# Model-based Distances

Recall that 1 weakness of the Hamming distance = proportion of sites that differ

is that it fails to account for back substitutions or multiple substitutions on an edge ———•  ———•  and unless sequences are

$$A \to G \to A \qquad A \to C \to G$$

$\underline{\text{very}}$ closely related, the Hamming distance tends to under estimate the amount of evolutionary distance between 2 taxa.

To address this, we use our models JC, K2P, K3ST, GTR, GM to "correct" the Hamming distance and to account for unseen changes.

Example: Jukes - Cantor model and Jukes- Cantor distance

Parameters: single edge thought of as path from taxon $a$ to $b$ in tree

$$\vec{p_0} = (.25 \ .25 \ .25 \ .25)$$

$$Q = \begin{pmatrix} -\alpha & \alpha/3 & \alpha/3 & \alpha/3 \\ & -\alpha & & \\ & & -\alpha & \\ & & & -\alpha \end{pmatrix}, \text{ branch length } t, \text{ and Markov matrix } M(t) =$$

$$M = \begin{pmatrix} 1-a & a/3 & a/3 & a/3 \\ & 1-a & & \\ & & 1-a & \\ & & & 1-a \end{pmatrix} \qquad \text{with } a(t) = \frac{3}{4}\left(1 - \frac{4}{3}\alpha t\right)$$

Assuming the JC model, the expected pattern frequency array is

$$P = \text{diag}([.25 \ .25 \ .25 \ .25]) M(t) = \begin{pmatrix} \frac{1}{4}(1-a) & a/12 & a/12 & a/12 \\ a/12 & \frac{1}{4}(1-a) & a/12 & a/12 \\ a/12 & & \ddots & \\ a/12 & & & \end{pmatrix} \quad \begin{array}{l} \text{with } a \text{ as} \\ \text{above} \end{array}$$

[$\underline{\underline{No}}$ data here. yet.]

If we interpret $\alpha$ = rate in units $\frac{\#\text{substitutions}}{\text{time}}$, then

$$\alpha t = \left( \frac{\#\text{ substitutions per site}}{\text{time } t} \right) (\text{time } t)$$

↙ (model at single site)

$$\boxed{\alpha t = \#\text{ of substitutions over time } t}$$

including those hidden ones

This is (will be) the Jukes Cantor distance. once we solve for it.

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

$$\Rightarrow \boxed{\alpha t = -\frac{3}{4} \ln\left(1 - \frac{4}{3} a\right)} \quad = \#\text{ of subst. per site over elapsed time } t$$

Since we don't have $a$ in hand, we must estimate it from data,

i.e from the empirical pattern freq. array

In theory, $S_0$ and $S_1$ disagree with probability $12\left(\frac{a}{12}\right) = a$

Sum of off-diagonal entries

Thus, we estimate $a$ with $\hat{a} = \frac{\#\text{ of sites with non-constant pattern}}{\#\text{ of sites } n}$

$$= \text{Hamming distance !}$$

Defn: The Jukes-Cantor distance $d_{JC}(S_0, S_1)$ between aligned sequences is

$$d_{JC} = -\frac{3}{4} \ln\left(1 - \frac{4}{3}\hat{a}\right) \qquad \hat{a} = d_{Hamming}(S_0, S_1)$$

Ex.

$S_0$:  A  A  A  $\boxed{C}$  G      G  C  $\boxed{A}$  $\boxed{T}$  G

$S_1$:  A  A  A  $\boxed{T}$  G      G  C  $\boxed{T}$  $\boxed{A}$  G

$$d_{Hamming}(S_0, S_1) = \frac{3}{10} = .3$$

$$d_{JC}(S_0, S_1) = -\frac{3}{4} \ln\left(1 - \frac{4}{3}(.3)\right) \approx .38$$

Several Comments:

- $d_{JC} = .38 > .3 = d_{Hamming}$   to account for hidden mutations

- $d_{JC}(S_0, S_1) = -\frac{3}{4}\ln\left(1 - \frac{4}{3}\hat{a}\right)$   $\Rightarrow$   $0 \le \hat{a} < \frac{3}{4}$ for the log to make sense

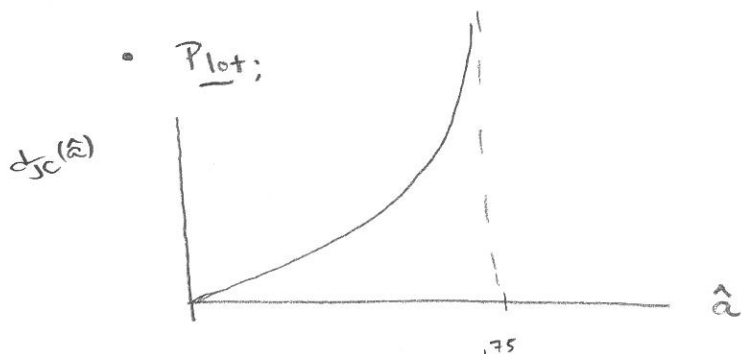  This makes sense: if 2 sequences are generated at random using the root distribution ($\frac{1}{4}$ $\frac{1}{4}$ $\frac{1}{4}$ $\frac{1}{4}$) then $S_0$ and $S_1$ will agree roughly $\frac{1}{4}$ of the time and disagree roughly $\hat{a} = \frac{3}{4}$ of the sites   $d(S_0, S_1) >> 0$

$S_0$:  _ _ _

$S_1$:  _ _ _

- Plot:



$d_{JC}(\hat{a})$

.75

$\hat{a}$

Thus as $\hat{a}$ gets close to .75, $d_{JC}(S_0, S_1) \rightarrow \infty$, huge distances

"Saturated sequences" informally, you can not differentiate them from 2 randomly selected sequences

- Since JC is time-reversible, on a tree or more generally we can consider paths in trees.

$t_1 \bigwedge t_2 \equiv \frac{t_1 + t_2}{S_0 \quad S_1}$

$S_0 \quad S_1$

Similar methods can be used to derive distance formulas for K2P, K3ST, GTR.

$$d_{K2P}(S_1, S_2) = -\frac{1}{2} \ln(1 - 2\hat{b} - \hat{e}) - \frac{1}{4} \ln(1 - 2\hat{e})$$

$\hat{b}$ = proportion of observed transitions

$\hat{e}$ = proportion of observed transv.

$$d_{k3}(S_1, S_2) = -\frac{1}{4}\left( \ln(1 - 2\hat{b} - 2\hat{c}) + \ln(1 - 2\hat{b} - \hat{d}) + \ln(1 - 2\hat{c} - \hat{d}) \right)$$

$\hat{b}$

$\hat{c}$

$\hat{d}$

are best estimates for $M_{K3P} = \begin{pmatrix} * & b & c & d \end{pmatrix}$

$\underset{=}{OR}$

$P = diag([\text{ }^{1}/_{4} \text{ }^{1}/_{4} \text{ }^{1}/_{4} \text{ }^{1}/_{4}]) M_{K3P}$

If $\hat{c} = \hat{d}$, then this simplifies to K2P.

There is also a more general GTR distance $\rightarrow$ see book. Requires

1) normalization          2) Knowledge of Trace of a matrix.

Why normalize?

Eg. JC model

Probability:          ( Model )

$$d_{JC}(S_1, S_2) = -\frac{3}{4} \ln\left(1 - \frac{4}{3}a\right)$$

$\parallel$

$\alpha t$

$\underbrace{\phantom{xxxx}}$

Intertwined

Estimate from data.

$$d_{JC}(S_0, S_1) = -\frac{3}{4} \ln\left(1 - \frac{4}{3}\hat{a}\right)$$

$$\alpha t = 2\alpha\left(\frac{t}{2}\right) = 4\alpha\left(\frac{t}{4}\right) = c\alpha\left(\frac{t}{c}\right)$$

we can compute the product $\alpha t$

but neither $\alpha$ or $t$ separately

However, $\alpha$ represents a total mutation rate: in JC

$$\frac{1}{4}\alpha + \frac{1}{4}\alpha + \frac{1}{4}\alpha + \frac{1}{4}\alpha = \alpha$$

rate leaving
State A

G    C    T

rate leaving state

$$Q = \begin{pmatrix} -1 & 1/3 & 1/3 & 1/3 \end{pmatrix} \text{ etc.}$$

Normalize so that $\alpha = 1$

$$\alpha t = 1(t) = \frac{\#\ of\ subst.}{t} \cdot t$$

$$= t \quad \#\ of\ substitutions$$

i.e.    $t$ measures the expected number of substitutions over the elapsed time.

MATLAB eg's