Treespaces
oooooo

Mean tree
oooo

Convergence Assessment
ooooo

# Statistics for phylogenetic time trees

Lars Berling

11/02/2025

Treespaces
oooooo

Mean tree
oooo

Convergence Assessment
ooooo

## The goal

Treespaces
oooooo

Mean tree
oooo
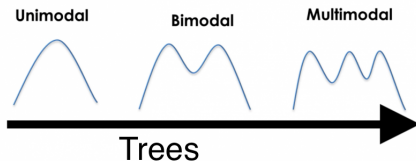
Convergence Assessment
ooooo

# The goal



**Bayesian Phylogenetic Inference via MCMC**

- Key object of interest is often the rooted tree topology.
- **MCMC Output**: Sample of trees (typically thousands)
- Challenges in estimating mean and variance in treespace due to its high-dimensional, non-Euclidean nature.
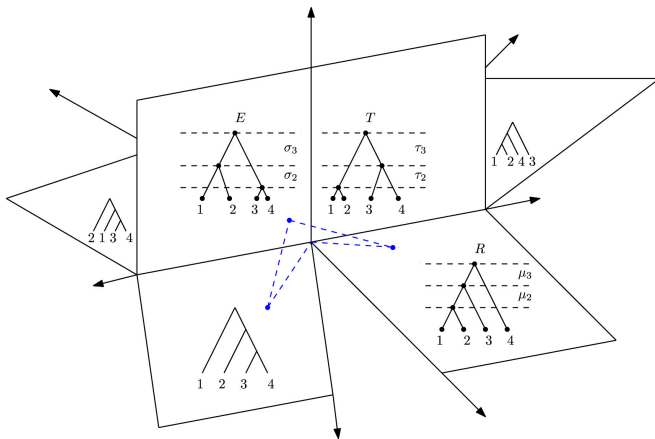
Treespaces
●○○○○○

Mean tree
○○○○

Convergence Assessment
○○○○○

# Treespaces

## $\tau$-space



Figure: Three-dimensional projection of a part of $\tau$ space with 4 taxa.

Treespaces
○○●○○○

Mean tree
○○○○

Convergence Assessment
○○○○○

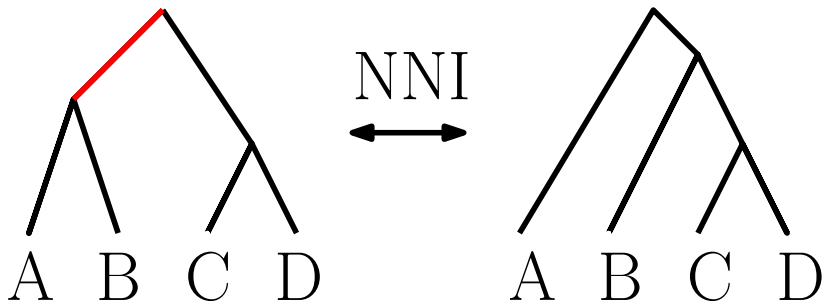# The downfall of stratified spaces



Figure: Geodesics and cone paths

**Treespaces**
○○○●○○

Mean tree
○○○○

Convergence Assessment
○○○○○

## The space of ranked trees

## The space of ranked trees

Treespaces
○○○●○○

Mean tree
○○○○

Convergence Assessment
○○○○○

# The space of ranked trees

**Treespaces**
○○○●○○

Mean tree
○○○○

Convergence Assessment
○○○○○

## The space of ranked trees



rank NNI

A B C D      A B C D      A B C D

### The RNNI graph

This graph is the **treespace** of ranked trees, **R**anked **N**earest **N**eighbour **I**nterchange space.

**Treespaces**
○○○○●○

Mean tree
○○○○

Convergence Assessment
○○○○○

## The shortest path problem

### Induced distance

The minimal **number of rearrangement operations** to transform one tree into another

Equivalent: find a **path of minimal length** in the RNNI graph

**Treespaces**
○○○○●○

Mean tree
○○○○

Convergence Assessment
○○○○○

## The shortest path problem

### Induced distance

The minimal **number of rearrangement operations** to transform one tree into another

Equivalent: find a **path of minimal length** in the RNNI graph

### Theorem

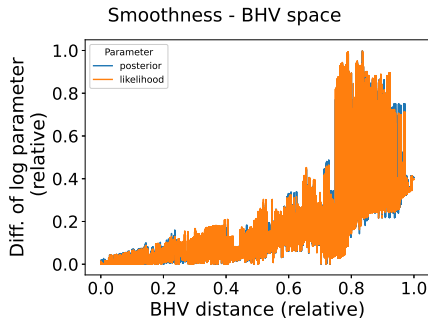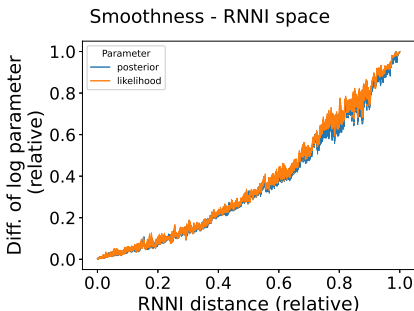The **FindPath** algorithm computes a shortest path in RNNI, with time complexity $O(n^2)$. [a]

_____

[a] Collienne, Lena, and Alex Gavryushkin. "Computing nearest neighbour interchange distances between ranked phylogenetic trees." Journal of Mathematical Biology 82.1-2 (2021): 8.

**Treespaces**
○○○○○●

Mean tree
○○○○

Convergence Assessment
○○○○○

## Probability distributions are 'continuous'

Comparing **probability distributions** in BHV and RNNI tree space on **one data set**. **x**-Axis displays **tree metric distance** (relative), **y**-axis displays **difference in probability** (relative)
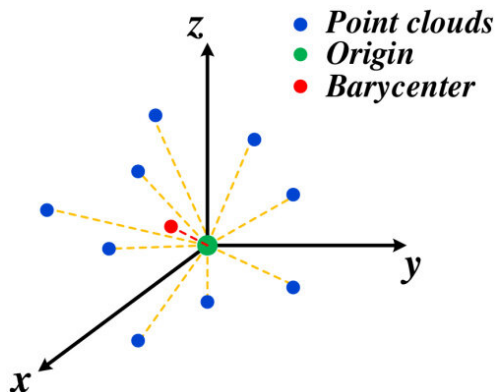
# Probability distributions are 'continuous'

Comparing **probability distributions** in BHV and RNNI tree space on **one data set**. **x**-Axis displays **tree metric distance** (relative), **y**-axis displays **difference in probability** (relative)
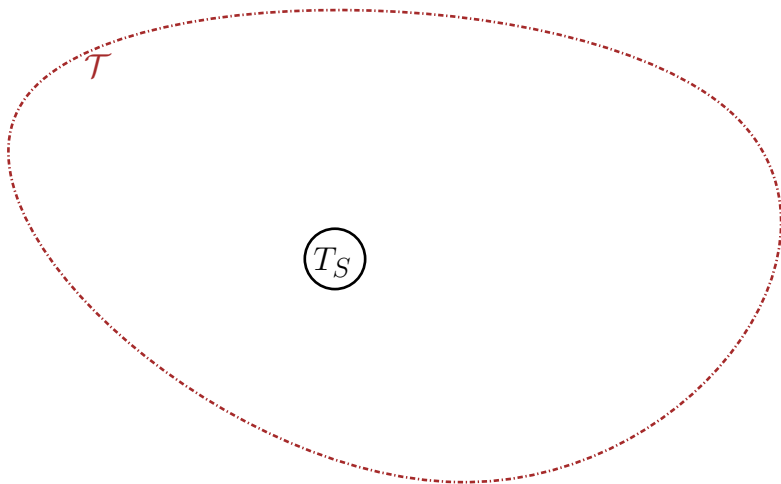


Smoothness - RNNI space

Smoothness - BHV space

Treespaces
oooooo

Mean tree
●ooo

Convergence Assessment
ooooo

# Mean tree

Treespaces
oooooo

Mean tree
o●oo

Convergence Assessment
ooooo

# Geometric means



- **Point clouds**
- **Origin**
- **Barycenter**

### Fréchet variance

$Var(t)_{\mathcal{T}} = \sum_{t_i \in \mathcal{T}} d(t_i, t)^2$

Treespaces
oooooo

Mean tree
oooo

Convergence Assessment
ooooo

# The algorithm

Treespaces
oooooo

Mean tree
ooeo

Convergence Assessment
ooooo

# The algorithm

Treespaces
○○○○○○

Mean tree
○○●○

Convergence Assessment
○○○○○

# The algorithm

Treespaces
○○○○○○

Mean tree
○○●○

Convergence Assessment
○○○○○

# The algorithm

Treespaces
○○○○○○

Mean tree
○○●○

Convergence Assessment
○○○○○

## The algorithm

Treespaces
oooooo

Mean tree
oooo

Convergence Assessment
ooooo

# The algorithm

Treespaces
○○○○○○

**Mean tree**
○○●○

Convergence Assessment
○○○○○

# The algorithm

Treespaces
oooooo

Mean tree
ooo●

Convergence Assessment
ooooo

# Comparing Likelihood to MCC



## MCC

Maximum Clade Credibility Tree from treeannotator (BEAST).

Treespaces
oooooo

Mean tree
oooo

Convergence Assessment
●oooo

# Convergence Assessment

Treespaces
○○○○○○

Mean tree
○○○○

Convergence Assessment
○●○○○

# Convergence

- Sampling from the stationary distribution
- $\rightarrow$ Parameter trace no trend
- **E**ffective **S**ample **S**ize at least 200 (rule of thumb)



Figure: Not converged



Figure: Converged

Treespaces
oooooo

Mean tree
oooo

Convergence Assessment
o●oooo

## Convergence

- Sampling from the stationary distribution
- → Parameter trace no trend
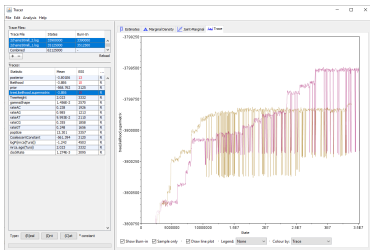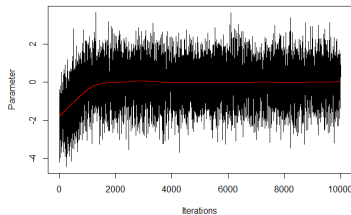- **E**ffective **S**ample **S**ize at least 200 (rule of thumb)



Figure: Not converged



Figure: Converged

Treespaces
oooooo

Mean tree
oooo

Convergence Assessment
ooo●oo

## Do two sets have the same underlying distribution?

independently sampled sets of trees

Treespaces
○○○○○○

Mean tree
○○○○

Convergence Assessment
○○●○○
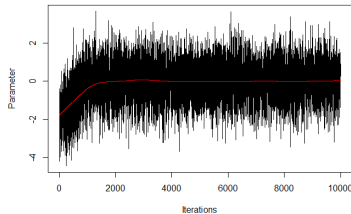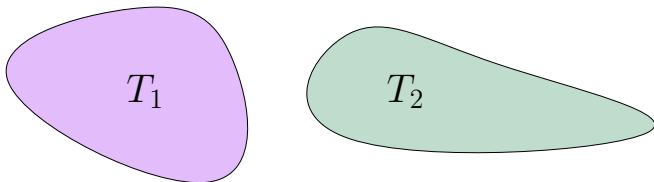
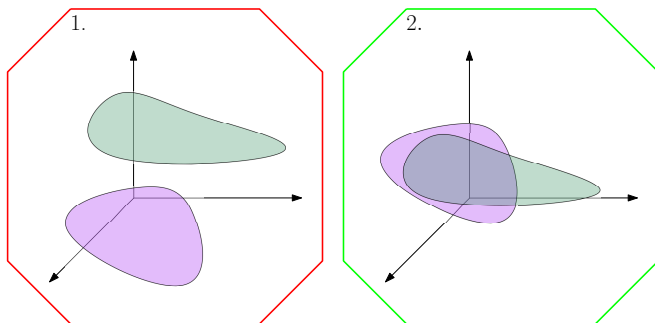# Do two sets have the same underlying distribution?

independently sampled sets of trees

Treespaces
○○○○○○

Mean tree
○○○○

Convergence Assessment
○○○●○

## Gelman Rubin diagnostic for trees



$\star = $ new sample in $T_1$
$\cdots = $ RNNI distance

# Gelman Rubin diagnostic for trees

## Potential scale reduction factor

$$PSRF(t|\mathcal{T}_1, \mathcal{T}_2) = \sqrt{\frac{Var(t)_{\mathcal{T}_2}}{Var(t)_{\mathcal{T}_1}}}, t \in \mathcal{T}_1 \ ^{a}$$

---

[a]Inference from Iterative Simulation Using Multiple Sequences, A. Gelman and D. Rubin

## Fréchet variance (normalized)

$$Var(t)_{\mathcal{T}} = \frac{\Sigma_{t_i \in \mathcal{T}} d(t_i, t)^2}{|\mathcal{T}|}$$

Treespaces
oooooo

Mean tree
oooo

Convergence Assessment
ooo●o

# Gelman Rubin diagnostic for trees

# Gelman Rubin diagnostic for trees



Chain 1

Chain 2

☐ = above threshold

☐ = within threshold

### Further assessment of overlap

- **E**ffective **S**ample **S**ize at least 200 (rule of thumb)
- Further downstream analysis: Summarizing/ computing a mean tree

# Results on DS1-DS11

| | ESS-threshold | mean - 0.05 | | | | | | mean - 0.02 | | | | | | mean - 0.01 | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | C1 | C2 | C3 | C4 | C5 | C6 | C1 | C2 | C3 | C4 | C5 | C6 | C1 | C2 | C3 | C4 | C5 | C6 |
| DS1 | 200 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 |
| | 500 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| DS2 | 200 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 |
| | 500 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 |
| DS3 | 200 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| | 500 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| DS4 | 200 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 500 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| DS5 | 200 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 1 |
| | 500 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 |
| DS6 | 200 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| | 500 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| DS7 | 200 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| | 500 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| DS8 | 200 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 |
| | 500 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 1 |
| DS9 | 200 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| | 500 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 1 |
| DS10 | 200 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 1 |
| | 500 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| DS11 | 200 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| | 500 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

Orange cell = convergence, Blue cell = non-convergence

# Conclusion

## Summary and Outlook

- Developed a mean tree within the RNNI treespace
- $\rightarrow$ Theoretical foundation is still an open problem
- Convergence assessment of samples of trees
- $\rightarrow$ Multimodal distributions?
- $\rightarrow$ Advanced characterization of tree distributions

# Conclusion

### Software

- BEAST2 package: ASM
  https://github.com/rbouckaert/asm
- Python package: tetres
  https://github.com/bioDS/tetres

### References

- Lars Berling; Lena Collienne; Alex Gavryushkin, Estimating the mean in the space of ranked phylogenetic trees, Bioinformatics (2024)
- Lars Berling; Remco Bouckaert; Alex Gavryushkin, An Automated Convergence Diagnostic for Phylogenetic MCMC Analyses, IEEE/ACM TCBB (2024)