

①

$$P(\sigma_i | T) = P(T_1=A_1, T_2=C_1, T_3=G_1, T_4=G | T)$$

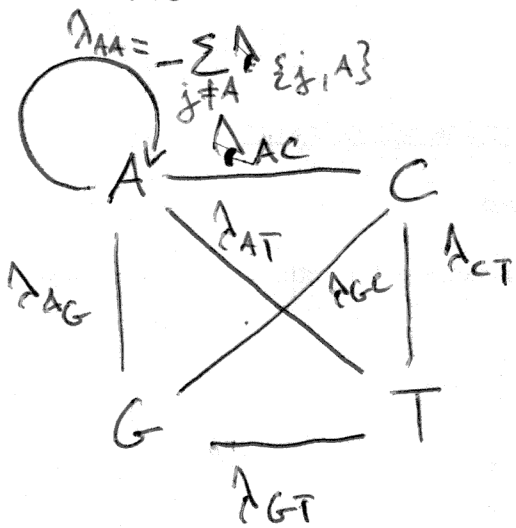
↑
site i

$$= \sum_{x \in X} \sum_{y \in X} \sum_{z \in X} P(T_1=A_1, T_2=C_1, T_3=G_1, T_4=G_1, u_3=z, u_2=y, u_1=x | T) \quad x \in \{A, C, G, T\}$$

$$= \sum_{x \in X} \sum_{y \in X} \sum_{z \in X} P(u_3=z) P(z \rightarrow x) P(z \rightarrow y) P(x \rightarrow A) P(x \rightarrow C) P(y \rightarrow G) P(y \rightarrow x)$$

↑ factor in the b_i
into all $P(\rightarrow)$

How can we get $P(x \rightarrow c | b_5)$



	A	C	G	T
A	*			
C	λ_{AC}	*		
G	λ_{AG}	λ_{CG}	*	
T	λ_{AT}	λ_{CT}	λ_{GT}	*

SYMMETRIC

note matrix R

equilibrium frequencies

$$\vec{\pi} = (\pi_A, \pi_C, \pi_G, \pi_T) \quad \sum_{x \in X} \pi_x \stackrel{!}{=} 1$$

$-\sum \lambda_{ii} \stackrel{!}{=} 1 \Rightarrow$ only 5 free parameters in the subst. matrix

$\vec{\pi}$ yields 3 additional free parameters

②

I JC: Jukes - Cantor model

$$JC: \begin{pmatrix} * & & & \\ a & * & & \\ a & a & * & \\ a & a & a & * \end{pmatrix} \rightarrow \pi = (0.25, 0.25, 0.25, 0.25)$$

$$\pi_i = \pi_j \quad \forall i, j$$

II Felsenstein 81 model: $\pi_i \neq \pi_j$

III Kimura 2 parameter model

$$\pi_i = \pi_j \quad \begin{pmatrix} * & & & \\ \beta & * & & \\ \alpha & \beta & * & \\ \beta & \alpha & \beta & * \end{pmatrix}$$

IV HKY 85 $\pi_i \neq \pi_j$ and Kimura 2 parameter model matrix

V GTR: General Time Reversible model

$$\begin{pmatrix} * & & & \\ \alpha & * & & \\ \beta & \gamma & * & \\ \delta & \varepsilon & \zeta & * \end{pmatrix}$$

time reversibility requires entries of a

$$\pi_i \cdot q_{ij} = \pi_j \cdot q_{ji}$$

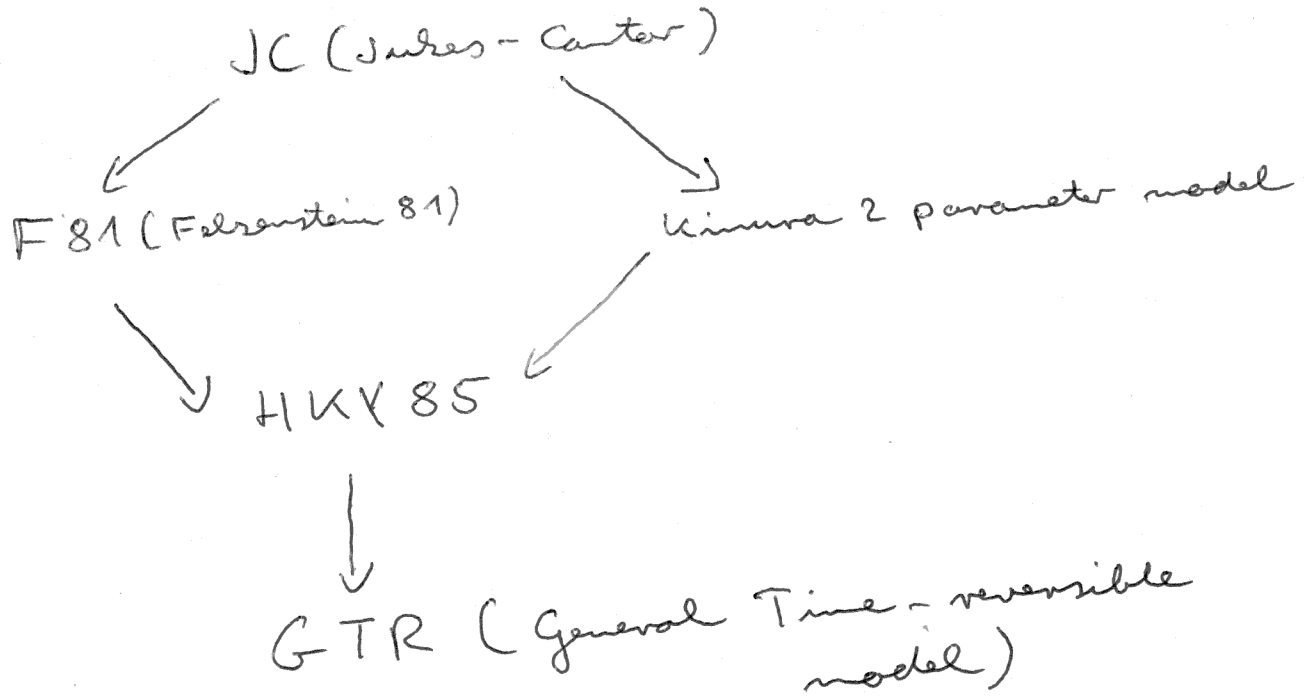
Q matrix

$$Q \text{ matrix: } Q = \text{diag}(\vec{\pi}) \cdot R$$

$$\begin{pmatrix} \pi_A & 0 & 0 & 0 \\ 0 & \pi_C & 0 & 0 \\ 0 & 0 & \pi_G & 0 \\ 0 & 0 & 0 & \pi_T \end{pmatrix}$$

2a

Model hierarchy and nested models



③

Q matrix of GTR model

$$Q = \begin{pmatrix} * & \pi_c \lambda_{ac} & \dots \\ \pi_a \lambda_{ac} & * & \dots \\ \pi_a \lambda_{ag} & & * & \dots \\ \pi_a \lambda_{at} & & & * \end{pmatrix}$$

$$\pi_a (\pi_c \lambda_{ac}) = \pi_c (\pi_a \lambda_{ac})$$

time reversibility holds

Likelihood ratio test

$$LHR \quad D = -2 \ln \left(\frac{LH(M_1)}{LH(M_2)} \right)$$

Model 1
 Model 2
 difference in #
 params. of M_1, M_2
 compare to χ^2 distribution

Only use this for nested models

AKAIKE INFORMATION CRITERION

$$AIC_i = -2 \ln(LH(M_i)) + 2 \cdot (p_i)$$

number of free parameters in model M_i

Models do not need to be nested
 Model i

$$BIC_i = -2 \ln(LH(M_i)) + p_i \cdot \ln(n)$$

↑
 size of the data
 # sites or
 # sites * # bases

④

Back to Likelihood calculations

$$LH = \sum \sum \sum p(x) P(x \rightarrow y | \theta, \mu)$$

model (Q-matrix)

↑
branch length

JC model $\lambda_{AC} = \lambda_{AT} = \dots = \lambda_{GT} = \alpha \cdot \frac{1}{3}$

$$\lambda_{AA}$$

$$\alpha \rightarrow \frac{4}{3} \alpha$$

$$\begin{array}{l} \frac{1}{3} \alpha \rightarrow C \\ \rightarrow G \\ \rightarrow T \end{array}$$

$$\begin{array}{l} \frac{1}{3} \alpha \rightarrow C \\ \rightarrow G \\ \rightarrow T \\ \rightarrow A \end{array}$$

Poisson distribution for calculating $P(t)$ from Q

$$k \text{ events in time } t = \frac{(at)^k}{k!} e^{-\frac{4}{3}at}$$

$$k=0 \rightarrow \text{no event}$$

$$e^{-\frac{4}{3}at} = P(\text{no event})$$

$$P(\text{some event}) = 1 - e^{-\frac{4}{3}at}$$

$$\text{for large } t \quad P(\text{some event}) = \frac{1}{4}$$

$$P(A \rightarrow C | t) = \frac{1}{4} (1 - e^{-\frac{4}{3}at})$$

Calculating $P(t)$ for a Q matrix of a GTR model \rightarrow see slides

in general we compute: $P(t) = e^{Qt}$

5

$$P_{ij}(t) = e^{Qt}$$

$$Q = U \Lambda U^{-1}$$

↑ remember: $\text{diag}(\vec{\pi}) \cdot R$

TRICK use Q'

$$Q' := \text{diag}(\sqrt{\pi'})^{-1} Q \text{diag}(\sqrt{\pi'})$$

$$Q' = U' \Lambda U'^T$$

↑ diagonal of Eigenvalues

$$\Leftrightarrow \text{diag}(\sqrt{\pi'})^{-1} Q \text{diag}(\sqrt{\pi'}) = U' \Lambda U'^T$$

$$\Leftrightarrow Q = \underbrace{(\text{diag}(\sqrt{\pi'}) U')}_{=U} \Lambda \underbrace{(U'^T \text{diag}(\sqrt{\pi'})^{-1})}_{=U^{-1}}$$

$$\Leftrightarrow Q = U \Lambda U^{-1}$$

$$P(t) = e^{Qt} = \sum_{k=0}^{\infty} \frac{1}{k!} (Qt)^k = \sum_{k=0}^{\infty} \frac{1}{k!} \underbrace{(U \Lambda U^{-1})^k}_{U \Lambda^k U^{-1}} \cdot t^k$$

$$\begin{aligned} & U \Lambda U^{-1} U \Lambda U^{-1} U \Lambda \\ & \quad \uparrow \quad \uparrow \\ & \text{Identity Matrix} \quad \text{Matrix} \\ & = U \Lambda^k U^{-1} \end{aligned}$$

$$= U \left(\sum_{k=0}^{\infty} \frac{1}{k!} \Lambda^k t^k \right) U^{-1}$$

remember: $\Lambda = \text{diag}(\Delta_i)$

$$\text{hence } P(t) = U e^{\text{diag}(\Delta_i) \cdot t} U^{-1}$$

6

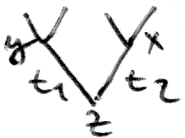
Protein substitution matrices
Dayhoff et. al.: PAM matrices
counting approach

Use only closely related sequences such that A → T → C does not happen!

$$\begin{pmatrix} * & 3 & 1 & 0 & \dots \\ \dots & * & 2 & 2 & \dots \\ \dots & \dots & \dots & 2 & \dots \\ \dots & \dots & \dots & \dots & * \end{pmatrix}$$

← each substitution is counted as one change

MAXIMUM LIKELIHOOD BRANCH LENGTH OPTIMIZATION



$$\sum \sum \sum \pi_z P(z \rightarrow y | t_1) \cdot P(z \rightarrow x | t_2) \cdot \underbrace{LH_1(y) \cdot LH_2(x)}$$

$$= \sum_x \sum_y \pi_y P(y \rightarrow x | t_1 + t_2) \cdot C(x, y) \cdot C(x, y)$$

remains constant

⇒ we can move z around the branch connecting x and y and will get the same score

Newton-Raphson procedure, second order Taylor-Approximation

$$\tilde{P}(t) = P(a) + P'(a)(t-a) + \frac{1}{2} P''(a)(t-a)^2 + \epsilon$$

$$= \underbrace{(\dots)}_{c_0} + \underbrace{(P'(a) - \frac{2}{2} P''(a) a)}_{c_1} t + \underbrace{(\frac{1}{2} P''(a) t^2)}_{c_2}$$

$$= c_0 + c_1 \cdot t + c_2 \cdot t^2$$

$$P' = c_1 + 2c_2 \cdot t \stackrel{!}{=} 0$$

minimization

→ NEXT:

$$\Leftrightarrow t = -\frac{c_1}{2c_2} = -\frac{P'(a) - P''(a) \cdot a}{2 \cdot \frac{1}{2} P''(a)} = \underbrace{a - \frac{P'(a)}{P''(a)}}_{\text{circle}}$$

⑦ Derivatives of the phylogenetic likelihood function:

$$P(t) = \sum_x \sum_y C_{(x,y)} \cdot P(y \rightarrow x | t)$$
$$= \sum_x \sum_y C_{(x,y)} U e^{\lambda_{xy} t} U^{-1}$$

\Rightarrow

$$P'(t) = \sum_x \sum_y C_{(x,y)} U \lambda_{xy} e^{\lambda_{xy} t} U^{-1}$$

$$P''(t) = \sum_x \sum_y C_{(x,y)} U \lambda_{yx}^2 e^{\lambda_{xy} t} U^{-1}$$