

1º Simpósio  **LAVIREO**

Modelos Evolutivos Teoria e Prática

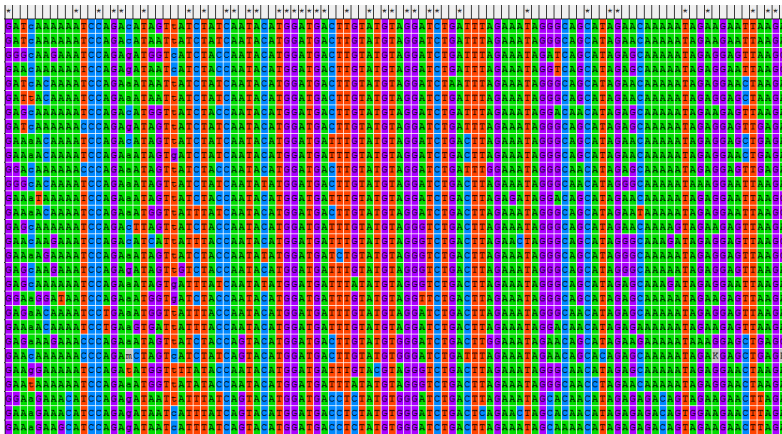
Edson Delatorre

Lab. de Genética Molecular de Microrganismos
Instituto Oswaldo Cruz/FIOCRUZ

delatorre.ioc@gmail.com

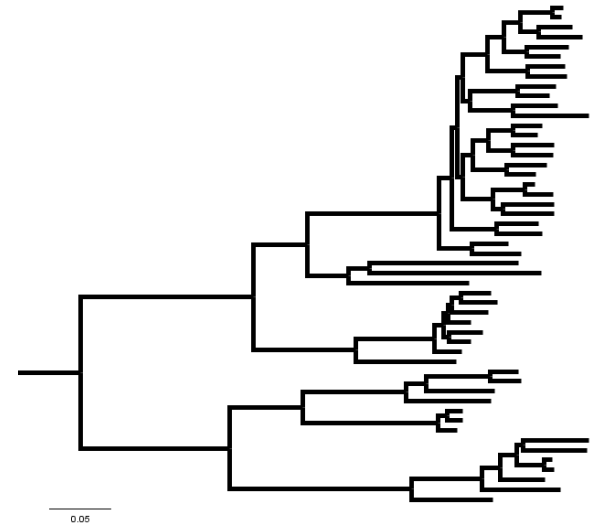
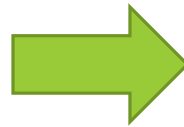
Rio de Janeiro – RJ, janeiro de 2019

Modelando a evolução



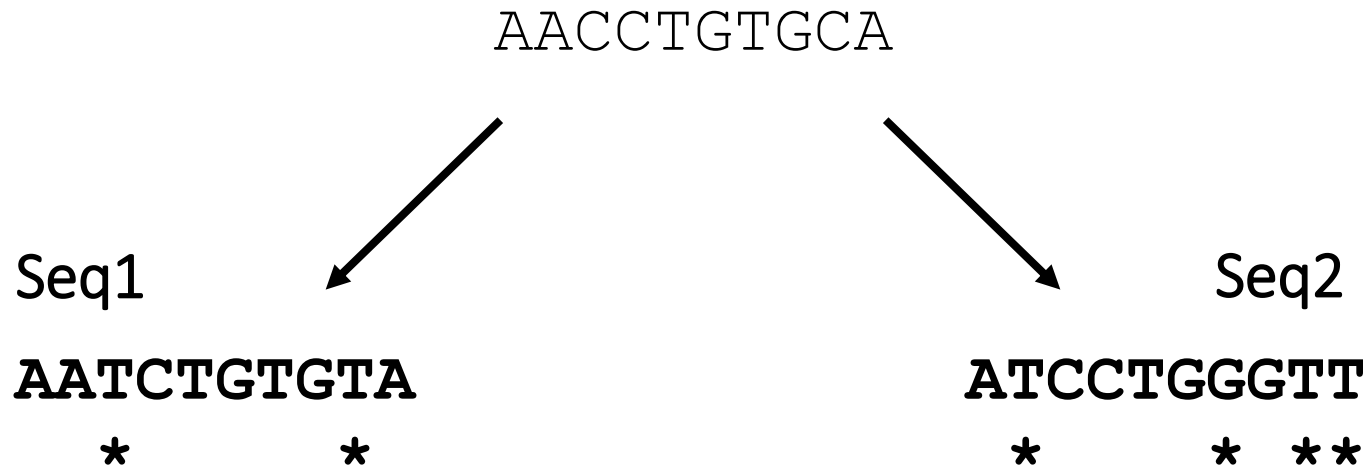
Alinhamento

Modelo



Árvore filogenética

Sequência ancestral



Seq1 **AATCTGTGTA**

Seq2 **ATCCTGGGTT**

* * * *

This block shows the two descendant sequences, Seq1 and Seq2, with their positions aligned. Seq1 is 'AATCTGTGTA' and Seq2 is 'ATCCTGGGTT'. The positions are aligned as follows: A (Seq1) under A (Seq2), A (Seq1) under T (Seq2), T (Seq1) under C (Seq2), C (Seq1) under T (Seq2), T (Seq1) under G (Seq2), G (Seq1) under G (Seq2), T (Seq1) under G (Seq2), and A (Seq1) under T (Seq2). Asterisks are placed under the positions where the sequences differ: positions 2, 4, 5, and 8.

p-distance:

número de diferenças nucleotídicas por sítio

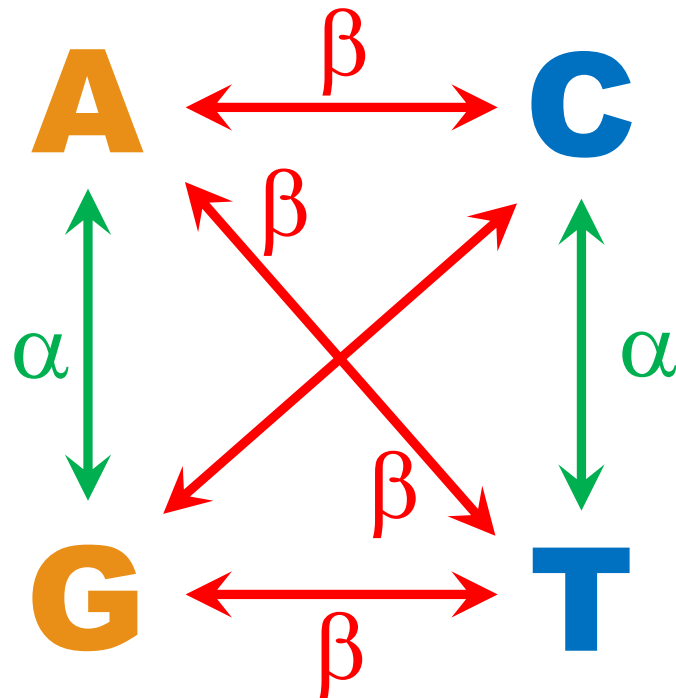
Seq1	AATCTGTGTA	$p\text{-distance} = 0.4$
Seq2	ATCCTGGGTT	
	** * *	

Normalmente subestima a “distância verdadeira”:

d - distância genética

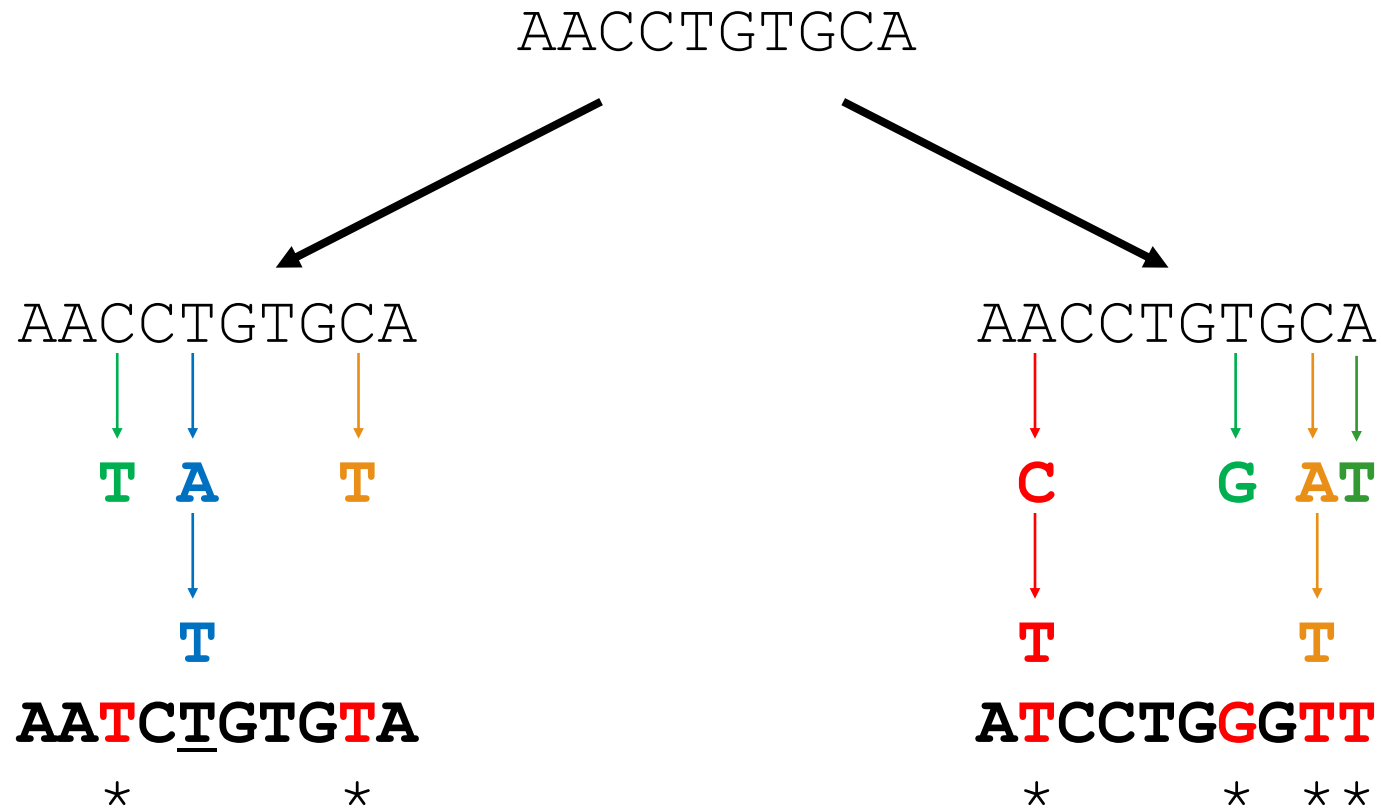
Transições (α): $Pu \leftrightarrow Pu$ (A,G) ou $Py \leftrightarrow Py$ (C,T)

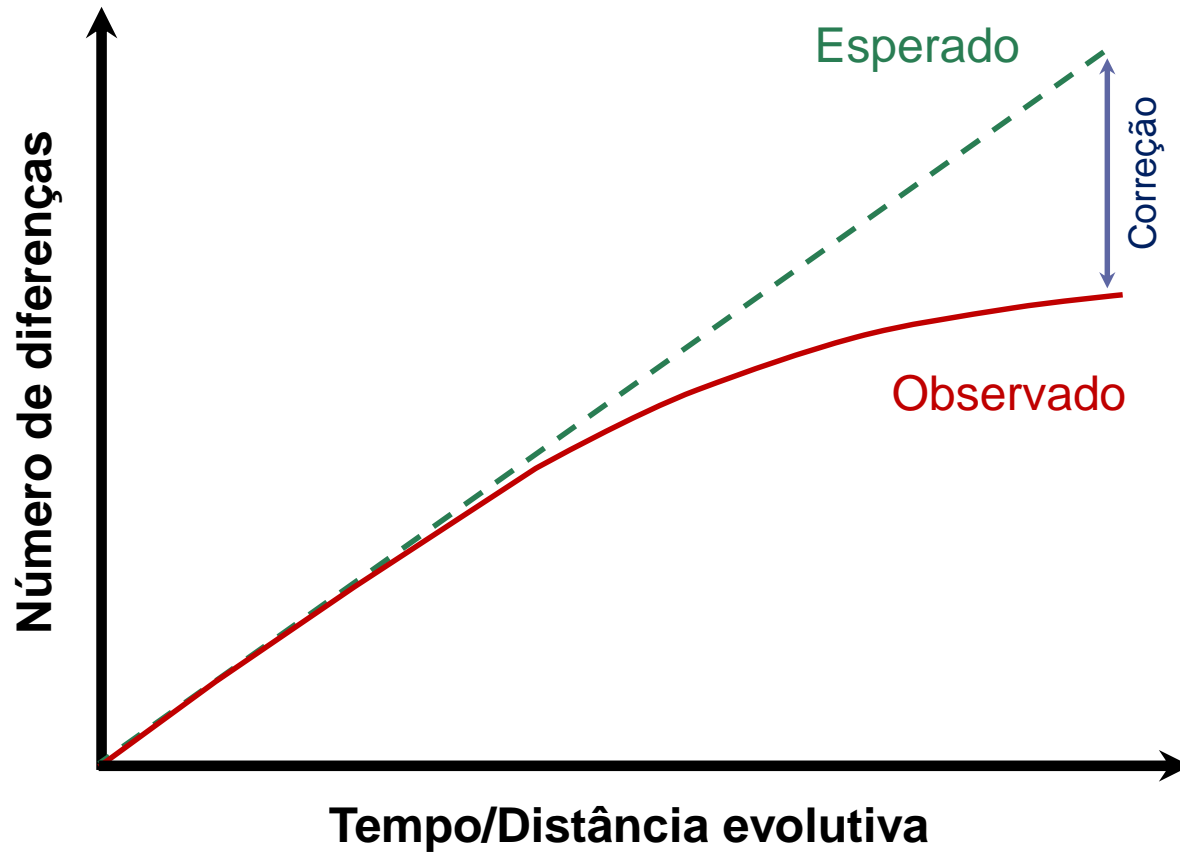
Transversões (β): $Pu \leftrightarrow Py$



Transições são pelo menos 2 vezes mais frequentes que **transversões**

Sequência ancestral



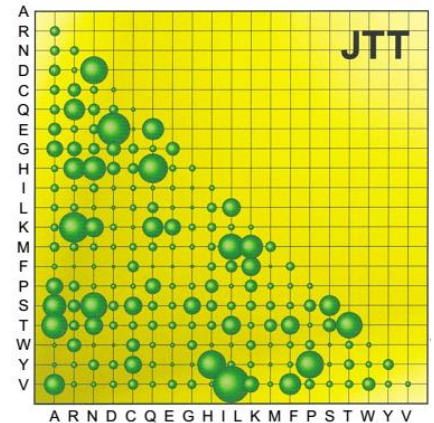


**Modelos
de
substituição**

Modelos evolutivos

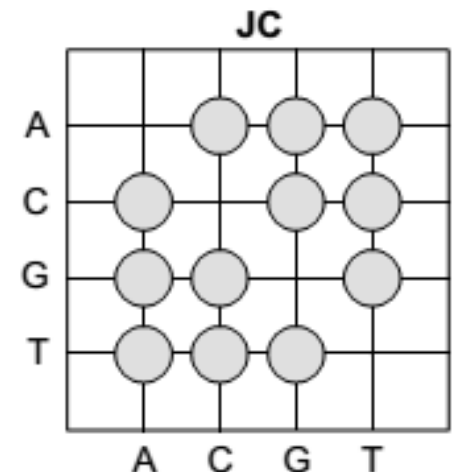
- Empíricos:

Propriedades calculadas através de comparações de um grande número de sequências observadas.



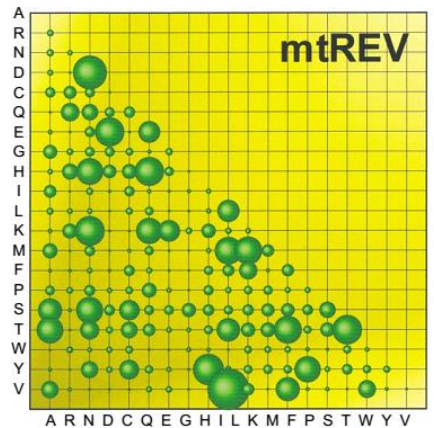
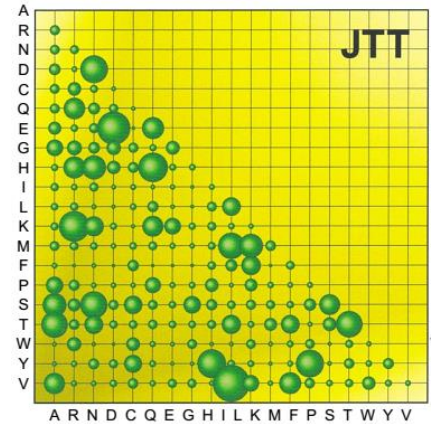
- Paramétricos:

Valores dos parâmetros são estimados a partir do dataset utilizado em cada análise particular



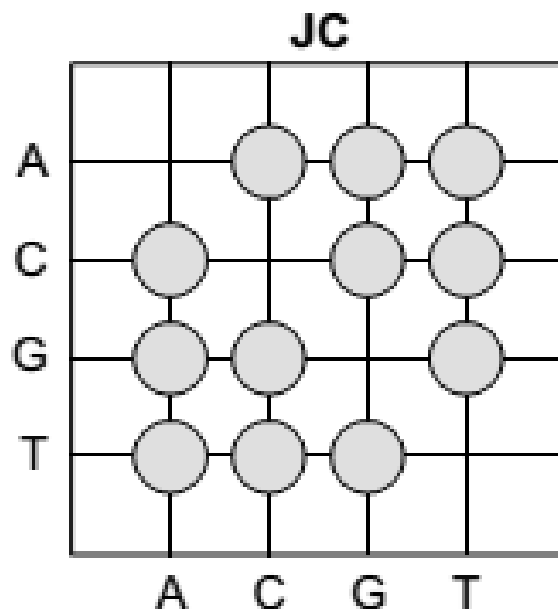
Modelos de substituição de aminoácidos

- Derivados da contagem simples de substituições de aminoácidos em grandes bancos de dados de sequências.
- Análise de sequências de proteínas intimamente relacionadas.
- Construção de matrizes específicas para grupos de proteínas.



Modelos de substituição de nucleotídeos

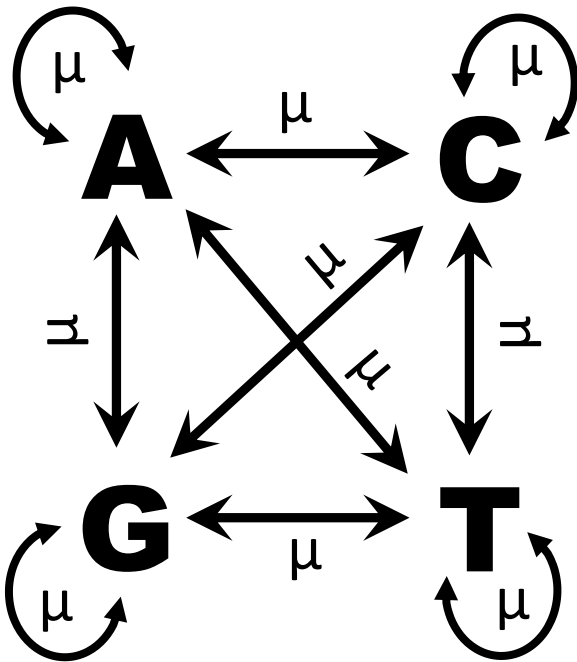
- Parâmetros
 - Frequência de bases
 - Taxa de substituição



$$\mathbf{P}_t = \begin{bmatrix} P_{AA} & P_{AC} & P_{AG} & P_{AT} \\ P_{CA} & P_{CC} & P_{CG} & P_{CT} \\ P_{GA} & P_{GC} & P_{GG} & P_{GT} \\ P_{TA} & P_{TC} & P_{TG} & P_{TT} \end{bmatrix}$$

Matriz

Modelo de Jukes e Cantor



Qualquer nucleotídeo i pode permanecer o mesmo

$i \longrightarrow i$ taxa = μ

ou mudar para qualquer dos outros nt ($j=3$)

$i \longrightarrow j$ taxa = $3 \cdot \mu$

Fórmula de Jukes e Cantor

$$P = \frac{3}{4} (1 - e^{-4/3 \mu t})$$

Resolvendo por μt

$$\mu t = - \frac{3}{4} \ln (1 - \frac{4}{3} P)$$

↑
Distância genética
estimada (**d**)

↑
Distância observada
P-distance

AACCTGTGCA

Seq1

AATCTGTGTA

*

*

Seq2

ATCCTGGGTT

*

*

**

Seq1 **AATCTGTGTA**

Seq2 **ATCCTGGGTT**

**

*

*

P-distance = **0.4**

$$d \text{ (JC)} = - 3/4 \ln [1 - 4/3 (0.4)] = \mathbf{0.5716}$$

Substituições de nt como processos de Markov homogêneos

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} \text{A} & \text{C} & \text{G} & \text{T} \end{matrix} \\ \begin{matrix} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{matrix} & \begin{pmatrix} -\mu(\textcolor{brown}{a}\pi_{\text{C}} + \textcolor{violet}{b}\pi_{\text{G}} + \textcolor{teal}{c}\pi_{\text{T}}) & \textcolor{brown}{a}\mu\pi_{\text{C}} & \textcolor{violet}{b}\mu\pi_{\text{G}} & \textcolor{teal}{c}\mu\pi_{\text{T}} \\ \textcolor{brown}{a}\mu\pi_{\text{A}} & -\mu(\textcolor{brown}{a}\pi_{\text{A}} + \textcolor{red}{d}\pi_{\text{G}} + \textcolor{teal}{e}\pi_{\text{T}}) & \textcolor{red}{d}\mu\pi_{\text{G}} & \textcolor{teal}{e}\mu\pi_{\text{T}} \\ \textcolor{violet}{b}\mu\pi_{\text{A}} & \textcolor{red}{d}\mu\pi_{\text{C}} & -\mu(\textcolor{violet}{b}\pi_{\text{A}} + \textcolor{red}{d}\pi_{\text{C}} + \textcolor{violet}{f}\pi_{\text{T}}) & \textcolor{violet}{f}\mu\pi_{\text{T}} \\ \textcolor{teal}{c}\mu\pi_{\text{A}} & \textcolor{teal}{e}\mu\pi_{\text{C}} & \textcolor{violet}{f}\mu\pi_{\text{G}} & -\mu(\textcolor{teal}{c}\pi_{\text{A}} + \textcolor{teal}{e}\pi_{\text{C}} + \textcolor{violet}{f}\pi_{\text{G}}) \end{pmatrix} \end{matrix}$$

- Premissas
- As taxas de mudança da base i para a base j são independentes da base que ocupou o sítio i anteriormente.
 - As taxas de substituição não mudam com o passar do tempo (homogeneidade).
 - As frequências relativas de A, C, G e T (π_{A} , π_{C} , π_{G} , π_{T}) estão em equilíbrio (estacionariedade).

Substituições de nt como processos de Markov homogêneos

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} \text{A} & \text{C} & \text{G} & \text{T} \end{matrix} \\ \begin{matrix} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{matrix} & \left(\begin{array}{cccc} -\mu(\text{a}\pi_{\text{C}} + \text{b}\pi_{\text{G}} + \text{c}\pi_{\text{T}}) & \text{a}\mu\pi_{\text{C}} & \text{b}\mu\pi_{\text{G}} & \text{c}\mu\pi_{\text{T}} \\ \text{a}\mu\pi_{\text{A}} & -\mu(\text{a}\pi_{\text{A}} + \text{d}\pi_{\text{G}} + \text{e}\pi_{\text{T}}) & \text{d}\mu\pi_{\text{G}} & \text{e}\mu\pi_{\text{T}} \\ \text{b}\mu\pi_{\text{A}} & \text{d}\mu\pi_{\text{C}} & -\mu(\text{b}\pi_{\text{A}} + \text{d}\pi_{\text{C}} + \text{f}\pi_{\text{T}}) & \text{f}\mu\pi_{\text{T}} \\ \text{c}\mu\pi_{\text{A}} & \text{e}\mu\pi_{\text{C}} & \text{f}\mu\pi_{\text{G}} & -\mu(\text{c}\pi_{\text{A}} + \text{e}\pi_{\text{C}} + \text{f}\pi_{\text{G}}) \end{array} \right) \end{matrix}$$

As taxas de substituição do nucleotídeo $i \rightarrow j = j \rightarrow i$

Em geral: $f = 1$

a, b, c, d, e são estimados a partir dos dados por ML.

Matriz Q para o modelo Jukes e Cantor (JC)

$$\text{Matriz de taxas} = \begin{pmatrix} - & \alpha & \alpha & \alpha \\ \alpha & - & \alpha & \alpha \\ \alpha & \alpha & - & \alpha \\ \alpha & \alpha & \alpha & - \end{pmatrix} \rightarrow a=b=c=d=e=f=1$$
$$\pi_A = \pi_T = \pi_C = \pi_G = 0.25$$

$$Q = \begin{pmatrix} -3/4\mu & 1/4\mu & 1/4\mu & 1/4\mu \\ 1/4\mu & -3/4\mu & 1/4\mu & 1/4\mu \\ 1/4\mu & 1/4\mu & -3/4\mu & 1/4\mu \\ 1/4\mu & 1/4\mu & 1/4\mu & -3/4\mu \end{pmatrix}$$

Modelos evolutivos

<u>Modelo</u>	<u>Matrix de taxas</u>	<u>Parâmetros</u>	<u>Param. livres</u>
K2P	$\begin{pmatrix} - & 1 & k & 1 \\ 1 & - & 1 & k \\ k & 1 & - & 1 \\ 1 & k & 1 & - \end{pmatrix}$	$\pi_A = \pi_C = \pi_G = \pi_T = 1/4$ Razão Ti/Tv $(b = e = k; a = c = d = f = 1)$	1
HKY85	$\begin{pmatrix} - & 1 & k & 1 \\ 1 & - & 1 & k \\ k & 1 & - & 1 \\ 1 & k & 1 & - \end{pmatrix}$	$\pi_A \neq \pi_C \neq \pi_G \neq \pi_T$ Razão Ti/Tv $(b = e = k; a = c = d = f = 1)$	4
TN93	$\begin{pmatrix} - & 1 & b & 1 \\ 1 & - & 1 & e \\ b & 1 & - & 1 \\ 1 & e & 1 & - \end{pmatrix}$	$\pi_A \neq \pi_C \neq \pi_G \neq \pi_T$ Razão Ti/Tv e razão Pu(Ti)/Py(Ti) $(b \neq e; a = c = d = f = 1)$	5
GTR	$\begin{pmatrix} - & a & b & c \\ a & - & d & e \\ b & d & - & 1 \\ c & e & 1 & - \end{pmatrix}$	$\pi_A \neq \pi_C \neq \pi_G \neq \pi_T$ $(a \neq b \neq c \neq d \neq e; f = 1)$	8

Heterogeneidade da taxa de evolução entre caracteres

- Todos os modelos assumem homogeneidade de taxas ao longo dos sítios
- Realidade:
 - Nt em diferentes posições dentro dos códonos possuem taxas diferentes (normalmente $3^\circ > 1^\circ > 2^\circ$)
 - Regiões hipervariáveis de proteínas (hotspots)

Subestimação das distâncias genéticas

A distribuição Gama

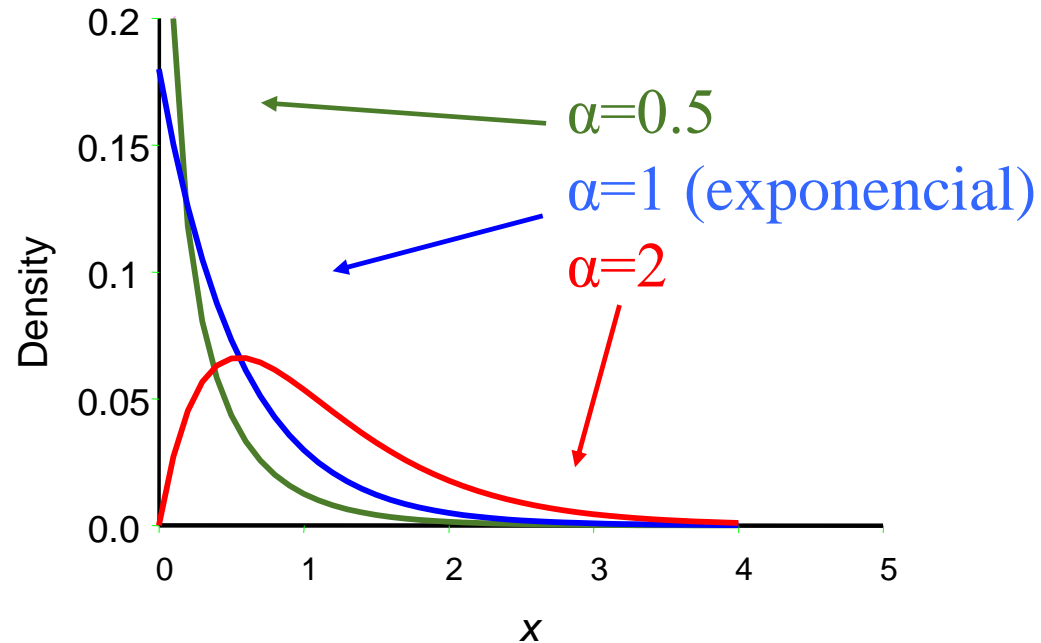
Função de densidade:

A forma da distribuição depende somente de um parâmetro *alpha*:

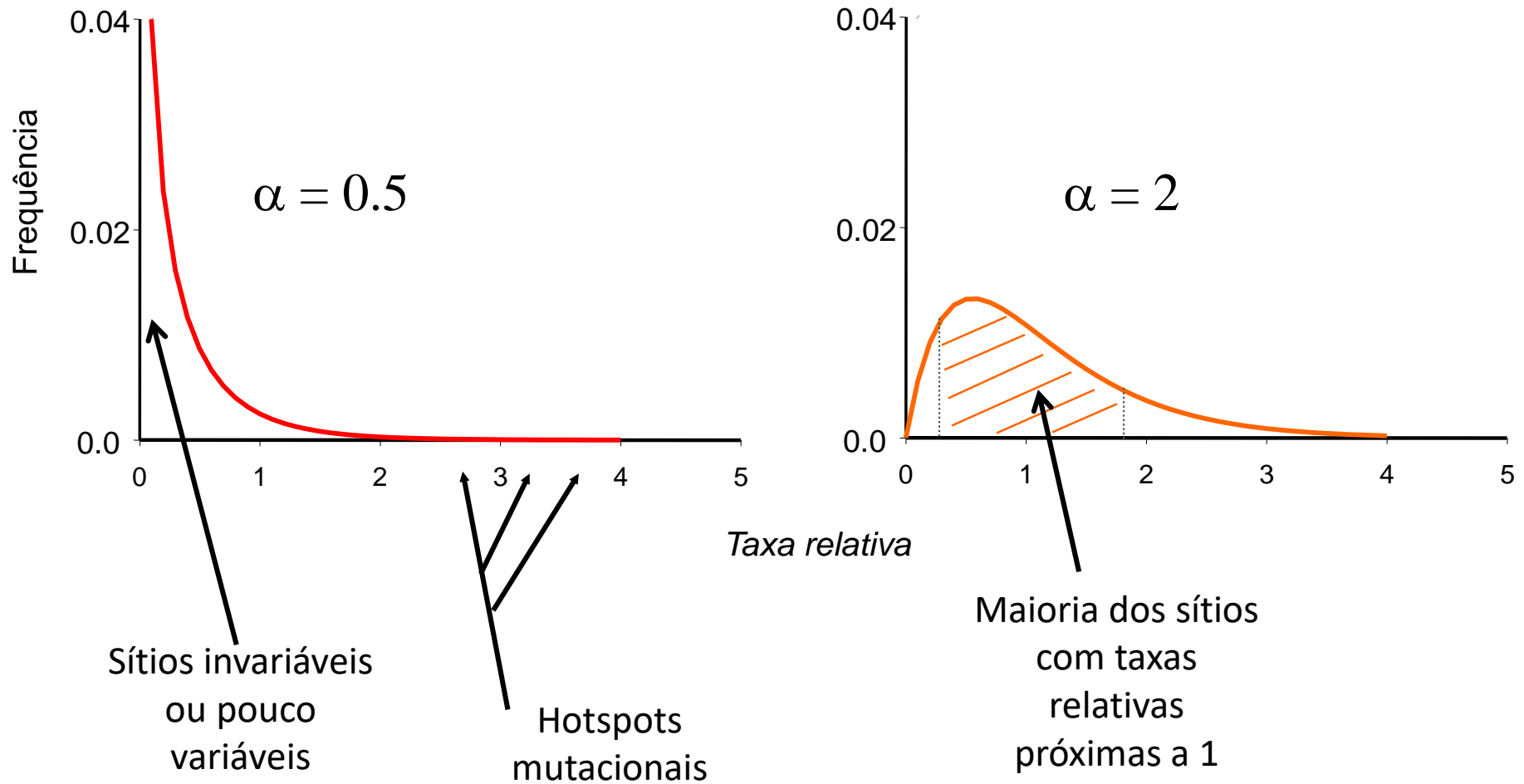
$\alpha < 1 \Rightarrow$ forma L

$\alpha > 1 \Rightarrow$ forma de sino (\sim Normal)

$\alpha \rightarrow \infty \Rightarrow$ um único valor no eixo x



Heterogeneidade de taxas através da distribuição Γ

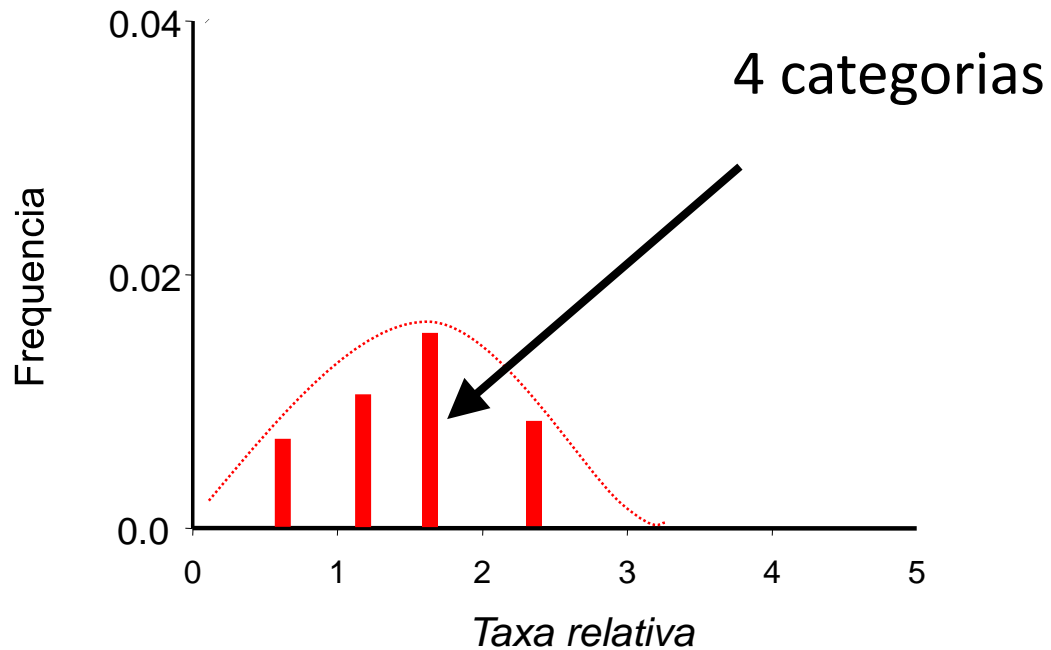


Modelo JC + Γ

$$d = \frac{3\alpha}{4} \left\{ \left[1 - \frac{4}{3} p \right]^{-1/\alpha} - 1 \right\}$$

A taxa de substituição varia de acordo com a distribuição Γ

Distribuição Γ discreta



- Em geral, quanto maior as categorias, melhor a aproximação.
- 4-8 categorias são normalmente suficientes.

Como escolher o melhor modelo?

- Modelos evolutivos são sempre simplificações, muitas vezes com premissas para tornar um problema complexo em algo passível de ser computado.
- Melhor modelo é aquele que se ajusta aos dados e permite fazer previsões acuradas.
- Em geral, quanto mais complexo o modelo, melhor o ajuste do mesmo aos dados.

Parâmetros

Acurácia



Teste da razão de verossimilhança (LRT)

- Próprio para a comparação de hipóteses evolutivas aninhadas

$$\Delta = 2 (\log_e L_1 - \log_e L_0)$$

L_1 – MV sob o modelo com mais parâmetros (complexo) → hipótese alternativa (H_1)

L_0 – MV sob o modelo com menos parâmetros (simples) → hipótese nula (H_0)

Δ se distribui como X^2 , com número de graus de liberdade igual a diferença entre os números de parâmetros livres entre os dois modelos.

Modelos

Equal base frequencies (3 df)

	JC	F81	K80	JC vs F81	HKY	SYM	GTR
Base frequencies	π	$\pi_A \pi_C \pi_G \pi_T$	π		$\pi_A \pi_C \pi_G \pi_T$	π	$\pi_A \pi_C \pi_G \pi_T$
Substitution rates	ρ	ρ	$\alpha\beta$		$\alpha\beta$	$\mu_1\mu_2\mu_3\mu_4\mu_5\mu_6$	$\mu_1\mu_2\mu_3\mu_4\mu_5\mu_6$

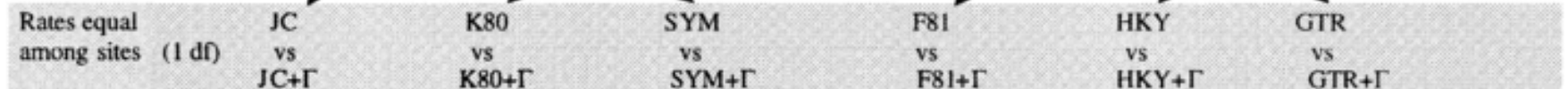
Transition rate equals
Transversion rate (1 df)



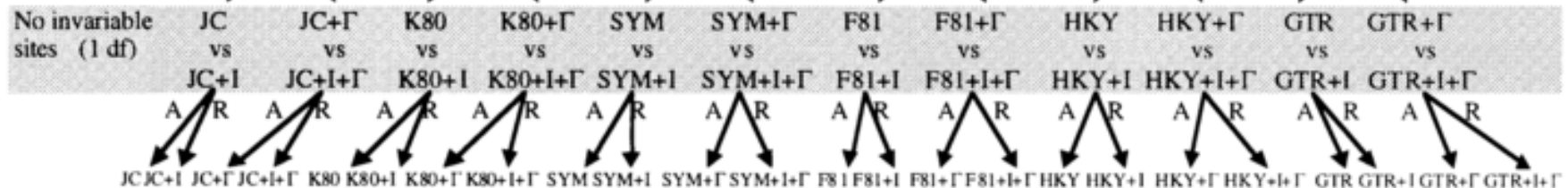
Equal transition rates and
Equal transversion rates (4 df)



Rates equal
among sites (1 df)



No invariable
sites (1 df)



Teste da razão de verossimilhança (LRT)

$$\Delta = 2 (\log_e L_1 - \log_e L_0)$$

- Se LR é significativo ($p < 0,05$ ou $< 0,01$): a inclusão de parâmetros adicionais no modelo alternativo aumenta significativamente a verossimilhança dos dados.
- Quando Δ é próximo a zero ($p > 0,05$): a hipótese alternativa não se ajusta aos dados significativamente melhor do que a hipótese nula.

Seleção de modelos com jmodeltest

- Seleção de modelos baseada na LRT hierárquica.
- Akaike Information Criterion

$$AIC = -2 \ln L + 2k$$

- Corrected AIC

$$AIC_c = AIC + 2k(k+1)/(N-k-1)$$

- Bayesian Information Criterion

$$BIC = -2 \ln L + k \ln N$$

L – verossimilhança do modelo, k – número de parâmetros estimáveis, N – tamanho da amostra