

Phylogenetics: Likelihood

COMP 571
Luay Nakhleh, Rice University

1

The Problem

- * **Input:** Multiple alignment of a set S of sequences
- * **Output:** Tree T leaf-labeled with S

2

Assumptions

- * **Characters are mutually independent**
- * **Following a speciation event, characters continue to evolve independently**

3

4

- * The likelihood of model M given data D , denoted by $L(M|D)$, is $p(D|M)$.
- * For example, consider the following data D that result from tossing a coin 10 times:
 - * HTTTTHTTTT

5

- * Model M1:
 - * A fair coin ($p(H)=p(T)=0.5$)
 - * $L(M1|D)=p(D|M1)=0.5^{10}$

6

- * Model M2:
 - * A biased coin ($p(H)=0.8, p(T)=0.2$)
 - * $L(M2|D)=p(D|M2)=0.8^2 \cdot 0.2^8$

7

* **Model M3:**

* A biased coin ($p(H)=0.1, p(T)=0.9$)

* $L(M3|D)=p(D|M3)=0.1^2 \cdot 0.9^8$

8

* The problem of interest is to infer the model **M** from the (observed) data **D**.

9

* The maximum likelihood estimate, or **MLE**, is:

$$\hat{M} \leftarrow \operatorname{argmax}_M p(D|M)$$

10

- * $\mathcal{D} = \text{HTTTTHTTTT}$
- * **M1:** $p(H) = p(T) = 0.5$
- * **M2:** $p(H) = 0.8, p(T) = 0.2$
- * **M3:** $p(H) = 0.1, p(T) = 0.9$
- * **MLE (among the three models) is M3.**

11

- * **A more complex example:**
 - * **The model M is an HMM**
 - * **The data \mathcal{D} is a sequence of observations**
 - * **Baum-Welch is an algorithm for obtaining the MLE M from the data \mathcal{D}**

12

- * **The model parameters that we seek to learn can vary for the same data and model.**
- * **For example, in the case of HMMs:**
 - * **The parameters are the states, the transition and emission probabilities (no parameter values in the model are known)**
 - * **The parameters are the transition and emission probabilities (the states are known)**
 - * **The parameters are the transition probabilities (the states and emission probabilities are known)**

Back to Phylogenetic Trees

13

- * What are the data D ?
- * A multiple sequence alignment
- * (or, a matrix of taxa/characters)

Back to Phylogenetic Trees

14

- * What is the (generative) model M ?
- * The tree topology
- * The branch lengths
- * The model of evolution (JC, ..)

Back to Phylogenetic Trees

15

- * What is the (generative) model M ?
- * The tree topology, T
- * The branch lengths, λ
- * The model of evolution (JC, ..), E

Back to Phylogenetic Trees

16

- * The likelihood is $p(D|T, \lambda, E)$.
- * The MLE is

$$(\hat{T}, \hat{\lambda}, \hat{E}) \leftarrow \operatorname{argmax}_{(T, \lambda, E)} p(D|T, \lambda, E)$$

Back to Phylogenetic Trees

17

- * In practice, the model of evolution is estimated from the data first, and in the phylogenetic inference it is assumed to be known.
- * In this case, given D and E , the MLE is

$$(\hat{T}, \hat{\lambda}) \leftarrow \operatorname{argmax}_{(T, \lambda)} p(D|T, \lambda)$$

Assumptions

18

- * Characters are independent
- * Markov process: probability of a node having a given label depends only on the label of the parent node and branch length between them t

Maximum Likelihood

19

- * Input: a matrix D of taxa-characters
- * Output: tree T leaf-labeled by the set of taxa, and with branch lengths λ so as to maximize the likelihood $P(D|T,\lambda)$

$P(D|T,\lambda)$

20

$$\begin{aligned} P(D|T,\lambda) &= \prod_{site\ j} p(D_j|T,\lambda) \\ &= \prod_{site\ j} \left(\sum_R p(D_j, R|T,\lambda) \right) \\ &= \prod_{site\ j} \left(\sum_R \left[p(\text{root}) \cdot \prod_{edge\ u \rightarrow v} p_{u \rightarrow v}(t_{uv}) \right] \right) \end{aligned}$$

- * What is $p_{i \rightarrow j}(t_{uv})$ for a branch $u \rightarrow v$ in the tree, where i and j are the states of the site at nodes u and v , respectively?

21

22

- * For the Jukes-Cantor model with the parameter μ (the overall substitution rate), we have

$$p_{i \rightarrow j}(t) = \begin{cases} \frac{1}{4}(1 + 3e^{-4\mu t}) & i = j \\ \frac{1}{4}(1 - e^{-4\mu t}) & i \neq j \end{cases}$$

23

- * If branch lengths are measured in expected number of mutations per site, ν (for JC: $\nu = (\mu/4 + \mu/4 + \mu/4)t = (3/4)\mu t$)

$$p_{i \rightarrow j}(\nu) = \begin{cases} \frac{1}{4}(1 + 3e^{-4\nu/3}) & i = j \\ \frac{1}{4}(1 - e^{-4\nu/3}) & i \neq j \end{cases}$$

24

- * The ML problem is NP-hard (that is, finding the MLE (τ, λ) is very hard computationally)
- * Heuristics involve searching the tree space, while computing the likelihood of trees
- * Computing the likelihood of a leaf-labeled tree τ with branch lengths can be done efficiently using dynamic programming

P(D,T,λ)

Let $C_j(x,v)$ = P(subtree whose root is v | $v_j=x$)

Initialization: leaf v and state x $C_j(x,v) = \begin{cases} 1 & v_j = x \\ 0 & \text{otherwise} \end{cases}$

Recursion: node v with children u,w

$$C_j(x,v) = \left[\sum_y C_j(y,u) \cdot P_{x \rightarrow y}(t_{vu}) \right] \cdot \left[\sum_y C_j(y,w) \cdot P_{x \rightarrow y}(t_{vw}) \right]$$

Termination:

$$L = \prod_{j=1}^m \left[\sum_x C_j(x, \text{root}) \cdot P(x) \right]$$

25

Running Time

- * Takes time $O(nk^2m)$, where n is the number of leaves in the tree, m is the number of sites, and k is the maximum number of states per site (for DNA, $k=4$)

26

Unidentifiability of the Root

- * If the base substitution model is reversible (most of them are!), then rooting the same tree differently doesn't change the likelihood.

27

Questions?

28
