

Maximum likelihood :-

Likelihood (Model) = Probability (Data | Model)

Data → sequence alignment

Model Parameters → Nucleotide frequencies, Substitution rates
tree topology, branch length

Result → Maximum Likelihood estimation of

- Topology
- branch length
- Model Parameters
- Overall Likelihood

Substitution Model → $P(i \rightarrow j) = \begin{bmatrix} P_{AA} & P_{AC} & \dots & P_{AT} \\ P_{CA} & P_{CC} & \dots & P_{CT} \\ \vdots & \vdots & \ddots & \vdots \\ P_{TA} & P_{TC} & \dots & P_{TT} \end{bmatrix}$

$L(i) = \left[\frac{P_{AA}^{t_1}}{P_{AA}^{t_1}} + \frac{P_{AC}^{t_1}}{P_{AC}^{t_1}} + \dots + \frac{P_{AT}^{t_1}}{P_{AT}^{t_1}} \right]$

Likelihood = $L = L(1) \cdot L(2) \cdot L(3) \dots L(N)$

$\ln[L] = \ln[L(1)] + \ln[L(2)] + \ln[L(3)] \dots \ln[L(N)]$

Bayesian Inference Method:

$$P(w_i | D) = \frac{P(D | w_i) P(w_i)}{\sum_{j=1}^n P(D | w_j) P(w_j)}$$

computationally heavy on large data

Monte Carlo Markov chain (MCMC)

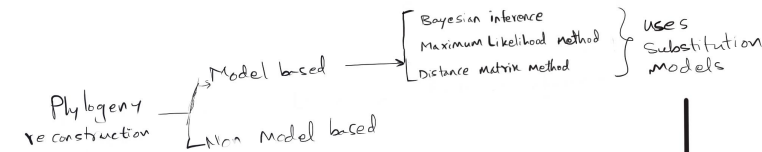


i) $\frac{P(w_i | D)}{P(w_{i-1} | D)} = \frac{P(D | w_i) P(w_i)}{P(D | w_{i-1}) P(w_{i-1})} \rightarrow \text{if } > 1 \text{ Accept } T^* \text{ as new Tree}$

ii) $\frac{P(w_i^* | D)}{P(w_i | D)} = \frac{P(D | w_i^*) P(w_i^*)}{P(D | w_i) P(w_i)} \rightarrow \text{if } > 1 \text{ Accept } w_i^* \text{ as new Parameters}$

iii) Repeat until Global/Local maximum is reached.

Extension: MCMC (MC³) Algorithm → Simultaneous multipoint Search



Maximum parsimony → Example:

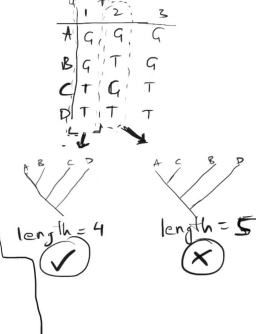
tree score: Max(Parsimony) ≈ Min(changes Mutations)

→ Constructs all possible trees for dataset

→ For each tree: Length = How many mutations are required?

→ Select the shortest length tree.

→ If more than 1 solution exists, they are equally good / equally parsimonious trees.



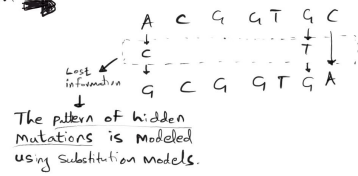
→ Fitch algorithm to simplify tree space search

→ other extensions

- PAUP
- MEGA (Branch and Bound)
- TNT

Nucleotide sequence data

Example:



$$D_{\text{new}} = \frac{3}{4} \ln \left(1 - \frac{4}{3} D_{\text{obs}} \right)$$

	A	B	C
A	D_{AA}	D_{AB}	D_{AC}
B	D_{BA}	D_{BB}	D_{BC}
C	D_{CA}	D_{CB}	D_{CC}

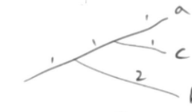
Jukes and Cantor Model used in Distance Matrix Method

	A	C	G	T
A	α	α	α	α
C	α	α	α	α
G	α	α	α	α
T	α	α	α	α

Relative rate matrix

$$P(t) = e^{Qt} = \begin{bmatrix} P_{AA} & P_{AC} & \dots & P_{AT} \\ P_{CA} & P_{CC} & \dots & P_{CT} \\ \vdots & \vdots & \ddots & \vdots \\ P_{TA} & P_{TC} & \dots & P_{TT} \end{bmatrix}$$

Probability that T will change to A given time 't'



$$Q = \sum_{i < j} (D_{ij} - d_{ij})^2$$

Minimise this function

least square optimality

→ Best tree = Smallest sum of squared errors

Minimum evolution optimality criteria

→ Best tree = Shortest tree

Neighbor Joining

- Clustering Algorithm
- Unique resultant tree

Other Models

- HKY85
- GTR

1 variable = α → Too few

2 variables = α, β → Sweet spot

3 variables = α, β, γ → Too many