

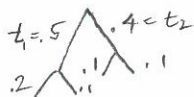
The continuous-time formulation:

parameters T, \vec{P} but replace 2) with

2a) a rate matrix $Q = \begin{pmatrix} q_{AA} & \dots & q_{AT} \\ \vdots & \ddots & \vdots \\ q_{TA} & \dots & q_{TT} \end{pmatrix}$

with non-negative
off diagonal entries
and row sum 0

2b) for each edge e in the tree, a non-negative branch length t_e .



The Markov transition matrix $M_e = e^{Qt_e} = (e^Q)^{t_e}$

An important assumption in the continuous time formulation is a Common Rate Matrix Q for all edges of the tree.

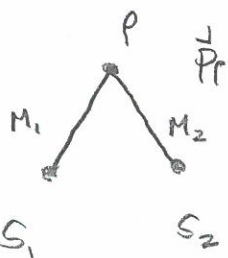
III) Computing the expected frequency arrays.

i.e. the $P(\text{leaves show pattern } i, j \dots \text{ in at leaves})$

→ Refer to page 10.5

Eg. 2-edge tree:

Parameters $(T, \vec{P}, \{M_1, M_2\})$



Let P be the expected joint frequency array

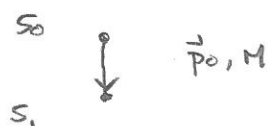
P is 4×4 with entries

$$P(i, j) = P(S_1 = i, S_2 = j) = \text{expected value of seeing pattern } \begin{matrix} i \\ j \end{matrix}$$

$S_1: i$

$S_2: j$

The joint distribution $P(i,j)$ for a 1-edge tree



$$P = \text{diag}(\vec{p}_0) M$$

Since

$$= \begin{pmatrix} P_A & 0 & 0 & 0 \\ 0 & P_G & 0 & 0 \\ 0 & 0 & P_C & 0 \\ 0 & 0 & 0 & P_T \end{pmatrix} \begin{pmatrix} P(s_1=A | s_0=A) & P(s_1=G | s_0=A) & P(s_1=C | s_0=A) & P(s_1=T | s_0=A) \\ P(s_1=A | s_0=G) & & & \\ P(s_1=A | s_0=C) & & \ddots & \\ P(s_1=A | s_0=T) & & & \end{pmatrix}$$

$$= \begin{pmatrix} P(s_0=A, s_1=A) & P(s_0=A, s_1=G) & P(s_0=A, s_1=C) & P(s_0=A, s_1=T) \\ P(s_0=G, s_1=A) & & & \\ P(s_0=C, s_1=A) & & \ddots & \\ P(s_0=T, s_1=A) & & & P(s_0=T, s_1=T) \end{pmatrix}$$

↑
P!

Two review items from linear algebra:

Matrix Transpose M^T

Eigenvalues and Eigenvectors: Suppose M is an $n \times n$ matrix, $\lambda \neq 0$, \vec{v} is $n \times 1$ vector, and

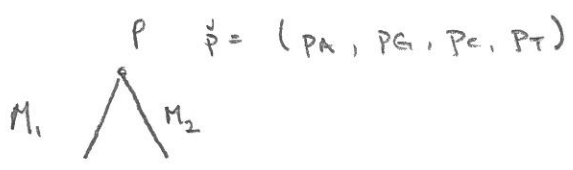
$$M\vec{v} = \lambda\vec{v}$$

↑ ↑
eigenvector eigenvalue

Eigenvector

Left eigenvectors ...

To do this, sum over all possible states at the root p



- A=1
- G=2
- C=3
- T=4

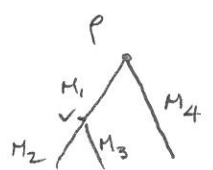
The probability of pattern $\begin{smallmatrix} i \\ j \end{smallmatrix}$ at the leaves is

$$P(s_1=i, s_2=j) = P(i,j) = \sum_{k=1}^4 p_k \underbrace{M_1(k,i) M_2(k,j)}_{\begin{smallmatrix} k \\ i \quad j \end{smallmatrix}}$$

Exercise for Math Students:

Show $P = M_1^T \text{diag}(\vec{p}_p) M_2$ for this 2-edge tree.

Eg. 3-edge tree:



Parameters $(T, \vec{p}_p, \{M_1, M_2, M_3\})$

must sum over states at p and v

$s_1 \quad s_2 \quad s_3$

$$P(A, A, G) = P(1, 1, 2) =$$

$$\sum_{j=1}^4 \sum_{i=1}^4 \vec{p}(i) M_1(i,j) M_2(j,A) M_3(j,A) M_4(i,G)$$

\uparrow
v

\uparrow
p

degree 5 polynomial with 16 summands

IV) Specific Markov models on trees used in phylogenetics

The Jukes-Cantor model

1 parameter model

$$\vec{P}_1 = (.25 \ .25 \ .25 \ .25)$$

uniform root distribution

all bases equally likely

$$\bullet \text{ rate matrix } Q = \begin{pmatrix} -\alpha & \alpha/3 & \alpha/3 & \alpha/3 \\ \alpha/3 & -\alpha & \alpha/3 & \alpha/3 \\ \alpha/3 & \alpha/3 & -\alpha & \alpha/3 \\ \alpha/3 & \alpha/3 & \alpha/3 & -\alpha \end{pmatrix}$$

$$\alpha > 0$$

$\alpha \equiv$ total rate at which a specific base is changing to any of the other 3

$\alpha/3 = \text{off-diagonal} \Rightarrow$ constant rate for all conversions

• for a tree branch of length t

$$M(t) = e^{Qt} = \begin{pmatrix} 1-a & a & a & a \\ a & 1-a & a & a \\ a & a & 1-a & a \\ a & a & a & 1-a \end{pmatrix}$$

$$\text{where } a = a(t) = \frac{3}{4} \left(1 - e^{-\frac{4}{3}\alpha t} \right)$$

The matrix exponential is actually easy to compute since

the matrix of eigenvectors is a Hadamard matrix and the eigenvalues

of A are $0, -\frac{4}{3}\alpha, -\frac{4}{3}\alpha, -\frac{4}{3}\alpha$

Finally, note that the JC model has a STABLE BASE DISTRIBUTION

$$(.25 \ .25 \ .25 \ .25)M = (.25 \ .25 \ .25 \ .25) \Rightarrow$$

Should see uniform distribution of states in all sequences.

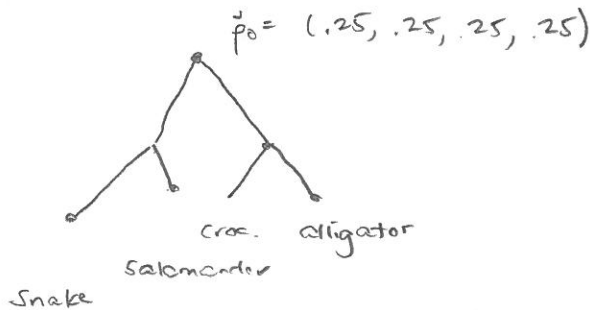
A stable base distribution \vec{p} for M is a λ eigenvector of M with eigenvalue $\lambda=1$.

$$\vec{p}M = \lambda \vec{p}$$

Detour: Review eigenvalues and eigenvectors possibly.

In particular, when the root distribution $\vec{p}_0 = \vec{p}$ = stable base distribution, then at all vertices of tree, the state distribution is $\vec{p}_0 = \vec{p}$. This is called a STATIONARY MODEL.

Eg. Jukes-Cantor JC



$$\vec{p}_{alligator} = (\quad)$$

$$\vec{p}_{salomander} = (\quad)$$

etc.

The KIMURA MODELS

$$\vec{p}_0 = (.25, .25, .25, .25)$$

K2P \equiv Kimura 2-parameter model

$$Q = \begin{pmatrix} * & \beta & \gamma & \gamma \\ \beta & * & \gamma & \gamma \\ \gamma & \gamma & * & \beta \\ \gamma & \gamma & \beta & * \end{pmatrix}$$

A, G, C, T order

$\beta \equiv$ transition rate

$2\gamma \equiv$ transversion rate

\uparrow

2 types. Eg $A \rightarrow C$ both equally likely
 $A \rightarrow T$

$$M_{K2P} = e^{Qt} = \begin{pmatrix} * & b & c & c \\ b & * & c & c \\ c & c & * & b \\ c & c & b & * \end{pmatrix}$$

$$\text{where } b = \frac{1}{4} (1 - 2e^{-2(\beta+\gamma)t} + e^{-4\gamma t})$$

$$c = \frac{1}{4} (1 - e^{-4\gamma t})$$

$*$ $= 1 - b - 2c$ Since row sums $= 1$.

Kimura 3 parameter model:

K3ST

14.

$$\vec{p}_0 = (.25 \quad .25 \quad .25 \quad .25) \quad \text{uniform}$$

$$Q = \begin{pmatrix} * & \beta & \gamma & \delta \\ \beta & * & \delta & \gamma \\ \gamma & \delta & * & \beta \\ \delta & \gamma & \beta & * \end{pmatrix} \quad \text{rate matrix} \quad \text{and}$$

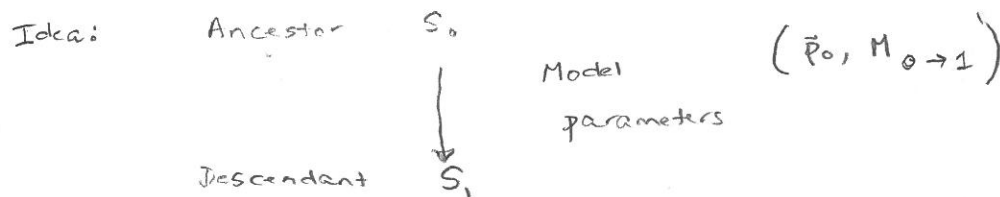
associated Markov matrices

$$M(t) = e^{Qt} = \begin{pmatrix} * & b & c & d \\ b & * & d & c \\ c & d & * & b \\ d & c & b & * \end{pmatrix}$$

Both Kimura models are stationary, i.e. the distribution of As, Gs, Cs, Ts should be uniform at all nodes of tree including leaves.

TIME REVERSIBLE MODELS including the

GENERAL TIME REVERSIBLE model. (GTR)



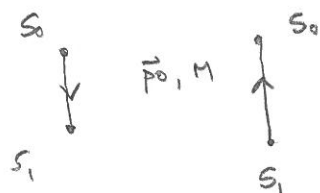
If the direction of the edge is reversed, then a TIME REVERSIBLE model uses the exact same parameters.



Clearly, a model can only be time reversible if it has a

Stable base distribution.

Mathematically, a model is time reversible if the joint distribution, i.e. the expected pattern frequency array does not change if you change the direction of an edge



If $P = P(i, j) = P(S_0 = i, S_1 = j)$, then

$$P = \underset{S_0 \rightarrow S_1}{\text{diag}(p_0) M} = \left[\underset{S_0 \rightarrow S_1}{\text{diag}(p_0) M} \right]^T = P^T$$

Defn: Process is time reversible, if

$$\text{diag}(p_0) M = M^T \text{diag}(p_0)$$

or in the continuous time formulation

$$\text{diag}(p_0) Q = Q^T \text{diag}(p_0)$$

Given a tree T^P , the GENERAL TIME REVERSIBLE model has numerical parameters

1) root distribution $p_p = (p_A, p_G, p_C, p_T)$

2) six rate parameters $\alpha, \beta, \gamma, \xi, \epsilon, \eta$

3) branch lengths t_e for each edge in T^P .

With these parameters, the common time reversible rate matrix is

$$Q = \begin{pmatrix} * & P_{A \rightarrow C} & P_{C \rightarrow G} & P_{T \rightarrow G} \\ P_{A \rightarrow C} & * & P_{C \rightarrow G} & P_{T \rightarrow G} \\ P_{A \rightarrow G} & P_{G \rightarrow C} & * & P_{T \rightarrow M} \\ P_{A \rightarrow G} & P_{G \rightarrow C} & P_{C \rightarrow M} & * \end{pmatrix}$$

* makes row sum = 0

and $M_e = e^{Q t_e}$ for each edge.

- GTR:
- 1) has \vec{p}_0 eigenvector \equiv stable base distribution
 - 2) most commonly used model in phylogenetic analyses

Other models

F81 \equiv arbitrary \vec{p}_0

$$\alpha = \beta = \gamma = \epsilon = \delta = \eta = 1$$

"JC without uniform root distribution"

HKY \equiv arbitrary \vec{p}_0

$$\beta = \gamma = \delta = \epsilon = 1 \quad \alpha = \eta = K$$

"K2ST' without uniform root distribution"