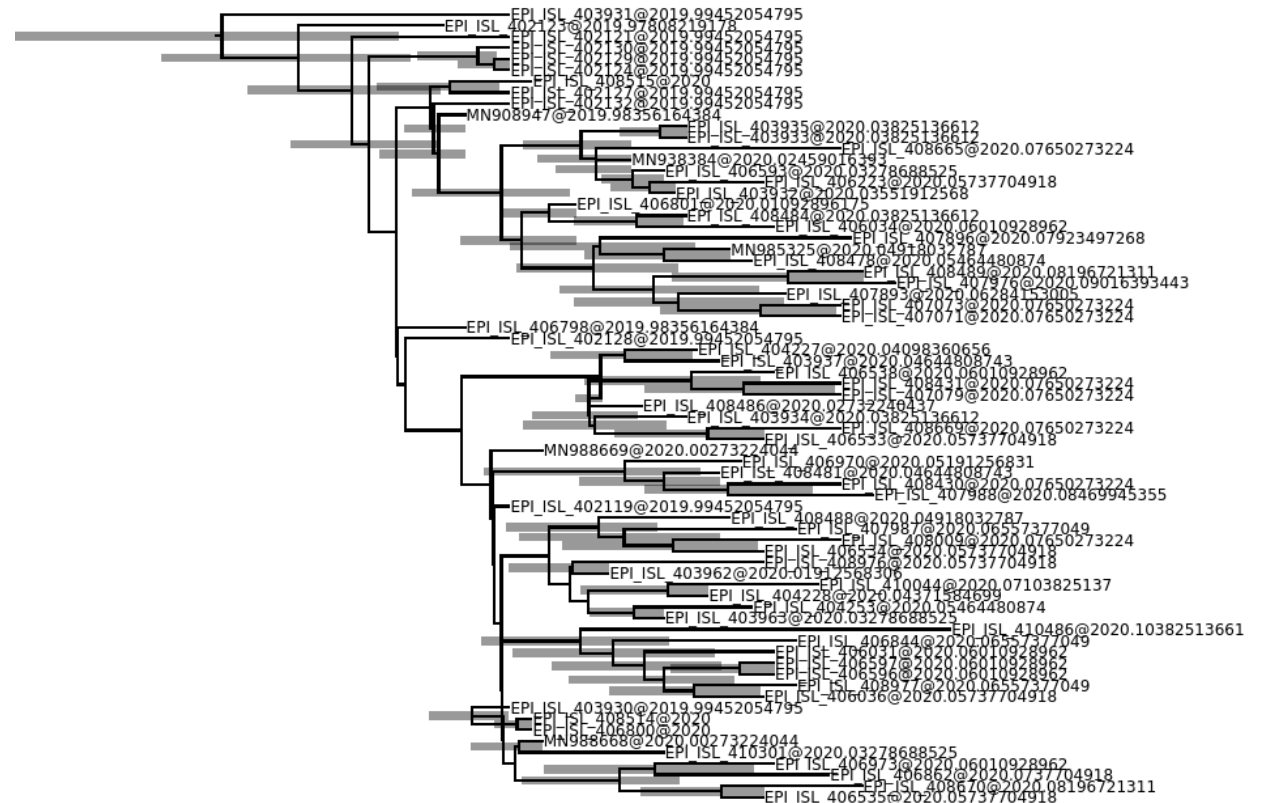
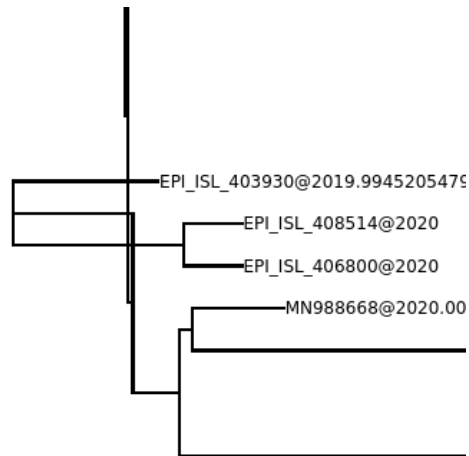
Node heights

- Keep: keep node heights of summary tree.
 - Can have credible intervals outside of nodes.



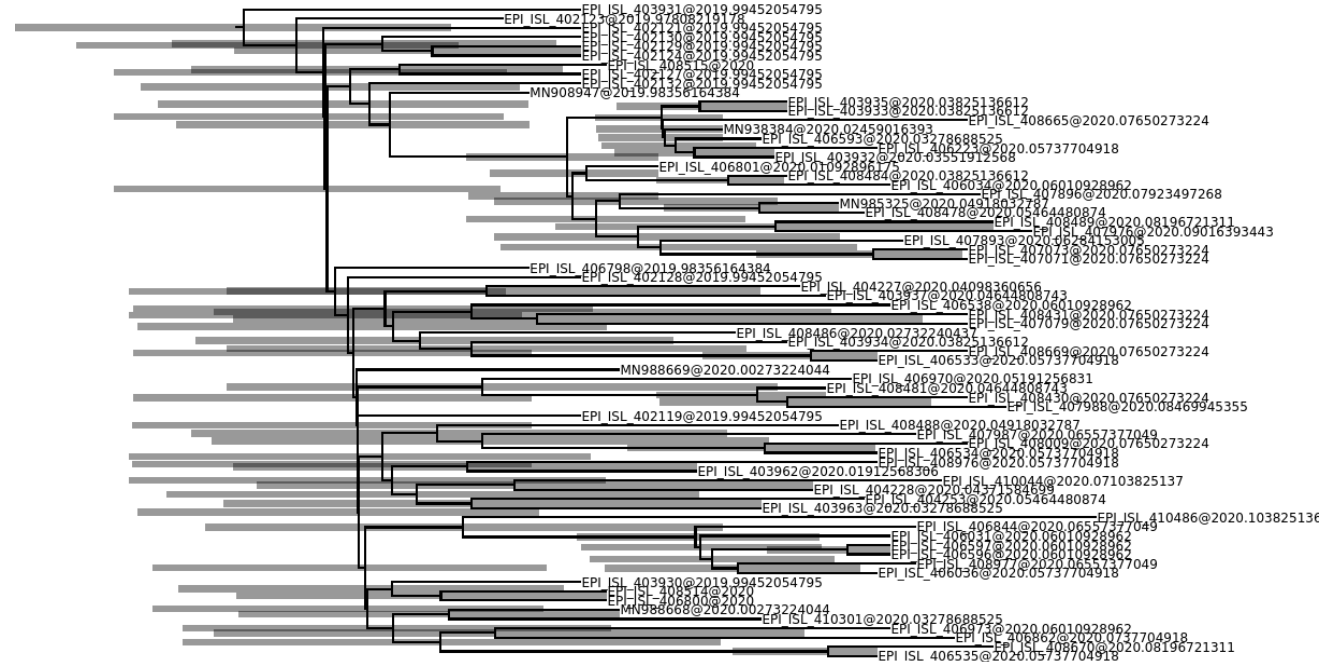
Node heights

- Mean: take average of nodes if present in other trees.
 - Can have negative branch lengths, especially for poorly supported nodes.



Node heights

- Common ancestor: take the mean age of taxa for all trees (regardless of monophyly).
 - Can lead to very deep divergences and wide credible intervals.
 - No negative branch lengths.



Concept summary

- Trees can be summarised in different ways and they all have compromises.
- Approaches to summarising node times can distort branch lengths (e.g. mean or common ancestor).
- The 'best' tree is not necessarily a reasonable summary (e.g. a very wide posterior distribution).