

# Al buio non si trova

## Biostatistics in the 21st century

---

Luiz Max Carvalho

[lmax.fgv@gmail.com](mailto:lmax.fgv@gmail.com)

Available from: <https://github.com/maxbiostat/presentations/>



# Plan for today

## Music

A little metaphor to get us in the mood

## Problem I: using historical data to the fullest

Optimal Bayesian dynamic borrowing of information

## Problem II: dealing with huge complex data

MCMC in tree space: a journey through a strange land

## A football metaphor

Because why not include a second metaphor?

## A (useful?) methaphor

Che gelida manina (La Boheme, Puccini, 1896)

Che gelida manina, se la lasci riscaldar.

*What a frozen little hand, let me warm it for you.*

Cercar che giova? Al buio non si trova.

*What's the use of looking? We won't find it in the dark.*

Ma per fortuna, è una notte di luna,

*But luckily, it's a moonlit night,*

E qui la luna... l'abbiamo vicina.

*and the moon is near us here.*

Roberto Alagna & Leontina Vaduva, Paris, 1995.

# Le dirò con due parole, chi son

## Personal

- ⊙ Born and raised in Petrópolis-RJ;
- ⊙ Eldest of three kids;
- ⊙ Married and father of a daughter;
- ⊙ Mais Querido supporter.

## Academic

- ⊙ BSc in Microbiology & Immunology (UFRJ, 2012);
- ⊙ PhD Evolutionary Biology (Edinburgh, 2018);
- ⊙ Post doctoral researcher at ENSP/Fiocruz (2019);
- ⊙ Lecturer (Assistant Professor) at EMAP since Jan/2020.

## Applications of Statistics/Mathematics

Applications in Epidemiology, (Molecular) Biology, Ecology, Psychology, Linguistics, etc.

## Applied Statistics

Markov Chain Monte Carlo, Model combination and selection, Statistical Phylogenetics.

# Problem I: efficiently utilising available information

Loads of historical data: how to build informative priors?

Let  $\mathbf{y}_0 = (y_{01}, \dots, y_{0n_0})$  and  $\mathbf{y} = (y_1, \dots, y_n)$  be **historical** and **current** data, respectively.

Question: how do I build a prior that

- ⊙ Uses information in  $\mathbf{y}_0$  efficiently but also
- ⊙ Does not lead to borrowing too much information when the data sets are not compatible?

Applications: clinical trials, quality control, policy-making.

## Normalised power prior<sup>1</sup>

$$\tilde{\pi}(\theta, a_0 \mid \mathbf{y}_0) = \frac{L(\mathbf{y}_0 \mid \theta)^{a_0} \pi(\theta \mid \eta) \pi_A(a_0 \mid \phi)}{c(a_0; \eta, \phi)}$$

- ⊙ How pick  $\pi_A$  such that prediction error (say) is minimised?
- ⊙ How to **efficiently** compute

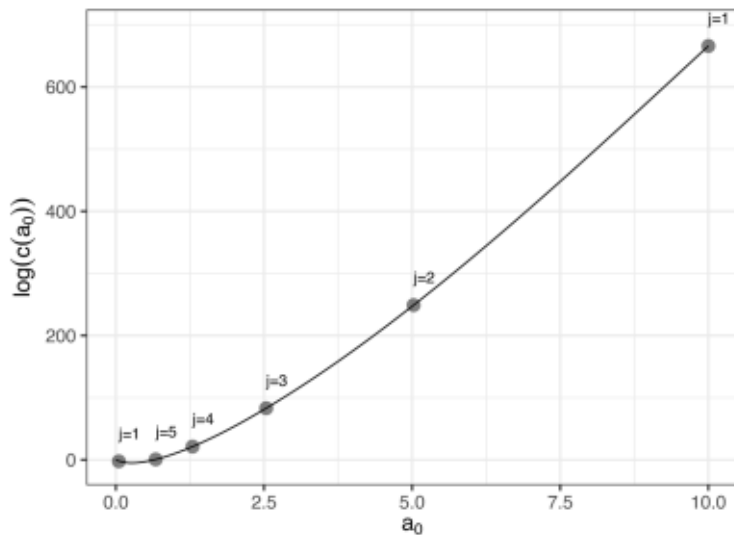
$$c(a_0; \eta, \phi) = \int_{\Theta} L(\mathbf{y}_0 \mid t)^{a_0} \pi(t \mid \eta) d\mu(t)$$

by leveraging its special properties as function of  $a_0$ ?

---

<sup>1</sup><https://doi.org/10.1002/sim.9124>

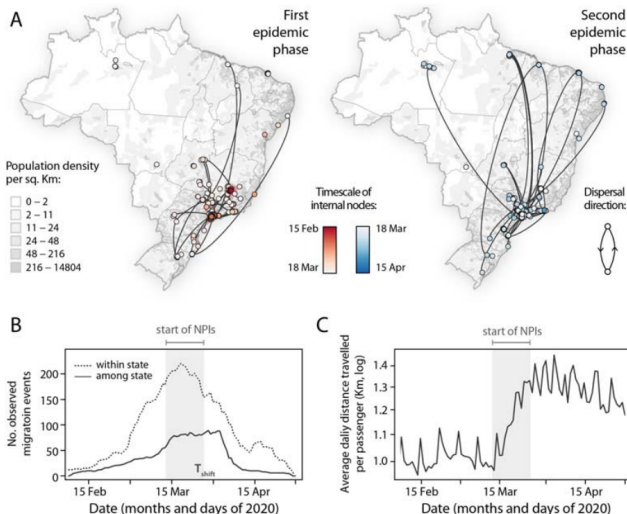
# Approximating the normalising constant





# Problem II: dealing with huge complex data

## Where did this virus come from?<sup>2</sup>



<sup>2</sup><https://doi.org/10.1126/science.abd2161>

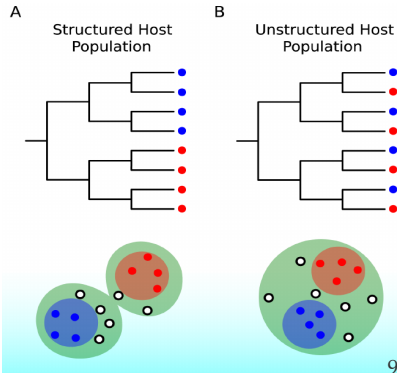
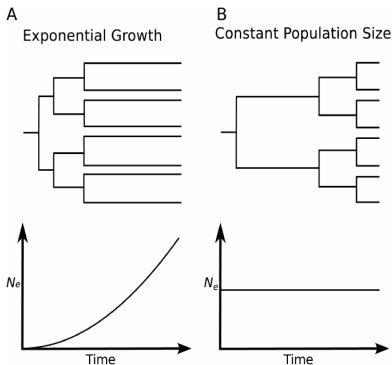
# Motivation

## Phylogenetics of fast-evolving viruses

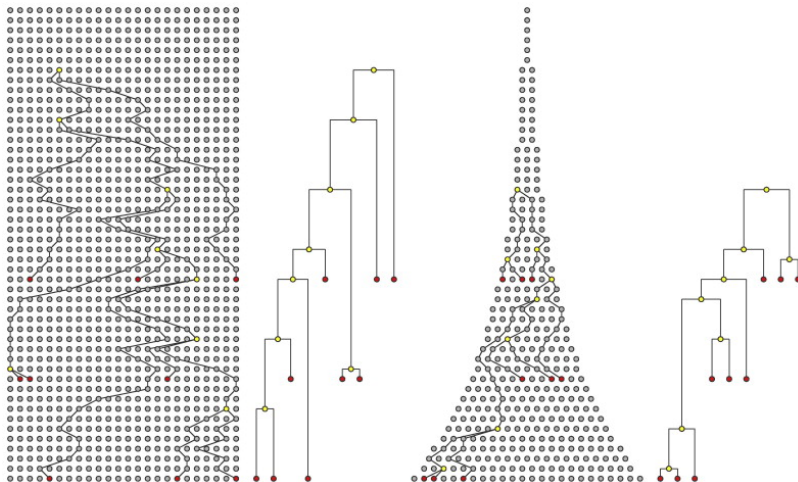
Inferring spatial and temporal dynamics from genomic data:

### Phylogenies\*!

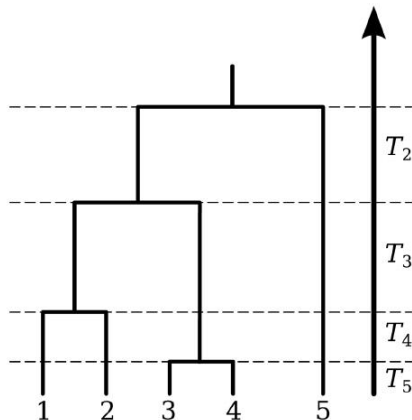
\* plus complicated models



# Trees and the coalescent



# Central object: time-calibrated trees



Let  $T_n$  denote the time for  $n$  lineages to *coalesce*, i.e., merge into one ancestral lineage, in a population of size  $N_e$ . Then:

$$\Pr(T_n = t) = \lambda_n e^{-\lambda_n t}$$

$$\lambda_n = \binom{n}{2} \frac{1}{N_e} = \binom{n}{2} \frac{1}{N_e \tau}$$

where  $N_e$  is the effective population size and  $\tau$  is the generation time. Let  $T_{\text{mrca}}$  denote the age of the most recent common ancestor:

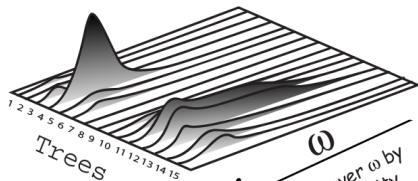
$$\begin{aligned}\mathbb{E}[T_{\text{mrca}}] &= \mathbb{E}[T_n] + \mathbb{E}[T_{n-1}] + \dots + \mathbb{E}[T_2] \\ &= 1/\lambda_n + 1/\lambda_{n-1} + \dots + 1/\lambda_2 \\ &= 2N_e \left(1 - \frac{1}{n}\right)\end{aligned}$$

Figure: Figure 4 from [Volz et al. \(2013\)](#).

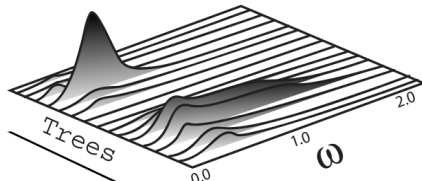
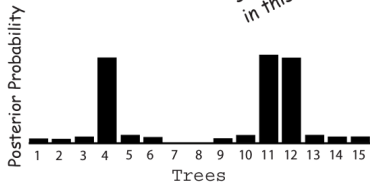
$$p(t, \mathbf{b}, \boldsymbol{\omega} | D) = \frac{f(D | t, \mathbf{b}, \boldsymbol{\omega}) \pi(t, \mathbf{b}, \boldsymbol{\omega})}{\sum_{t_i \in T_n} \int_B \int_{\Omega} f(D | t_i, \mathbf{b}_i, \boldsymbol{\omega}) \pi(t_i, \mathbf{b}_i, \boldsymbol{\omega}) d\boldsymbol{\omega} d\mathbf{b}_i} \quad (1)$$

- ⊙  $D$ : observed sequence (DNA) data;
- ⊙  $T_n$ : set of all binary ranked trees;
- ⊙  $\mathbf{b}_k$ : set of branch lengths of  $t_k \in T_n$  ( $\mathbb{R}_+^{2n-2}$ , kind of) ;
- ⊙  $\boldsymbol{\omega}$ : set of parameters of interest such as substitution model parameters, migration rates, heritability coefficients, etc.

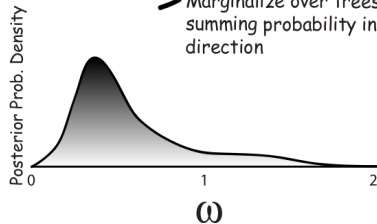
# The end product



← Marginalize over  $\omega$  by summing probability in this direction

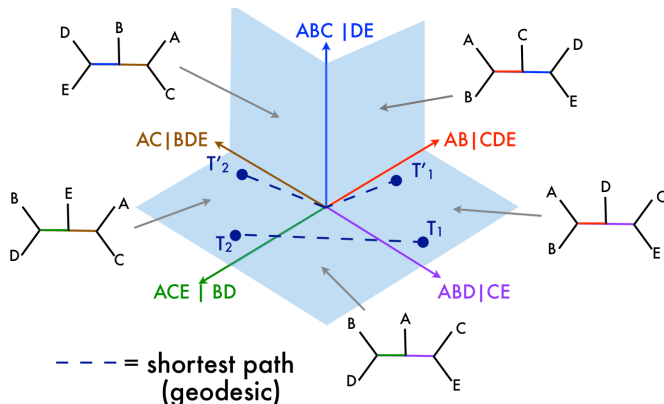


→ Marginalize over trees by summing probability in this direction



This place is weird..

## Traversing cubic complexes efficiently<sup>3</sup>



Applications: Molecular Epidemiology, Evolutionary Biology.

<sup>3</sup><https://youtu.be/h9bWRQ6aeKA>

# (Adaptive) Metropolis-Hastings for trees

General MH setup.

Let  $\tau = (t, \mathbf{b})$  denote a tree with topology  $t$  and branch lengths  $\mathbf{b}$ . For two trees  $\tau$  and  $\tau'$ , denote the transition kernel by  $q_\gamma(\tau|\tau') := \Pr(\tau' \rightarrow \tau|\gamma)$ .

Accepting with probability

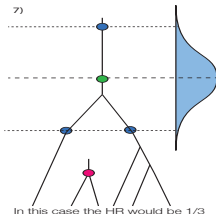
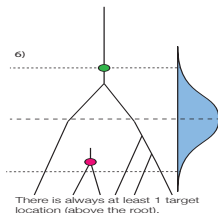
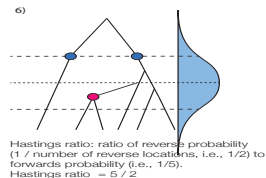
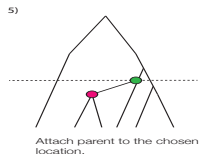
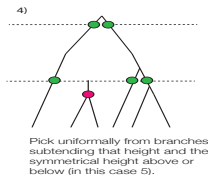
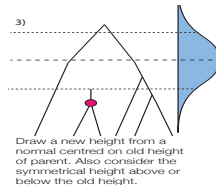
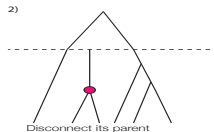
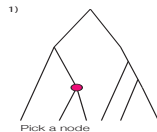
$$A_\gamma(\tau|\tau') = \min \left( 1, \frac{p(\tau', \boldsymbol{\omega}|D)q_\gamma(\tau|\tau')}{p(\tau, \boldsymbol{\omega}|D)q_\gamma(\tau'|\tau)} \right)$$

leads to the desired target.

**Note:** Here  $\gamma > 0$  is a so-called tuning parameter.



# STL – illustration



Carvalho (2019), Chapter 2.

## Remark

*Assume strictly positive branch lengths. Then SubTreeLeap induces an irreducible Markov chain on  $T_n$ .*

**Sketch:** Starting at  $x \in T_n$ , notice there exists  $\delta_y^\star > 0$  such that  $P(x \rightarrow y \mid \delta_y^\star) > 0$  for any tree  $y \in T_n$  in the SPR neighbourhood of  $x$ .

## Theorem

*Assume the target satisfies  $p(A) > 0$  for all  $A \subset \Psi$ . Then, SubTreeLeap induces an ergodic Markov chain on  $\Psi$ .*

**Sketch:** Employ the remark to get to the case where  $d_{\text{SPR}}(x, y) = 0$  and then establish Harris recurrence.

# Open problems in MCMC for phylogenies

Open problems:

- How can we construct more efficient proposals? How to exploit structure? **Geometry!**
- How to quantify exploration of the target? (Custom) **Tools!**
- Optimal scaling: what's the **optimal** acceptance probability?

## Another coat of golden paint?



"Why put another layer of gold paint on the Bentley when you are losing the engine?"

Zinedine Zidane, about Claude Makélélé leaving Real Madrid .

# Take home

## A light in the dark

Maths gives us methods with provable guarantees

## Computational methods are key

Learn to program and learn Computational Statistics<sup>4</sup>

## Maths works

Today we've employed: combinatorics, probability theory, basic calculus, optimisation and classical **and** Bayesian Statistics.

## We've got loads to do!

Biomedical statistics is where most of the cool data and problems are.

---

<sup>4</sup>Here's a place to start:

THE  
END