


National Institute of Allergy and Infectious Diseases


Phylogenetics and Sequence Analysis

Lecture 7: Molecular Evolutionary Analysis of Pathogens Using BEAST

Kurt Wollenberg, PhD
Phylogenetics Specialist
Bioinformatics and Computational Biology Branch (BCBB)
Fall 2015





Kurt Wollenberg, PhD
Phylogenetics Specialist
Bioinformatics and Computational Biology Branch (BCBB)





Course Organization

- Building a clean sequence
- Collecting homologs
- Aligning your sequences
- Building trees
- **Further Analysis**



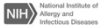
Lecture Organization

- What's so special about pathogens?
- What is BEAST?
- A short tour of Bayesian MCMC analysis
- An overview of the BEAST package
- **BEAST Analysis Demo**
- Odds and ends



What's so special about pathogens?


- Short generation time
- Rapid evolution
- Genotypes - easy, phenotypes - hard
- Large populations
- Structured populations
- Rigorous temporal sampling of genotypes

 National Institute of Allergy and Infectious Diseases

NIAD

What is BEAST?

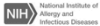
- **B**ayesian **E**volutionary **A**nalysis **S**ampling **T**rees
- A collection of programs for performing Bayesian MCMC analysis of molecular sequences
- Can incorporate sample time information
- Can perform a broad range of other evolutionary analyses using sequence data.

 National Institute of Allergy and Infectious Diseases

NIAD

What is Bayesian analysis?

- Calculation of the probability of parameters (tree, substitution model) given the data (sequence alignment)
- $p(\theta|D) = (\text{Likelihood} \times \text{prior}) / \text{probability of the data}$
- $p(\theta|D) = p(D|\theta)p(\theta)/p(D)$

 National Institute of Allergy and Infectious Diseases

NIAD

What is Bayesian analysis?

Likelihood that this die is unbiased?

Roll	Observed	Expected
1	5	6
2	5	6
3	6	6
4	6	6
5	10	6
6	15	6

NIH National Institute of Allergy and Infectious Diseases

Bayesian Analysis

Exploring the posterior probability distribution

Posterior probabilities of trees and parameters are approximated using Markov Chain Monte Carlo (MCMC) sampling

Markov Chain: A statement of the probability of moving from one state to another

NIH National Institute of Allergy and Infectious Diseases

What is MCMC?

Markov Chain Monte Carlo

Markov chain

One link in the chain

Monte Carlo

Choosing a link

NIH National Institute of Allergy and Infectious Diseases

What is MCMC?

Markov Chain Monte Carlo: accept or reject?

$$\Pr(\text{accept}) = \min \left(1, \frac{\Pr(a)}{\Pr(A)} \times \frac{\Pr(X|a)}{\Pr(X|A)} \times \frac{\Pr(A|a)}{\Pr(a|A)} \right)$$

NIH National Institute of Allergy and Infectious Diseases

What is BEAST?

The Programs:

- BEAUti - Creating XML input files
- BEAST - MCMC analysis of molecular sequences
- Tracer - Viewing MCMC output
- LogCombiner - Combining output files
- TreeAnnotator - Generate the consensus tree
- FigTree - Drawing a tree

NIH National Institute of Allergy and Infectious Diseases

Different types of BEAST analyses


- Calculating a Bayesian coalescent phylogeny
- **Calculating a Time-Stamped Bayesian coalescent**
- Phylogeographic analysis (time and location data)
- Estimated population dynamics (Bayesian skyride and skygrid)
- Combined gene and species phylogeny estimate (*BEAST)

NIH National Institute of Allergy and Infectious Diseases


Defining your analysis

- Prior knowledge of tree?
- Calibrating nodes?
- Substitution model?
- Effective population sizes?
- What priors to use?

NIAID



Setting up the analysis: BEAUTi



Partition Name	File Name	Type	Size	Data Type	Site Model	Clock Model	Partition Tree
0000	0000.nex	35	628	nucleotide	0000	1 0000	0 0000


Data: 35 taxa, 1 partition, fix clock rate to 1.0 in nucleotide_group. Generate BEAUTi file...

NIAID

Setting up the analysis: BEAUTi

- Incorporate sample times
- Substitution model parameters
- Strict or relaxed clock?
- Tree prior
- Substitution model priors
- Adjustments from previous runs (operators)
- Setting the chain

NIAID



Running the analysis: BEAST

- Load your input file
- Use BEAGLE?
- That's it

NIAID

Evaluating the analysis: Tracer

- Check for convergence
- Viewing parameter estimates from multiple runs
- Extract parameter estimates and statistics

NIAID


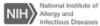
Evaluating the analysis: Tracer

- Check for convergence
- Viewing parameter estimates from multiple runs
- Extract parameter estimates and statistics




NIAID

Evaluating the analysis: Tracer

- What if my analysis didn't converge?
- Can I make multiple simultaneous runs?
 - Swarm on Biowulf


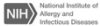


Merging output files: LogCombiner



Merging output files: LogCombiner

- Log files vs Tree files
- Selecting files
- Specifying burn-in (number of states)
- Specifying subsampling
- Specifying output file



Merging output files: LogCombiner

- Burn in?

NIH National Institute of Allergy and Infectious Diseases

NIAD

Calculating the tree: TreeAnnotator

NIH National Institute of Allergy and Infectious Diseases

NIAD

Calculating the tree: TreeAnnotator

- Burn in? Number of trees or the number of steps.
- Tree Type: MCC, Max sum of CC, or target
- Node heights: target, mean, or median
- Specify input and output files

NIH National Institute of Allergy and Infectious Diseases

NIAD

Drawing trees: FigTree

The screenshot shows the FigTree v1.4.0 interface. The main window displays a phylogenetic tree with a scale bar at the bottom ranging from 700.0 to 0.0. On the left, there is a 'Layers' panel with various options like 'Zoom', 'Appearance', 'Time Scale', 'Node Labels', etc. A vertical 'NIAID' logo is positioned on the right side of the slide.

Drawing trees: FigTree

- Specifying additional values (esp. posterior probabilities)
- Tree appearance
- Ordering branches
- Re-rooting
- Exporting graphics

The slide lists five bullet points related to FigTree. At the bottom left is the NIH logo with the text 'National Institute of Allergy and Infectious Diseases'. A vertical 'NIAID' logo is on the right side.

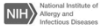

Running BEAST

DEMO

The slide features the title 'Running BEAST' and a large 'DEMO' text in the center. At the bottom left is the NIH logo with the text 'National Institute of Allergy and Infectious Diseases'. A vertical 'NIAID' logo is on the right side.

Running BEAST: Demo

- Site Models
- Substitution Models
 - HKY - Unequal base frequencies and transition/transversion rate ratio
 - Must specify prior and initial estimates for transition/transversion rate ratio
 - GTR - Unequal base frequencies and each substitution has its own rate parameter
 - Must specify prior and initial estimates for each substitution rate (relative to C-T rate)

Running BEAST: Demo

- Site Models
 - Substitution models
 - Site heterogeneity models
- Get estimates from the program jModelTest



```

Model selected: HKY+G
-nl = 1276.8109
K = 9
AIC = 3371.6218

Base frequencies:
freqA = 0.2259
freqC = 0.3199
freqG = 0.2405
freqT = 0.2137

Substitution model:
Rate matrix
R(A) [A-C] = 0.2494
R(B) [A-G] = 0.0655
R(C) [A-T] = 0.7435
R(D) [C-G] = 0.3907
R(E) [C-T] = 0.0655
R(F) [G-T] = 1.0000

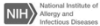

Among-site rate variation
Proportion of invariable sites (I) = 0.6508
Variable sites (V)
Gamma distribution shape parameter = 0.5913
    
```

Running BEAST: Demo

Site Models

- Site heterogeneity models
 - Gamma
 - Modeling rate of change using a discrete gamma distribution
 - Invariant
 - Percent of non-variable sites in the data

Running BEAST: Demo

Site Models

- Site heterogeneity models
 - Gamma

Mean = $k\theta$
Shape parameter = θ
Coefficient of Variation = $1/\theta$

NIH National Institute of Allergy and Infectious Diseases

NIAD

Setting up the analysis: Models

Testing Models and Priors

Does the relaxed clock fit the data?

Examining the uclid.stdev distribution

NIH National Institute of Allergy and Infectious Diseases

NIAD

Running BEAST

Testing Models and Priors

Path Sampling/Stepping Stone analysis

- Estimation of marginal likelihoods under different analysis parameters.
- Invoke on MCMC tab in BEAUti.
- Separate runs necessary for each changed parameter.
- Runs a complete MCMC analysis, then the X PS/SS iterations.

NIH National Institute of Allergy and Infectious Diseases

NIAD

Running BEAST

Testing Models and Priors

Path Sampling/Stepping Stone analysis

NIAD

Running BEAST

Testing Models and Priors

Path Sampling/Stepping Stone analysis log marginal likelihoods

	Path Sampling	Stepping Stone
HKY/strict clock	-4725.85	-4728.68
HKY+gi/strict	-4515.99	-4518.05
HKY+gi/LN relaxed	-4436.10	-4438.75
GTR/strict clock	-4746.62	-4749.14
GTR+gi/strict	-4526.87	-4529.05
GTR+gi/LN relaxed	-4548.39	-4551.22

NIAD

Running BEAST: swarm on Biowulf

- Requires a .swarm file
 - A text file containing

```

beast beastJob_1.xml > beastJob_1out.txt
sleep 2; beast beastJob_2.xml > beastJob_2out.txt
sleep 4; beast beastJob_3.xml > beastJob_3out.txt
sleep 6; beast beastJob_4.xml > beastJob_4out.txt
sleep 8; beast beastJob_5.xml > beastJob_5out.txt
    
```

- Run in command line

```

[username]$ swarm -f beastInput.swarm -module BEAST
    
```

NIAD

Seminar Follow-Up Site

For access to past recordings, handouts, slides visit this site [from the NIH network](http://collab.niaid.nih.gov/sites/research/SIG/Bioinformatics/): <http://collab.niaid.nih.gov/sites/research/SIG/Bioinformatics/>

1. Select a Subject Matter

Recommended Browsers:

- IE for Windows.
- Safari for Mac (Firefox on a Mac is incompatible with NIH Authentication technology)

Login

- If prompted to log in use "NIH" in front of your username

2. Select a Topic

View:

- Seminar Details
- Handout and Reference Docs
- Relevant Links
- Seminar Recording Links

NIAD

37



Retrieving Slides/Handouts

This lecture series

NIAD

38



Retrieving Slides/Handouts

This file

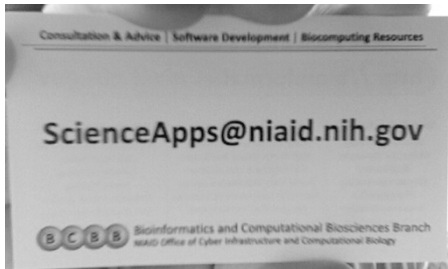
This lecture

NIAD

39



Questions?



NIH National Institute of Allergy and Infectious Diseases

NIAD

40

Thank you



NIH National Institute of Allergy and Infectious Diseases

NIAD

41
