

All we observe is differences between sequences:

ATG GAA ...

ATG GAA ...

We want to infer evolutionary distance t . We do this in terms of p the probability of a difference given z :

Jukes-Cantor: $p = \frac{3}{4} e^{-\frac{4}{3}ut} + \frac{1}{4}$

General: $\vec{p} = \vec{p}_0 e^{-utW}$

Equilibrium frequencies:

$$\vec{\pi} = \vec{\pi} W$$

Likelihood:

$\Pr(\text{data}|\text{model})$

Simple example: $\text{sequence 1} \quad \text{sequence 2}$
 $AC \xrightarrow{ut} AA$

In this case, the model is just the value of ut plus any free parameters is substitution model. For instance, for HKY, the substitution model would have 4 parameters: β , π_A , π_C , π_G .

Here we will just use Jukes-Cantor, so the model is just ut .

$$\begin{aligned} \Pr(AC|AA, ut) &= \Pr(A|A, ut) \cdot \Pr(C|A, ut) \leftarrow \text{what does this line assume?} \\ &= p \cdot \left(\frac{1-p}{3} \right) \\ &= \left(\frac{3}{4} e^{-\frac{4}{3}ut} + \frac{1}{4} \right) \left(\frac{1 - [\frac{3}{4} e^{-\frac{4}{3}ut} + \frac{1}{4}]}{3} \right) \end{aligned}$$

ut	$\Pr(AC AA, ut)$
0	0
0.1	0.028
0.5	0.077
0.824	0.083
1.0	0.082
2.0	0.070
4.0	0.063

maximum likelihood estimate \rightarrow

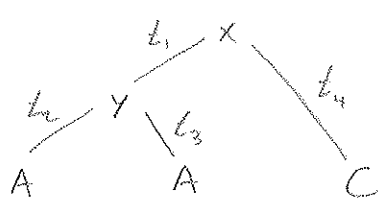
For not much data (2 nbs), likelihood is not sharply peaked. Would be more peaked with more data.

Even for best model, likelihood $\ll 1$. Why?

More than 2 sequences:

(2)

Alignment of 1 nt sequences: A
A
C



Here is one tree topology T.
Model is new

{substitution model, T, l_1, l_2, l_3, l_4 }

Given this topology T:

$$\Pr(\text{data} | \text{model}) = \sum_x \sum_y \Pr(A, A, C, x, y | l_1, l_2, l_3, l_4)$$

$$= \sum_x \sum_y \Pr(x) \cdot \Pr(y | x, l_1) \cdot \Pr(A | y, l_2) \cdot \Pr(A | y, l_3) \cdot \Pr(C | x, l_4)$$

(this sum will get very large for large trees, Felsenstein pruning algorithm, or dynamic programming)

$$= \sum_x \Pr(x) \cdot \Pr(C | x, l_4) \cdot \sum_y \Pr(y | x, l_1) \cdot \Pr(A | y, l_2) \cdot \Pr(A | y, l_3)$$

For Jukes-Cantor, what is $\Pr(x)$?

$$\Pr(x) = \pi_x = 1/4$$

For Jukes-Cantor, what is $\Pr(C | x, l_4)$?

$$\Pr(C | x, l_4) = \begin{cases} p = \frac{3}{4} e^{-\frac{4}{3} l_4} + \frac{1}{4} & \text{if } C = x \\ 1 - p = \frac{3}{4} - \frac{3}{4} e^{-\frac{4}{3} l_4} & \text{if } C \neq x \end{cases}$$

To find maximum likelihood, must optimize parameters (branch lengths, model parameters) for each topology
plus search over tree topologies.

Model comparison:

Should we use a more complex model? For instance Kimura 2-parameter vs. Jukes-Cantor?

Likelihood ratio test (nested models)

Akaike Information Criterion (AIC): $AIC = -2 \cdot \ln(L) + 2 \cdot \text{parameters}$

Bayes Theorem:

$$\begin{aligned} \Pr(\text{model} | \text{data}) &= \frac{\Pr(\text{data} | \text{model}) \cdot \Pr(\text{model})}{\Pr(\text{data})} \\ &= \frac{\Pr(\text{data} | \text{model}) \cdot \Pr(\text{model})}{\sum_{\text{model}'} \Pr(\text{data} | \text{model}') \cdot \Pr(\text{model}')} \end{aligned}$$

Advantage: Tells us what we really want, $\Pr(\text{model} | \text{data})$. No overfitting problem, given appropriate priors.

Disadvantages: 1) Hard to compute (less of a problem now)
2) what is prior, $\Pr(\text{model})$?

How to compute?

Markov Chain Monte Carlo (MCMC)

- 1) start w/ model m_1 ,
- 2) pick a new model m_2 . If using Metropolis's method, choose m_2 such that proposal rate $m_1 \rightarrow m_2 = m_2 \rightarrow m_1$. Typically local steps (m_2 similar to m_1).
- 3) Compute $R = \frac{\Pr(m_2 | \text{data})}{\Pr(m_1 | \text{data})} = \frac{\Pr(\text{data} | m_2) \cdot \Pr(m_2)}{\Pr(\text{data} | m_1) \cdot \Pr(m_1)}$
- 4) If $R \geq 1$, move to model m_2 . Otherwise move to m_2 w/ probability R .
- 5) Repeat step (2)

If MCMC is iterated repeatedly, the fraction f_i of samples when the chain is at model m_i is $f_i = \Pr(m_i | \text{data})$.

To see this, consider models m_i and m_j with $\Pr(m_i | \text{data}) \geq \Pr(m_j | \text{data})$

At equilibrium,

$$f_i \cdot \Pr(m_i \rightarrow m_j) = f_j \cdot \Pr(m_j \rightarrow m_i)$$

$$\text{Now } \Pr(m_i \rightarrow m_j) = R = \frac{\Pr(m_j | \text{data})}{\Pr(m_i | \text{data})}$$

$$\Pr(m_j \rightarrow m_i) = 1$$

$$\text{So: } f_i \cdot \frac{\Pr(m_j | \text{data})}{\Pr(m_i | \text{data})} = f_j \Rightarrow \frac{f_i}{f_j} = \frac{\Pr(m_i | \text{data})}{\Pr(m_j | \text{data})}$$

In MCMC, how do we know when we have run chain "long enough"?

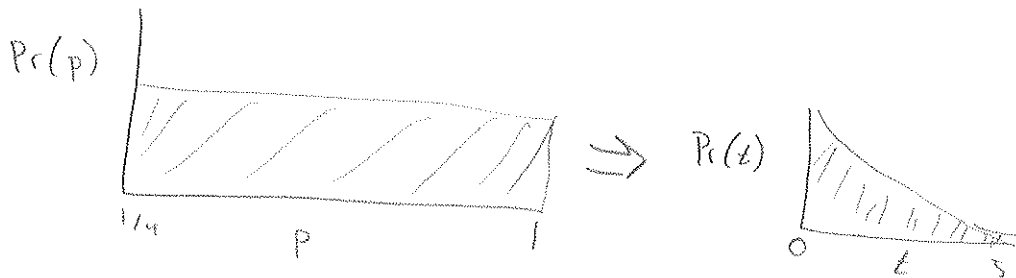
About the prior:

④

Jukes-Cantor: $p = \frac{3}{4} e^{-\frac{4}{3} \mu t} + \frac{1}{4}$

Choosing a "Flat" prior over parameters can be dangerous...

Flat prior on p



Flat prior on t

