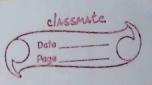
## 2023113019: Raunak Schratia 1TB: 12



A phylogenetic tree is a visual representation of the evolutionary
relationships arrong organisms inferred from
requences, The process of creating aconstituting
nee has many sups:
1. Sequence Data Collection: Gather Many
sequences from multiple organisms.
a. Sequence Alignment: Align the shipments
ummilaritus à différences.
3. Model Celection: Choose a suitable evolutionery model to
account for the sequence variations.
Parsimony (MP), Maximum likelihood (ML), or Baylsian
Inference (RT) to construct the phylogenetic tree:
5. True Evaluation: disess the tree's reliability: & accuretacy by
ming metrics such as bootstrap values or posterior probabilities
· I( a with ) flooring booking with the chariful
One metriod of thee construction, Bayerian inference shall be
discussed in much more detail here.
In a Bayerian analysis, inferences of phylogeny are based
ight the posterior probabilities of phylogenetic trees. The
justerior probability of the ith phylogenetic tree (+ Ti)

f(T; |X) = f(X|T[) f(T;) = f(X|T[) f(X|T;) f(T;)

(x/Z;)= (, f(x/Z;, v, 0) (v,0) dvdE



this summation is over all B(5) trees that are possible for a species [6(5) = (25-5)] { variotely & 25-3(5-3)!  $\frac{2(s) = (2s - 3)!}{2^{s-2}(s - 2)!} \left\{ \text{ nooted } \right\}, \text{ the integral is over all combinations of branch lengths } (v) & substitution parameters.}$ The prior for phylogenetic trees is  $f(T_i)$  & is assuably set to  $f(T_i) = 1$   $g(S_i)$ The prior on that branch lengths (v) & O me is denoted by of (v, o) Typically, the likelihood function [f(X[Ti, v,o]) is calculated under the assumption that substitution occur acc. to a time-homogeneous poisson process The summation & integrals see in a layerian analysis conner typically be walusted analytically. MRBAYES, the popular afterne that popularied by uses Markon Chain Monte Carlo (MCMC) to approximate the posterior probabilities of trees. meme to The basic MCMC algorithm works as follows: Is I new state for the chair is proposed very a

e. The acceptance probability for this new state is calculated. The acceptance probability is equal to the minimum of one or the likelihood nation. Welihood of new state to old state I times the prior ratio [prior

stochastic mechanism,

proposal nation (probability of proposing the dd elate to that of the new state). 3. A uniform pandom variable is drawn blu (0,1) of the no. is < acceptance probability, new state is

accepted, else rejected

This process is repealed thousands, or even millions of times. The proportion of the time any single tree is visited during the course of the chain is a valid approximation of its poeterior probability. Metropolis complet meme (mc3) is also & common alternative and in BT.

To summarise, there are 3 key concepts in BI:

Operor distribution. A probability distribution assigned to the phylogenetic tree, reflecting our mitial beliefs about @

1 Likelihood Function! Describes the probability of observing the molecular sequence dates given a fapology.

(E) Posterior Dutribution: Updated probability distribution

ets advantages are:

1 Robustness to Madel Hiolations , euch as relative branch length differences, & model violations.

Douantifies uncertainty in tree topology & beanch

3 at provides flexibility to accomplate various models of evolutionary processe:

CITATION: "MREAYES: Bayerian inference of phylogenetic trees" by John P. Huelsenbech & Frederik Ronquist (2001)