

# Inferencia Filogenética

## Types of Data

Distances

Nucleotide sites

Clustering algorithm	UPGMA neighbour-joining	
	Minimum evolution	Maximum parsimony Maximum likelihood

Tree-building method

Optimality criterion

# Maximum likelihood

The best explanation for the observed outcome

$$\text{Prob}(D | H)$$

A maximum likelihood:

Likelihood de H

Model of sequence evolution

Tree - topology and branch lengths

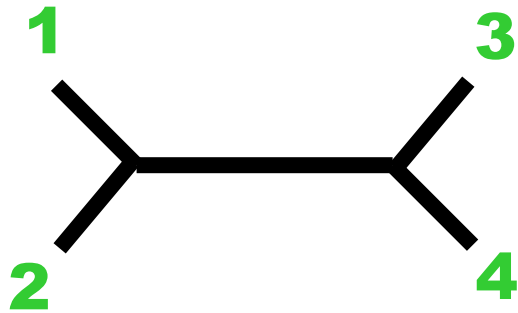
Observed data

- For a tree topology, what set of branch lengths makes the observed data most likely
- Which tree of all possible trees has the greatest likelihood

## Dois pressupostos:

- 1 – Evolution in different sites (on the given tree) is independent
- 2 – Evolution in different lineages is independent

# Exemplo simples: Dados, Modelo e Árvore



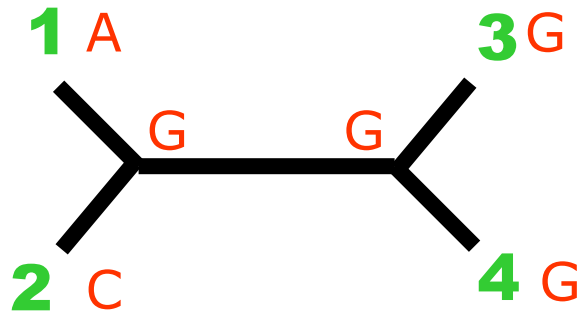
site 1 2 3 4 5

Species 1- **A T A T A**  
 Species 2- **A T C G C**  
 Species 3- **G C A G G**  
 Species 4- **G C C G G**

$$\begin{bmatrix}
 p_{AA} & \alpha & \alpha & \alpha \\
 \alpha & p_{CC} & \alpha & \alpha \\
 \alpha & \alpha & p_{GG} & \alpha \\
 \alpha & \alpha & \alpha & p_{TT}
 \end{bmatrix}$$

$p_{AA} = 1 - 3\alpha$

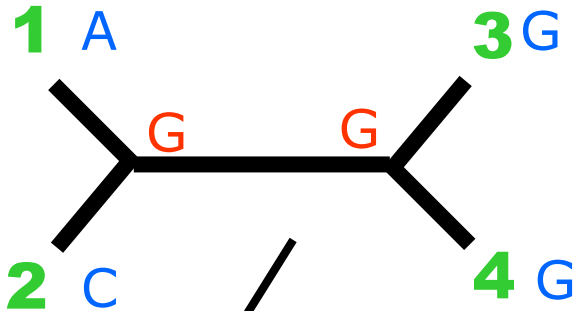
# Maximum likelihood



site	1	2	3	4	5
Species 1-	A	T	A	T	A
Species 2-	A	T	C	G	C
Species 3-	G	C	A	G	G
Species 4-	G	C	C	G	G

$$\alpha^2 \times (1 - 3\alpha)^3$$

Calcula-se a probabilidade de cada combinação para um dado site e somam-se todas



$$\alpha^2 \times (1 - 3\alpha)^3$$

G A

G C

G T

A G

A A

A C

A T

$$\alpha^3 \times (1 - 3\alpha)^2$$

C G

C A

C C

C T

T G

T A

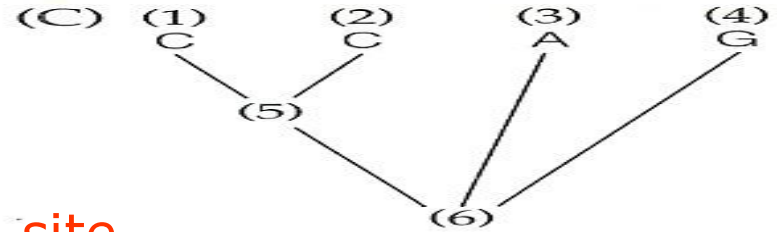
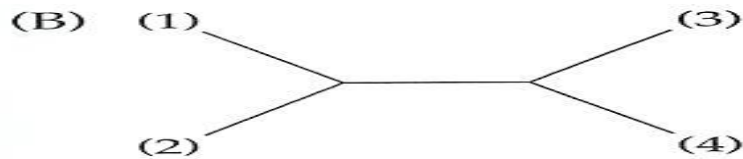
T C

T T

# Maximum likelihood

(A)

	1							$j$							$N$
(1)	C	...	G	G	A	C	A	C	G	T	T	T	A	...	0000
(2)	C	...	A	G	A	C	A	C	C	T	C	T	A	...	0000
(3)	C	...	G	G	A	T	A	A	G	T	T	A	A	...	0000
(4)	C	...	G	G	A	T	A	G	C	C	T	A	G	...	0000



Somam-se as combinações para cada site

(D)

$$L(j) = \text{Prob} \left( \begin{array}{c} C \quad C \quad A \quad G \\ \diagdown \quad \diagup \\ \quad A \quad \quad \quad \diagup \quad \diagdown \\ \quad \quad \quad \quad A \end{array} \right) + \text{Prob} \left( \begin{array}{c} C \quad C \quad A \quad G \\ \diagdown \quad \diagup \\ \quad C \quad \quad \quad \diagup \quad \diagdown \\ \quad \quad \quad \quad A \end{array} \right)$$

$$+ \dots + \text{Prob} \left( \begin{array}{c} C \quad C \quad A \quad G \\ \diagdown \quad \diagup \\ \quad G \quad \quad \quad \diagup \quad \diagdown \\ \quad \quad \quad \quad C \end{array} \right)$$

$$+ \dots + \text{Prob} \left( \begin{array}{c} C \quad C \quad A \quad G \\ \diagdown \quad \diagup \\ \quad T \quad \quad \quad \diagup \quad \diagdown \\ \quad \quad \quad \quad T \end{array} \right)$$

Repete-se para todos os sites

(E)

$$L = L(1) \cdot L(2) \cdot \dots \cdot L(N) = \prod_{j=1}^N L(j)$$

(F)

$$\ln L = \ln L(1) + \ln L(2) + \dots + \ln L(N) = \sum_{j=1}^N \ln L(j)$$

A soma dos log likelihoods dá-nos o log likelihood de uma dada árvore

# Maximum likelihood

Repete-se o cálculo para *todas* as árvores e escolhe-se aquela com maior probabilidade de ocorrer.

Searching methods para escolher as árvores que apresentam valores elevados de probabilidade