# Incorporating Heterogenous Data Sources in Phylogenetic Modeling

April Wright
LBRN Annual Meeting
January 18

# Big Data



Mazberry.com



Datanami.com

# Big Data

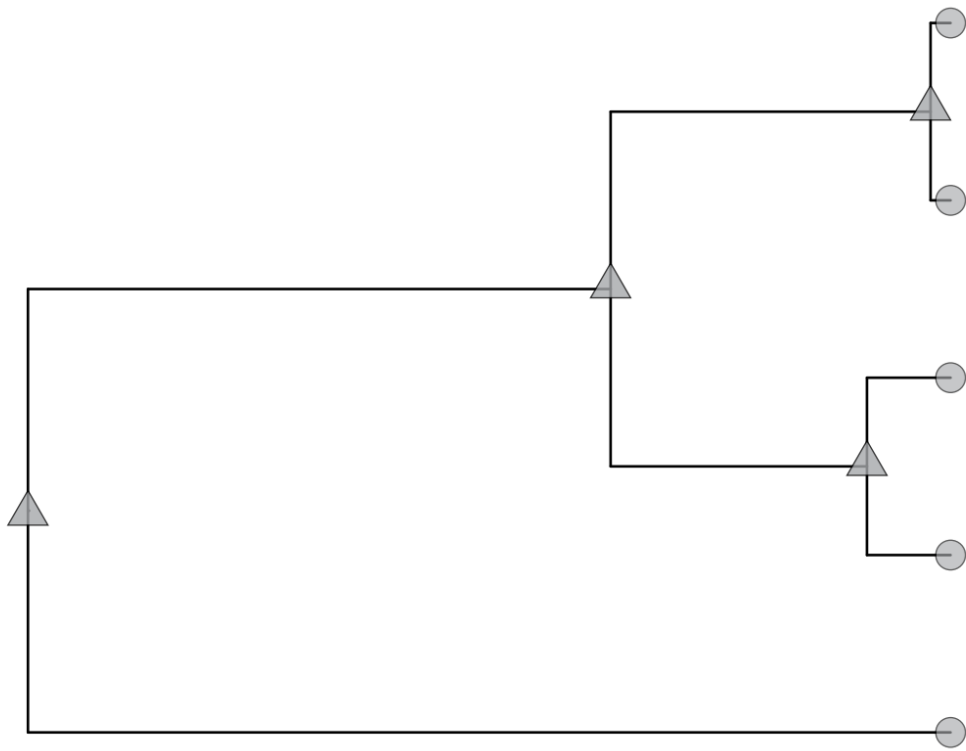Heterogeneous data = data coming from multiple sources, each with their own generating and collection process

Mzberry.com

Dtnmi.com

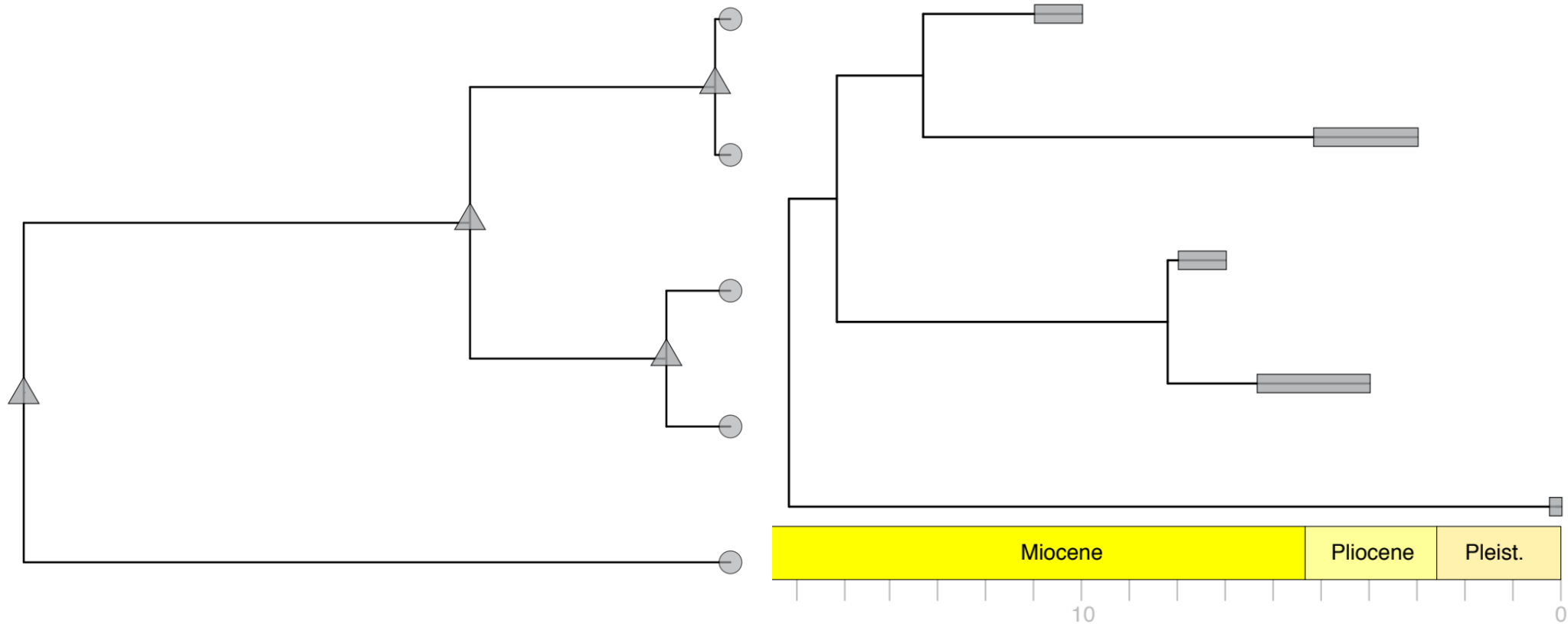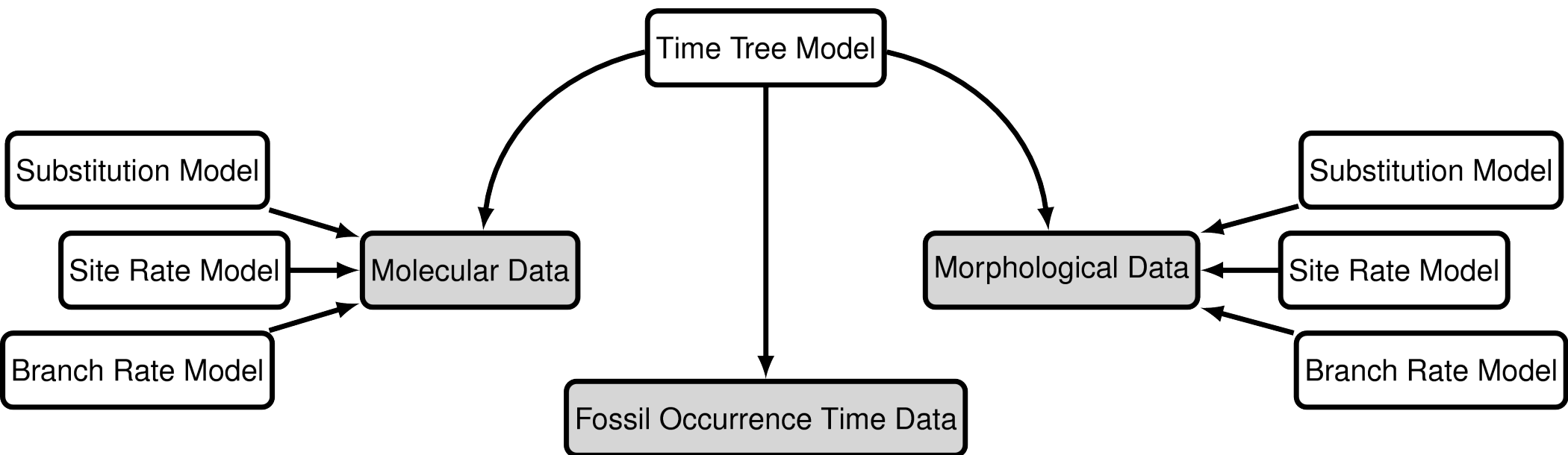Estimating phylogenetic trees usually involves modeling the evolution of nucleotide sequences over time

Wright and Warnock, in review

Scaling phylogeny to time involves adding in age information, and often phenotypic information

Wright and Warnock, in review

Wright, Pett and Heath

| Taxon 1 | A | C | T | A | C | T | C | G |   |
|---------|---|---|---|---|---|---|---|---|---|
| Taxon 2 | A | C | T | A | A | T | G | T | C |
| Taxon 3 | A | T | T | A | C | T | G | T | G |
| Taxon 4 | G | G | A | A | C | T | G | G | T |
| Taxon 5 | G | G | C | T | C | T | G | A | A |

Substitution Model

Site Rate Model

Branch Rate Model

Molecular Data

Wright, Pett
and Heath

Wright, Pett and Heath

| Taxon 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---------|---|---|---|---|---|---|---|---|---|
| Taxon 2 | 0 | 2 | 0 | 2 | 1 | 1 | 0 | 0 | 0 |
| Taxon 3 | 1 | 2 | 0 | 2 | 3 | 0 | 0 | 1 | 1 |
| Taxon 4 | 1 | 1 | 1 | 1 | 3 | 1 | 0 | 1 | 1 |
| Taxon 5 | 1 | 1 | 0 | 1 | 2 | 1 | 1 | 1 | 0 |

Morphological Data

Substitution Model

Site Rate Model

Branch Rate Model

Wright, Pett and Heath

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Taxon 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Taxon 2 | 0 | 2 | 0 | 2 | 1 | 1 | 0 | 0 | 0 |
| Taxon 3 | 1 | 2 | 0 | 2 | 3 | 0 | 0 | 1 | 1 |
| Taxon 4 | 1 | 1 | 1 | 1 | 3 | 1 | 0 | 1 | 1 |
| Taxon 5 | 1 | 1 | 0 | 1 | 2 | 1 | 1 | 1 | 0 |

Morphological Data

Substitution Model

Site Rate Model

Branch Rate Model

exponential

$\delta_\beta$ → $\beta_\beta$

discrete beta

$\delta_\alpha$ → $\alpha_\beta$

$\pi_i$ → site frequencies

exponential

$Q_i$ → site matrices

$i \in N$

to the phyloCTMC

Wright, Pett and Heath

exponential

$\delta_\mu$ → $\mu$    extinction rate

exponential

$\delta_\lambda$ → $\lambda$

speciation rate

exponential

$\psi$ ← $\delta_\psi$

fossilization rate

$\mathcal{T}$

$a$

uniform

$\phi$

$b$

origin time

$\rho$

sampling probability

to the phyloCTMCs

Wright, Pett and Heath

Time Tree Model

Substitution Model

Site Rate Model

Molecular Data

Branch Rate Model

Fossil Occurrence Time Data

Substitution Model

Morphological Data

Site Rate Model

Branch Rate Model

Wright, Pett and Heath

The complete outcome
of the diversification and
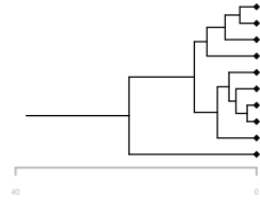sampling processes

The reconstructed tree

Model parameters

speciation (λ) = 0.1

Pure birth process

Wright and
Warnock, in
review

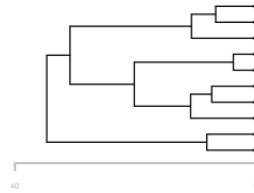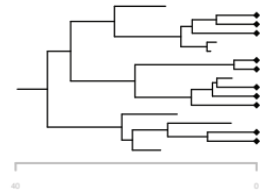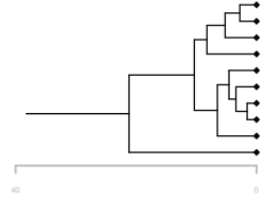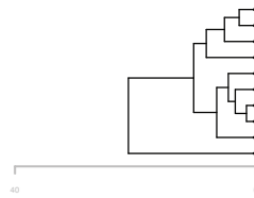The complete outcome of the diversification and sampling processes

The reconstructed tree

Model parameters

speciation ($\lambda$) = 0.1

Pure birth process

speciation ($\lambda$) = 0.1
extinction ($\mu$) = 0.05

Birth-death process

Wright and Warnock, in review

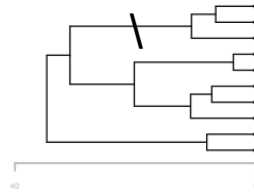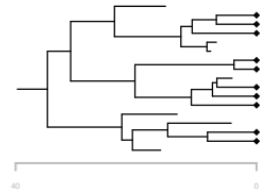The complete outcome of the diversification and sampling processes

The reconstructed tree
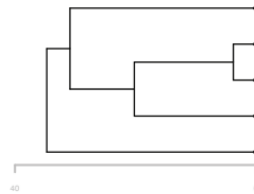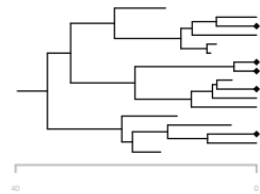
Model parameters

speciation (λ) = 0.1

Pure birth process

speciation (λ) = 0.1
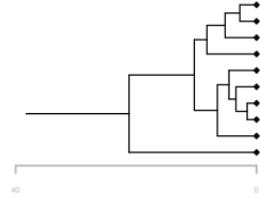extinction (μ) = 0.05

Birth-death process

speciation (λ) = 0.1
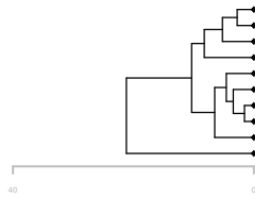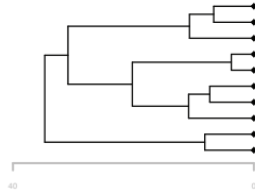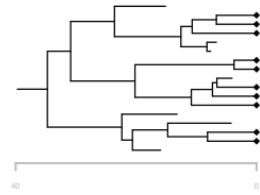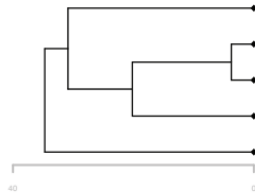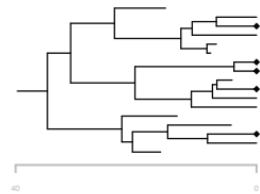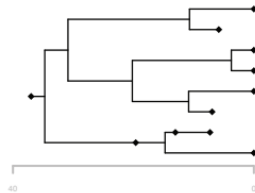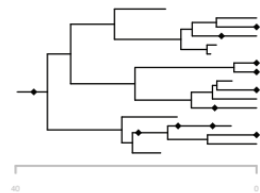extinction (μ) = 0.05
extant sampling (ρ) = 0.6

Birth-death sampling process

Wright and Warnock, in review

The complete outcome of the diversification and sampling processes | The reconstructed tree | Model parameters

speciation ($\lambda$) = 0.1

Pure birth process

speciation ($\lambda$) = 0.1
extinction ($\mu$) = 0.05

Birth-death process

speciation ($\lambda$) = 0.1
extinction ($\mu$) = 0.05
extant sampling ($\rho$) = 0.6

Birth-death sampling process

speciation ($\lambda$) = 0.1
extinction ($\mu$) = 0.05
extant sampling ($\rho$) = 0.6
fossil recovery ($\psi$) = 0.05

Fossilized birth-death process

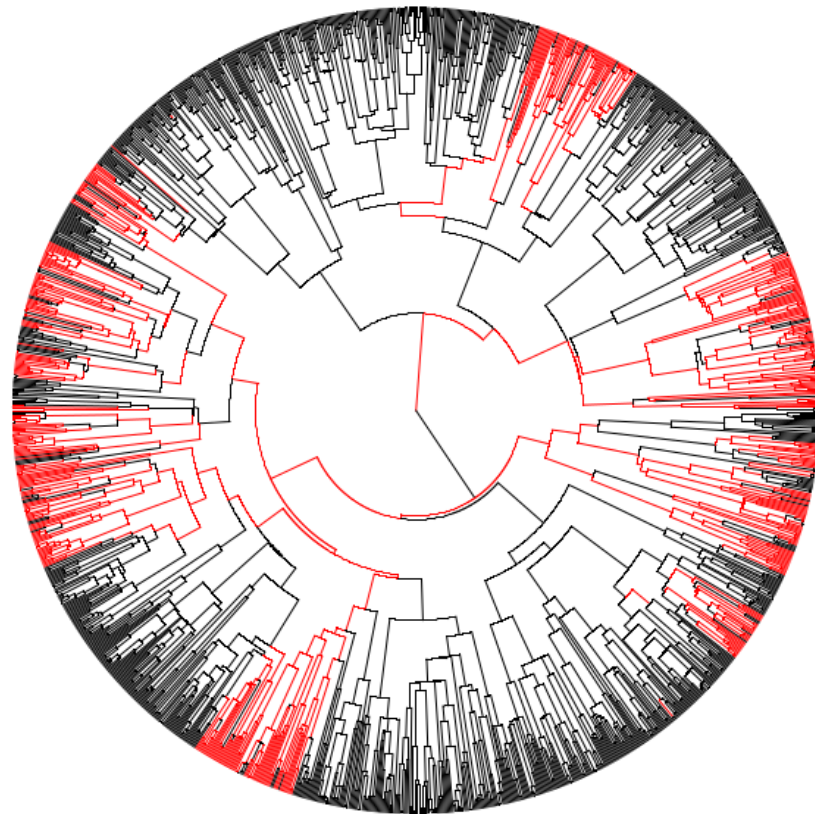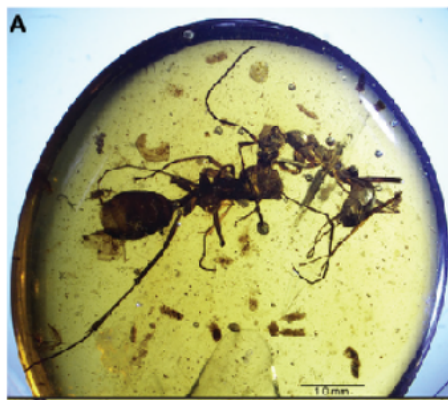Wright and Warnock, in review

# Why Ants?

Abundant molecular resources

666-tip multi-gene phylogeny from
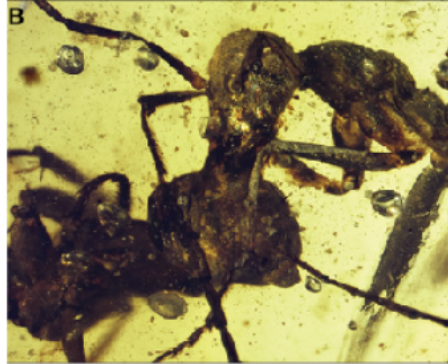Blanchard and Moreau (2017)

# Why Ants?

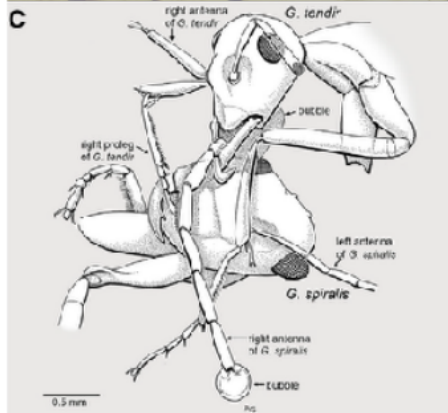# Why Ants?

Abundant morphological resources

Samples both with character data
and without character data

**Barden
&
Grimaldi, 2017**

# Why Ants?

# Why Ants?

# Why Ants?



# Vast, but biased resources

The data

0101...
1101...
0100...

Phylogenetic characters     Fossil ages

Tripartite model components

Substitution model     Clock model     Tree and tree model

Wright and Warnock, in review

# Systematic Conflict
## Molecular evidence supports each of these



Figure adapted from Borowiec et al 2019

# Systematic Conflict

## Molecular evidence supports each of these



## Morphology supports none of them

**Figure adapted from Borowiec et al 2019**

# Putting everything together

# Depending on the assumptions made about evolution, we recover support for every one of these topologies



**Figure adapted from Borowiec et al 2019**

If every hypothesis has some support, how do we know which is true?

© "Mike" Michael L. Baird  flickr.bairdphotos.com

Implied Gap

40
35
30
25

Wright and Lloyd, in review

Wright and Lloyd, in review

Not only can we support any
hypothesis we want, bad
hypotheses are often very close
to good ones

Implied Gap

Wright and Lloyd, in review

Accounting for Uneven Fossil Sampling

Accounting for Uneven Fossil Sampling

# Introduction to Posterior Prediction

## Assessing the fit of Normal distributions to trait data

Jeremy M. Brown and Christina L. Kolbmann

Last modified on October 10, 2019

# Training students to be data detectives

Computational Biology

## The why, when, and how of computing in biology classrooms [version 1; peer review: 1 approved]

April M. Wright [1], Rachel S. Schwartz [2], Jamie R. Oaks[3], Catherine E. Newman[4], Sarah P. Flanagan [5]

Author details

This article is included in the Bioinformatics Education and Training Collection collection.

# Training students to be data detectives

Computational Biology

## The why, when, and how of computing in biology classrooms [version 1; peer review: 1 approved]

April M. Wright [iD] [1], Rachel S. Schwartz [iD] [2], Jamie R. Oaks[3], Catherine E. Newman[4], Sarah P. Flanagan [iD] [5]

+ Author details

This article is included in the Bioinformatics Education and Training Collection collection.

Lessons learned in this class are being applied to our intro biology sequence

# Training students to be data detectives

# Training students to be data detectives

## Hands-On Learning with RevBayes

*Primary organizer:* Dr. April Wright, Southeastern Louisiana University

*Content:* This workshop will focus on using the phylogenetic estimation software RevBayes in an instructional setting. We will first introduce the graphical model framework used by the software. Graphical models can be used to introduce the fundamentals of probability, while also enabling transparent and flexible assembly of new phylogenetic models. Then, we will discuss using RevBayes in hands-on exercises for systematics and phylogenetics courses. Topics will include making use of the robust RevBayes tutorial library, tailoring pre-existing tutorials to your course, contributing your tutorials to the tutorial bank, using interactive computing notebooks, and integrating the software with R and Python.

# Training students to be data detectives

## Discrete morphology - Models and Tree Inference

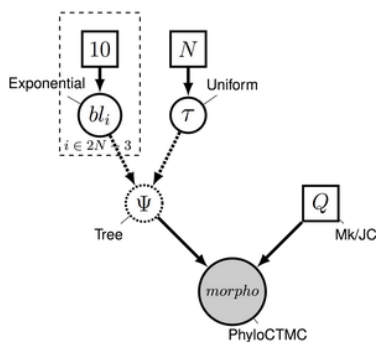April M. Wright

Source: `vignettes/module_05_TripartiteModel1_morph_change_models/05_RB_MCMC_Discrete_Morph.Rmd`

## Introduction to phylogenetic models of morphological evolution

Morphological data is commonly used for estimating phylogenetic trees from fossils. This tutorial will focus on estimating phylognetic trees from *discrete* characters, those characters which can be broken into non-overlapping character states. This type of data has been used for estimation of phylogenetic trees for many years. In the past twenty years, Bayesian methods for estimating phylogeny from this type of data have become increasingly common.

This tutorial will give an overview of common models and assumptions when estimating a tree from discrete morphological data. We will use a dataset from Zamora, Rahman, and Smith (2013). This dataset contains 23 extinct echinoderm taxa and 60 binary and multistate characters.

## Overview of Discrete Morphology Models

```
for (i in 1:n_branches) {
        bl[i] ~ dnExponential(10.0)
}
topology ~ dnUniformTopology(taxa)
psi := treeAssembly(topology, bl)


Q_morpho <- fnJC(2)

phyMorpho ~ dnPhyloCTMC( tree=psi,
Q=Q, type="Standard",
coding="variable" )

phyMorpho.clamp( morpho )
```

Graphical model showing the Mk model (left panel). Rev code specifying the Mk model is on the right-hand panel.

# Thank You!

- Jeremy Brown
- Christina Kolbmann
- Basanta Khakurel
- Courtney Grigsby
- Tyler Tran

- Rachel Warnock
- Tracy Heath
- Graeme Lloyd
- Corrie Moreau
- David Bapst

SOUTHEASTERN
LOUISIANA UNIVERSITY

LBRN

MS:

- Barido-Sottani J, Saupe E, Smiley TM, Soul, LC, Wright AM, Warnock RCM. In review. Seven rules for simulations in paleobiology.
- Wright AM, Lloyd, GT. In review. Bayesian analyses in phylogenetic paleontology: Interpreting the posterior sample. Preprint: https://github.com/graemetlloyd/ProjectWhalehead/blob/master/vignettes/MS.pdf
- Warnock RCM, Wright AM.  In review. Understanding the tripartite approach to Bayesian divergence time estimation. Preprint: https://www.overleaf.com/read/cbdxvgvxdkdq
- Wright AM, Schwartz RS, Oaks JM, Newman CM, and Flanagan SP. Accepted. The Why, When, and How of Computing in Biology Classrooms. Preprint: https://www.overleaf.com/read/cnnpbfzgzvfd
- Wright AM 2019. A systematist's guide to estimating Bayesian phylogenies from morphological data. Insect Systematics and Diversity 3: https://doi.org/10.1093/isd/ixz006.
- Wright AM 2019. treesiftr: An R package and server for viewing phylogenetic trees and data Journal of Open Source Education, 2(11), 35, https://doi.org/10.21105/jose.00035
- Devitt TJ, Wright AM, Cannatella DC, Hillis, DM. 2019. Species delimitation in endangered groundwater salamanders: Implications for aquifer management and biodiversity conservation. Proceedings of the National Academy of Sciences 116: 2624-2633

Courses:
- Biological Data Analysis: https://biologicaldataanalysis2019.github.io/2019/
- Introduction to Biodiversity Data Science: https://paleantology.github.io/GBIO153H/index.html
- Systematics: https://wrightaprilm.github.io/Systematics2020/index.html

Service:
Assistant Editor: Systematic Biology
Organizer, iEvoBio