

Jukes-Cantor model

Simplest case: characters are nucleotides, each nucleotide is equally likely to mutate to any other.

Let's say our initial probability distribution over nucleotides at a site is:

$$\vec{p}_0 = \begin{pmatrix} p_A \\ p_C \\ p_G \\ p_T \end{pmatrix}$$

Let w_{ij} be the probability that a single ~~mutation~~ substitution changes nucleotide from i to j :

$$w_{ij} = \begin{cases} 1/3 & \text{if } i \neq j \\ 0 & \text{if } i = j \end{cases}$$

We can write a matrix

$$\underline{W} = [w_{ij}] = \begin{pmatrix} 0 & 1/3 & 1/3 & 1/3 \\ 1/3 & 0 & 1/3 & 1/3 \\ 1/3 & 1/3 & 0 & 1/3 \\ 1/3 & 1/3 & 1/3 & 0 \end{pmatrix}$$

After exactly one ~~mutation~~ ^{substitution}, the

new probability distribution over nucleotides is

$$\underline{W} \vec{p}_0$$

So what about after some amount of time? We don't know exactly how many ~~mutations~~ ^{substitutions} have occurred, but assuming a molecular clock, the average will be ut where u is ~~rate~~ substitution rate and t is time.

If ~~mut~~ substitutions are independent, then $\Pr(m|ut) = e^{-ut} \frac{(ut)^m}{m!}$

So:

$$\vec{p}(t) = \Pr(m=0|ut) \cdot \vec{p}_0 + \Pr(m=1|ut) \underline{W} \vec{p}_0 + \Pr(m=2|ut) \cdot \underline{W}^2 \vec{p}_0 + \dots$$

$$= \sum_{m=0}^{\infty} \Pr(m|ut) \underline{W}^m \vec{p}_0$$

$$= \sum_{m=0}^{\infty} e^{-ut} \frac{(ut)^m}{m!} \underline{W}^m \vec{p}_0$$

$$= e^{-ut} \left(\sum_{m=0}^{\infty} \frac{(ut \underline{W})^m}{m!} \right) \vec{p}_0$$

$$= e^{-ut} e^{ut \underline{W}} \vec{p}_0$$

$$= e^{-ut(\underline{W} - \underline{I})} \vec{p}_0 = e^{-ut \underline{P}} \vec{p}_0 \quad \text{where } \underline{P} = \underline{W} - \underline{I} \text{ is transition matrix.}$$

or substitution matrix.

So for Jukes-Cantor model:

$$P_{ij} = \begin{cases} 1/3 & \text{if } i \neq j \\ -1 & \text{if } i = j \end{cases}$$

and

~~$\vec{p}(t) = e^{+mt} \vec{p}_0$~~

$$\vec{p}(t) = e^{+mt} \vec{p}_0$$

For this very simple \mathbb{P} , it turns out that there is an analytic expression for the matrix exponential:

$$\left[e^{mt \mathbb{P}} \right]_{ij} = \begin{cases} \frac{1}{4} + \frac{3}{4} e^{-mt} & \text{if } i = j \\ \frac{1}{4} - \frac{1}{4} e^{-mt} & \text{if } i \neq j \end{cases}$$

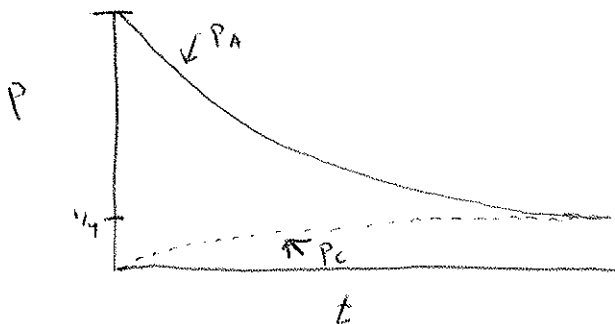
So let's say that

$$\vec{p}_0 = \begin{pmatrix} p_A = 1 \\ p_C = 0 \\ p_G = 0 \\ p_T = 0 \end{pmatrix}$$

Then at time t ,

$$p_A(t) = \frac{1}{4} + \frac{3}{4} e^{-mt}$$

$$p_C(t) = p_G(t) = p_T(t) = \frac{1}{4} - \frac{1}{4} e^{-mt}$$



$e^a e^b = e^{a+b}$

Using Jukes-Cantor for maximum likelihood estimate:

$$AA \xrightarrow{nt} AC$$

$$\Pr(AA|AC, nt) = \Pr(A|A, nt) \cdot \Pr(A|C, nt)$$

$$= \left(\frac{1}{4} + \frac{3}{4}e^{-nt}\right) \left(\frac{1}{4} - \frac{1}{4}e^{-nt}\right)$$

$$= \frac{1}{16} (1 + 3e^{-nt})(1 - e^{-nt})$$

~~$$= \frac{1}{16} (1 - e^{-nt} + 3e^{-nt} - 3e^{-2nt})$$~~

~~$$= \frac{1}{16} (1 + 2e^{-nt} - 3e^{-2nt})$$~~

~~$$= \frac{1}{16} (1 + 2e^{-nt} - 3e^{-2nt})$$~~

Not worth the algebra, doesn't clarify the point

nt	$\Pr(AA AC, nt)$
0	0
0.1	0.028
0.5	0.077
0.824	0.083
1.0	0.082
2.0	0.07
4.0	0.063

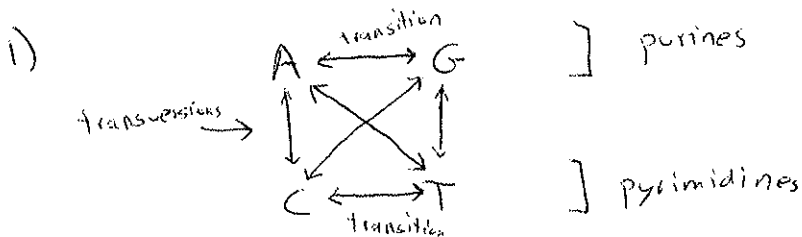
For not much data, likelihood is not sharply peaked, would be more peaked with more data.

Even for best value of nt , likelihood $\ll 1$

More general substitution models

In practice, we typically use models that are more complex than Jukes-Cantor and do the calculations using digital computing machines rather than pen and paper.

One example of such a model: HKY85. Captures 2 things not in Jukes-Cantor:



For most polymerases, transitions are more common than transversions. Why?

2) Empirically, most genes have unequal frequencies of the four nucleotides:

$$\phi_A \neq \phi_C \neq \phi_G \neq \phi_T$$

HKY85 substitution matrix

$$\begin{array}{c}
 \begin{array}{c} A \\ C \\ G \\ T \end{array} \begin{array}{c} \left[\begin{array}{cccc}
 & A & C & G & T \\
 A & * & \phi_C & \kappa\phi_G & \phi_T \\
 C & \phi_A & * & \phi_G & \kappa\phi_T \\
 G & \kappa\phi_A & \phi_C & * & \phi_T \\
 T & \phi_A & \kappa\phi_C & \phi_G & *
 \end{array} \right]
 \end{array}
 \end{array}$$

ϕ_x is equilibrium frequency of nucleotide x

κ is transition-transversion ratio. Is κ typically $>$ or $<$ 1?