

Project 2: Augmenting functional information about human genes using probabilistic phylogenetic modeling

George G. Vega Yon
vegayon@usc.edu

Duncan Thomas Paul D. Thomas Paul Marjoram Huaiyu Mi John Morrison

Department of Preventive Medicine
University of Southern California

November 14th, 2018

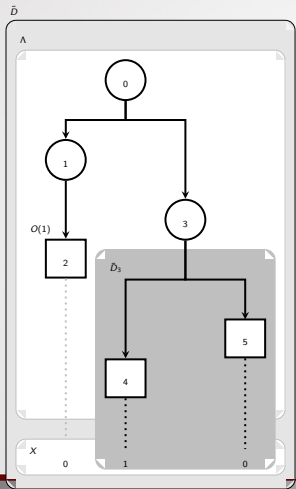
Recap: Model

1. A probabilistic model of gene function evolution,
2. The probability that the root node has the function is π ,
3. Conditional on its parent state, the probabilities that any given node has to either gain or lose a function are (μ_{01}, μ_{10}) ,
4. Finally, at the leaf node, the probability that a node with no function is mislabeled as having the function is ψ_{01} . Conversely, the probability that a node with a function is mislabeled as not having the function is ψ_{10} .

Parameter	Probability
π	The root node has the function
μ_{01}	Gaining a function
μ_{10}	Loosing a function
ψ_{01}	Mislabeling a 0
ψ_{10}	Mislabeling a 1

Table 1: Model parameters

Recap: Notation



Symbol	Description
$\Lambda \equiv (\mathcal{N}, \mathcal{E})$	Phylogenetic Tree.
$\mathbf{p}(n)$	Parent of node n .
$\mathbf{O}(n)$	Offspring of node n .
$\mathbf{X} \equiv \{x_n\}_{n \in \mathcal{N}}$	True annotations.
$\mathbf{Z} \equiv \{z_n\}_{n \in \mathcal{N}}$	Experimental annotations.
$D \equiv (\Lambda, \mathbf{X})$	Annotated Phylogenetic Tree.
$\tilde{D} \equiv (\Lambda, \mathbf{Z})$	Experimentally Annotated Phylogenetic Tree.
\tilde{D}_n	Induced Experimentally Annotated Subtree of node n .
\tilde{D}_n^c	Complement of \tilde{D}_n .

Table 2: Mathematical Notation

Changes from last year

From the formal (statistical) stand

- ▶ Prediction function: Right mathematical definition of the model prediction.
- ▶ New set of parameters: Propensity to report a finding.
- ▶ Flexible model specification: Definition of the likelihood function for different sets of parameters

By products generated during the implementation

- ▶ The `sLuRm` R package: A light-weight interface to slurm.
- ▶ Improvements on the `amcmc` R package, notably: automatic convergence.

Recap: The aphylo R package

Features:

- Provides a representation of *annotated* partially ordered trees.

Recap: The aphylo R package

Features:

- ▶ Provides a representation of *annotated* partially ordered trees.
- ▶ Interacts with the ape package (most used Phylogenetics R package with ~25K downloads/month)

Recap: The aphylo R package

Features:

- ▶ Provides a representation of *annotated* partially ordered trees.
- ▶ Interacts with the ape package (most used Phylogenetics R package with ~25K downloads/month)
- ▶ Implements the loglikelihood calculation of our model (with C++ under-the-hood).

Recap: The aphylo R package

Features:

- ▶ Provides a representation of *annotated* partially ordered trees.
- ▶ Interacts with the ape package (most used Phylogenetics R package with ~25K downloads/month)
- ▶ Implements the loglikelihood calculation of our model (with C++ under-the-hood).

Some new features

- ▶ Model specification via formula.
- ▶ Added the propensity to report discovery parameters.
- ▶ Two implementations of the prediction function (using a post-order algorithm as suggested by Prof. Suchard), and a brute force method... we use this for unit tests.
- ▶ (in the amcmc R package) Convergence monitoring and automatic stop of the MCMC algorithm

Nice visualizations

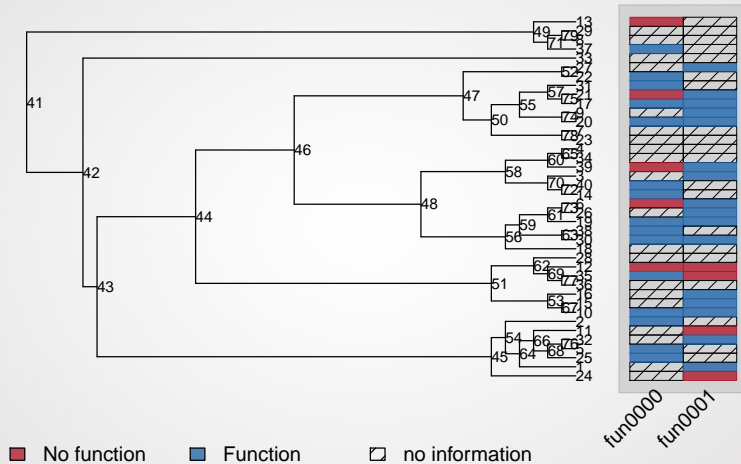


Figure 1: Annotated Phylogenetic Tree

$\text{Log } L(\psi_0, \psi_1, \mu_0, \mu_1, \Pi)$

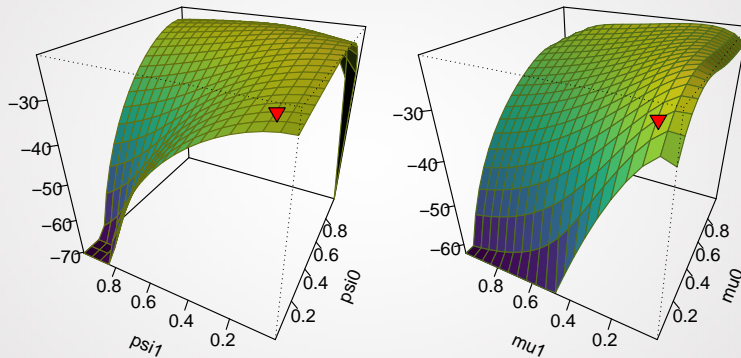


Figure 2: Surface of the likelihood of a given annotated tree.

Automatic specification of the likelihood function, e.g.

- ▶ $x \sim \mu$ baseline model
- ▶ $x \sim \mu + \psi + \Pi$ model including mislabeling and root node probabilities
- ▶ $x \sim \mu + \Pi$ same as before, but excluding mislabeling
- ▶ $x \sim \mu + \psi(1) + \Pi$ mislabeling of 1 is fixed
- ▶ $x \sim \mu + \psi(0, 1) + \Pi$ mislabeling of 0s and 1s is fixed

Flexible model specification

```
##  
## ESTIMATION OF ANNOTATED PHYLOGENETIC TREE  
##  
## Call: aphylo_mcmc(model = x ~ mu + psi + Pi, priors = bprior())  
## ll: -15.1028 ,  
## Method used: mcmc (748 iterations)  
## Leafs:  
## # of Functions 2  
##      Estimate Std. Err.  
## psi0    0.0998  0.0782  
## psi1    0.0955  0.0679  
## mu0     0.2379  0.0902  
## mu1     0.0499  0.0379  
## Pi      0.0888  0.0781
```



Results on the new specification (adventure)

The data generating process was $x \sim \mu + \psi + \eta + \Pi$ (η are the propensity to publication parameters).

Figure 4: Correct specification (includes 'eta')

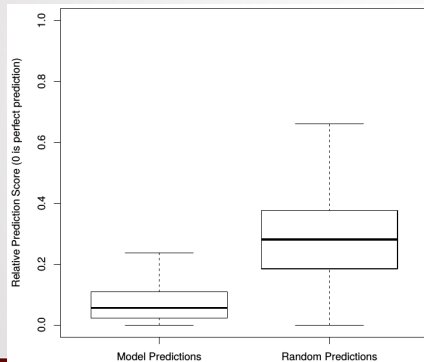
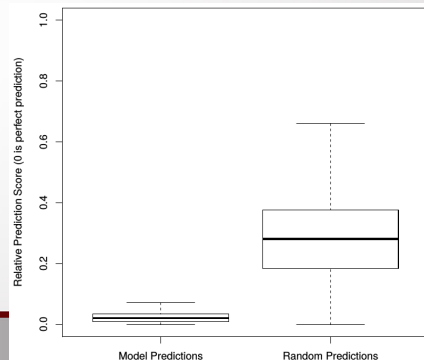


Figure 5: Miss specified model (does not include 'eta'). Missigness is confounded with propensity to fail to report



Concluding remarks

A parsimonious model of gene functions: easy to apply on a large scale (we already ran some simulations using all 13,000 trees from PantherDB... and it took us less than 1 week hour with 10 240 processors only).

Concluding remarks

A parsimonious model of gene functions: easy to apply on a large scale (we already ran some simulations using all 13,000 trees from PantherDB... and it took us less than 1 week hour with 10 240 processors only).

- ▶ Already implemented, we are currently in the stage of ~~writing the paper and setting up the simulation study~~ finishing and submitting the paper.

Concluding remarks

A parsimonious model of gene functions: easy to apply on a large scale (we already ran some simulations using all 13,000 trees from PantherDB... and it took us less than 1 week hour with 10 240 processors only).

- ▶ Already implemented, we are currently in the stage of ~~writing the paper and setting up the simulation study~~ finishing and submitting the paper.
- ▶ For the next steps, we are evaluating whether to include or how to include:

Concluding remarks

A parsimonious model of gene functions: easy to apply on a large scale (we already ran some simulations using all 13,000 trees from PantherDB... and it took us less than 1 week hour with 10 240 processors only).

- ▶ Already implemented, we are currently in the stage of ~~writing the paper and setting up the simulation study~~ finishing and submitting the paper.
- ▶ For the next steps, we are evaluating whether to include or how to include:
 - ▶ Type of node: speciation, duplication, horizontal transfer.
 - ▶ Branch lengths
 - ▶ Correlation structure between functions
 - ▶ ~~Using Taxon Constraints to improve predictions~~
 - ▶ Hierarchical model: Use fully annotated trees by curators as prior information.
- ▶ We are still unsure about how to procede with the software: R journal? Journal of Open Source Software? Journal of Statistical Software? Bioinformatics? etc.

Thank you!

Project 2: Augmenting functional information about human genes using probabilistic phylogenetic modeling

George G. Vega Yon
vegayon@usc.edu

Duncan Thomas Paul D. Thomas Paul Marjoram Huaiyu Mi John Morrison

Department of Preventive Medicine
University of Southern California

November 14th, 2018

13/1