# Statiscal Properties of Phylogenetic Tree Construction Methods

Consistency, Long Branch Attraction; Robustness to Model Violations, ...

Of the 3 methods of tree construction we have learned about --

MP, distance methods, ML — which is best?

...

For sure, one wants a statistical estimator to be CONSISTENT

(informally. if your data is perfectly in accord with the model

[or method], the method X reconstructs the true tree.)

For the formal definition, suppose you have chosen a specific model

with parameters $M_0 = (T, N)$  $N =$ numerical parameters and your data

consists of site pattern frequencies computed from independent

trials of the experiment. Ie. $n$ sites were generated under $M_0$.

independently.

Focusing only on the tree (easy extension to $(T, N)$) for notational ease,

then let $\hat{T}_n$ denote the estimator from your method based on

a sample of size $n$. Suppose $\varepsilon > 0$ is arbitrary. Then if

$$\lim_{n \to \infty} \text{Prob}\left( \| \hat{T}_n - T \| < \varepsilon \right) = \underline{\qquad}$$

↑ probabilistic quantification of when $\hat{T}_n$ will be with $\varepsilon$ of $T$

let the sample size go to ∞

measurement of how close $\hat{T}_n$ is to the true value $T$

then $\hat{T}_n$ is a CONSISTENT ESTIMATOR of $T \quad = (T, N)$
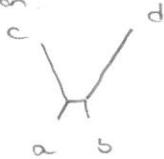
Formalization of a very _basic_ requirement for inference.

If $\hat{T}_{method}$ is not consistent, in practice no amount of data collection will help you to estimate $T$ with "method".
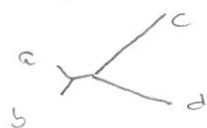
Which methods are consistent?

$\vdots$

Parsimony and Long Branch Attraction.

We know that a metric tree can be hard to infer

since $a, b$ are the closest, i.e. their DNA sequences should look the

most similar, but $a, c$ are sister. Indeed, NJ was introduced to

address this issue. We might expect parsimony to struggle for

such trees. For such a tree, a "method" might infer

a tree in which the long branches are attracted.

$$\equiv \text{Long Branch Attraction}$$

This phenomenon (LBA) extends to larger trees $n > 4$ and can

throw off inference if too remote an outgroup is contained in data.

Another way to view this is taxa $c, d$ are essentially independent

if those terminal branch lengths are long enough.

We will show that parsimony (MP) can be inconsistent

under a "2-state JC" model called the Cavender-Ferris-Neyman

CFN model. " Felsenstein Zone "

Details: Method:

Parsimony on a 4-taxon tree

Model: Explained below    CFN    2-state model

Data: Pattern frequencies    $xxyy$    $xyxy$    $xyyx$

↑    ↑    ↑
$n_1$    $n_2$    $n_3$

Counts from data

(order reversed in book.)

3-trees to choose from

$T_1$: $ab|cd$    $ps(T_1) = n_1 + 2n_2 + 2n_3 = 2n - n_1$

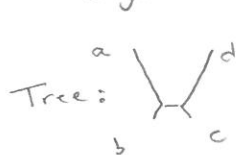$T_2$: $ac|bd$    $ps(T_2) = 2n_1 + n_2 + 2n_3 = 2n - n_2$
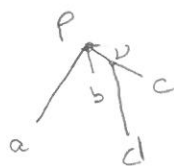
$T_3$: $ad|bc$    $ps(T_3) = 2n_1 + 2n_2 + n_3 = 2n - n_3$

where $n = $ # of informative sites

Parsimony Criterion: Choose $T_i$ with $n_i$ largest    (so $2n - n_i$ smallest)

End data analysis.

Begin: Generate sequences assuming the CFN model

Tree:  root →     Book: Here root at $a$

2-states.

Root Distribution is    $P_r = (.5, .5)$

2 Markov matrices one for short edges, one for long edges

$M_{long} = \begin{pmatrix} 1-q & q \\ q & 1-q \end{pmatrix}$    $M_{short} = \begin{pmatrix} 1-p & p \\ p & 1-p \end{pmatrix}$    $p, q \in (0, .5)$

With this model

$xxyy$      $P_1 = (1-q)^2 p(1-p)^2 + 2q(1-q)p(1-p)^2 + q^2 p^3$

$xyxy$      $P_2 = (1-q)^2 p^2(1-p) + 2q(1-q)p^2(1-p) + q^2(1-p)^3$

$xyyx$      $P_3 = (1-q)^2 p^3 + 2q(1-q)p(1-p)^2 + q^2 p(1-p)^2$

$\uparrow$

Work        (including HW)

If MP were to choose the true tree $T_1$, then since

$$\lim_{n \to \infty} \frac{n_i}{n} = P_i \qquad \text{it must be that} \qquad P_1 > P_2, P_3$$