

Phylogenetic inference and likelihood 2

Emily Jane McTavish

Life and Environmental Sciences
University of California, Merced

`ejmctavish@ucmerced.edu`, `twitter:snacktavish`

(With thanks to Mark Holder and Paul Lewis for slides)

Combining probabilities

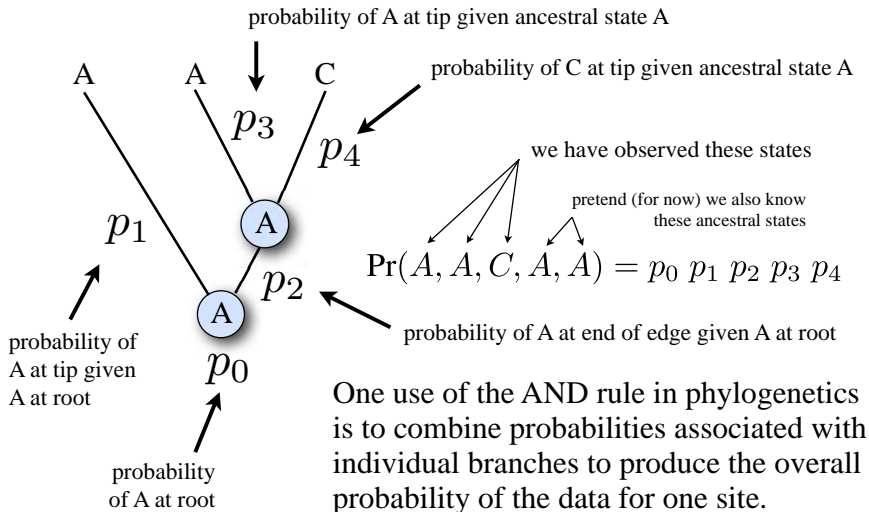
- *Multiply* probabilities if the component events must happen **simultaneously** (i.e. where you would naturally use the word AND when describing the problem)

Using 2 dice, what is the probability of



$$(1/6) \times (1/6) = 1/36$$

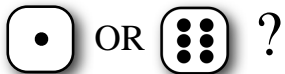
AND rule in phylogenetics



Combining probabilities

- *Add* probabilities if the component events are **mutually exclusive** (i.e. where you would naturally use the word OR in describing the problem)

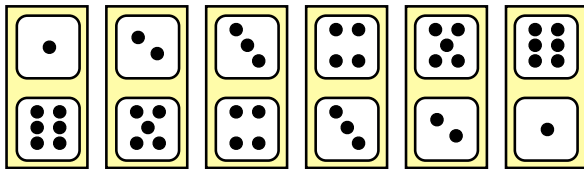
Using one die, what is the probability of



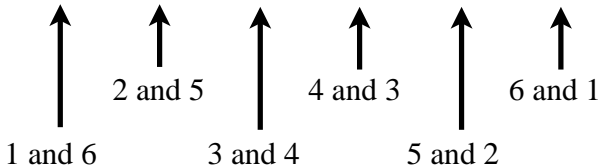
$$(1/6) + (1/6) = 1/3$$

Combining AND and OR

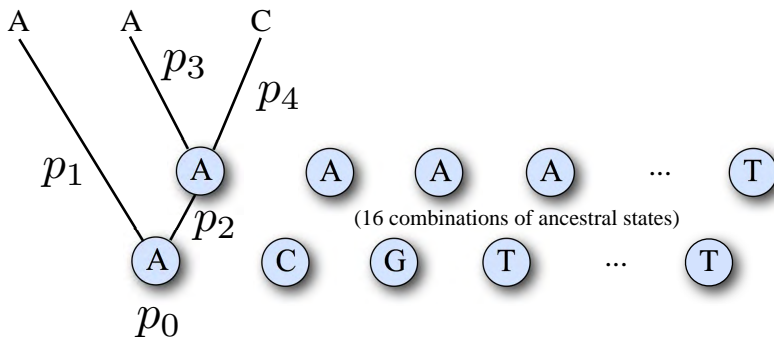
What is the probability that the sum of two dice is 7?



$$(1/36) + (1/36) + (1/36) + (1/36) + (1/36) + (1/36) = 1/6$$



Using both AND and OR in phylogenetics



AND rule used to compute probability of the observed data for *each combination* of ancestral states.

OR rule used to combine different combinations of ancestral states.

Independence

This is always true...

$$\underset{\text{joint probability}}{\Pr(A \text{ and } B)} = \Pr(A) \underset{\text{conditional probability}}{\Pr(B|A)}$$

If we can say this...

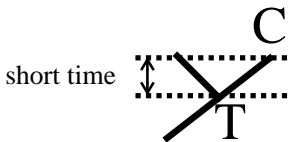
$$\Pr(B|A) = \Pr(B)$$

...then events A and B are **independent** and we can express the joint probability as the product of $\Pr(A)$ and $\Pr(B)$

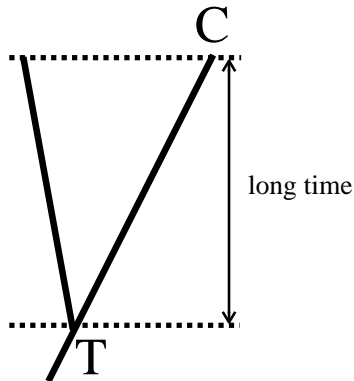
$$\Pr(A \text{ and } B) = \Pr(A) \Pr(B)$$

Non-independence in molecular evolution

The state present in the descendant is **not independent** of the state in the ancestor



less probable



more probable

Conditional Independence

Assume both A and B depend on C:

$$\Pr(A|C) \neq \Pr(A) \quad \Pr(B|C) \neq \Pr(B)$$

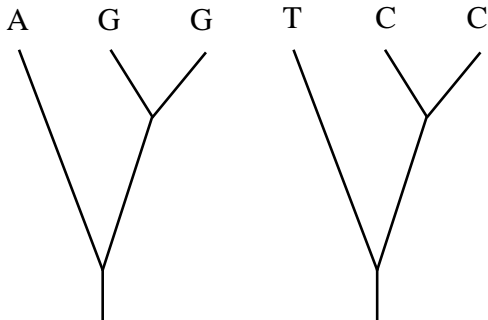
If we can say this...

$$\Pr(B|A,C) = \Pr(B|C)$$

...then events A and B are **conditionally independent** and we can express the joint (conditional) probability as the product of $\Pr(A|C)$ and $\Pr(B|C)$

$$\Pr(A \text{ and } B|C) = \Pr(A|C) \Pr(B|C)$$

Conditional independence in molecular evolution



The site data patterns AGG and TCC are assumed by most models to be conditionally independent.

The patterns both depend on the underlying tree (including edge lengths) and the substitution model.

$$\Pr(\text{AGG and TCC}|\text{tree, model}) = \Pr(\text{AGG}|\text{tree, model}) \Pr(\text{TCC}|\text{tree, model})$$

Likelihood calculated from a single sequence

$$\Pr(A) = \pi_A$$

$$\Pr(C) = \pi_C$$

$$\Pr(G) = \pi_G$$

$$\Pr(T) = \pi_T$$

First 32 nucleotides of the $\psi\eta$ -globin gene of gorilla:

GAAGTCCTTGAGAAATAAACTGCACACACTGG

$$\begin{aligned} L &= \pi_G \pi_A \pi_A \pi_G \pi_T \pi_C \pi_C \pi_T \pi_T \pi_G \pi_A \pi_G \pi_A \pi_A \pi_A \pi_T \pi_A \pi_A \pi_A \pi_C \pi_T \pi_G \pi_C \pi_A \pi_C \pi_A \pi_C \pi_A \pi_C \pi_T \pi_G \pi_G \\ &= \pi_A^{12} \pi_C^7 \pi_G^7 \pi_T^6 \end{aligned}$$

Note that we are assuming independence among sites here

$$\log L = 12 \log(\pi_A) + 7 \log(\pi_C) + 7 \log(\pi_G) + 6 \log(\pi_T)$$

We can already see by eye-balling this that a model allowing **unequal** base frequencies will **fit better** than a model that assumes **equal** base frequencies because there are about twice as many As as there are Cs, Gs and Ts.

Model ranking using LRT or AIC

Likelihood Ratio Tests (LRT) and the Akaike Information Criterion (AIC) provide two ways to evaluate whether an **unconstrained** model fits the data significantly better than a **constrained** version of the same model.

Find *maximum* $\log L$ under the *unconstrained* model:

$$\begin{aligned}\log L_{\text{unconstrained}} &= 12 \log(\pi_A) + 7 \log(\pi_C) + 7 \log(\pi_G) + 6 \log(\pi_T) \\ &= 12 \log(0.375) + 7 \log(0.219) + 7 \log(0.219) + 6 \log(0.187) \\ &= -43.1\end{aligned}$$

This model has 3 estimated parameters

Find *maximum* $\log L$ under the *constrained* model:

$$\begin{aligned}\log L_{\text{constrained}} &= 12 \log(\pi_A) + 7 \log(\pi_C) + 7 \log(\pi_G) + 6 \log(\pi_T) \\ &= 12 \log(0.25) + 7 \log(0.25) + 7 \log(0.25) + 6 \log(0.25) \\ &= -44.4\end{aligned}$$

This model has 0 estimated parameters

Likelihood Ratio Test (LRT)

Calculate the likelihood ratio test statistic:

$$\begin{aligned} R &= -2 [\log(L_{\text{constrained}}) - \log(L_{\text{unconstrained}})] \\ &= -2 [-44.4 - (-43.1)] \\ &= 2.6 \end{aligned}$$

(Note that the log-likelihoods used in the test statistic have been *maximized* under each model separately)

“unconstrained” does fit better than “constrained” ($-43.1 > -44.4$), but not significantly better ($P = 0.457$, chi-squared with 3 d.f.*)

*The number of degrees of freedom equals the difference between the two models in the number of estimated parameters. In this case, unconstrained has 3 parameters and constrained has 0, so $\text{d.f.} = 3 - 0 = 3$

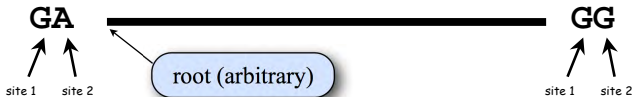
Comparing models in phylogenetics can be challenging, as topologies are not nested within one another.

We will discuss appropriate statistical approaches later in the course.

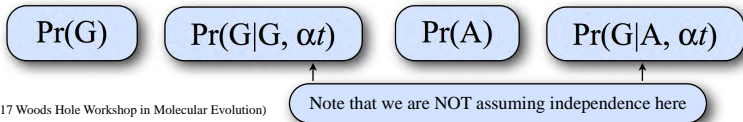
Likelihood of the simplest tree

sequence 1 sequence 2

To keep things simple, assume that the sequences are only 2 nucleotides long:



$$\begin{aligned}
 L &= L_1 L_2 \\
 &= \left[\begin{pmatrix} 1 \\ 4 \end{pmatrix} \quad \left(\frac{1}{4} + \frac{3}{4}e^{-4\alpha t} \right) \right] \left[\begin{pmatrix} 1 \\ 4 \end{pmatrix} \quad \left(\frac{1}{4} - \frac{1}{4}e^{-4\alpha t} \right) \right]
 \end{aligned}$$



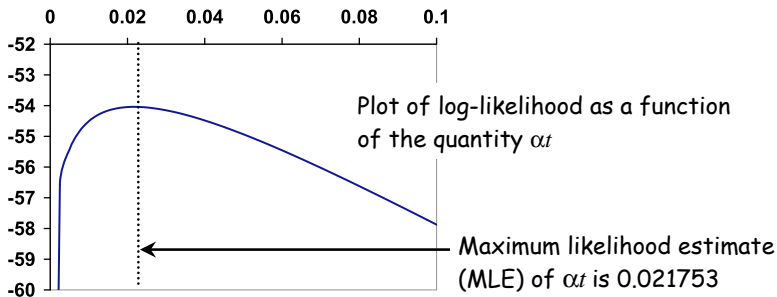
Maximum likelihood estimation

First 32 nucleotides of the $\psi\eta$ -globin gene of gorilla and orangutan:

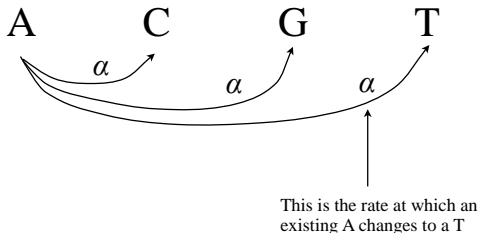
gorilla **GAAG**TCCTTGAGAAATAAACTGCACACACTGG

orangutan **GAAC**TCCTTGAGAAATAAACTGCACACACTGG

$$L = \left[\left(\frac{1}{4} \right) \left(\frac{1}{4} + \frac{3}{4} e^{-4\alpha t} \right) \right]^{30} \left[\left(\frac{1}{4} \right) \left(\frac{1}{4} - \frac{1}{4} e^{-4\alpha t} \right) \right]^2$$



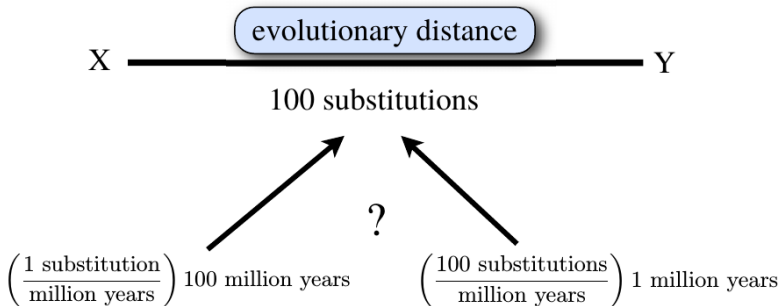
number of substitutions = rate \times time



Overall substitution rate is 3α , so the expected number of substitutions (v) is

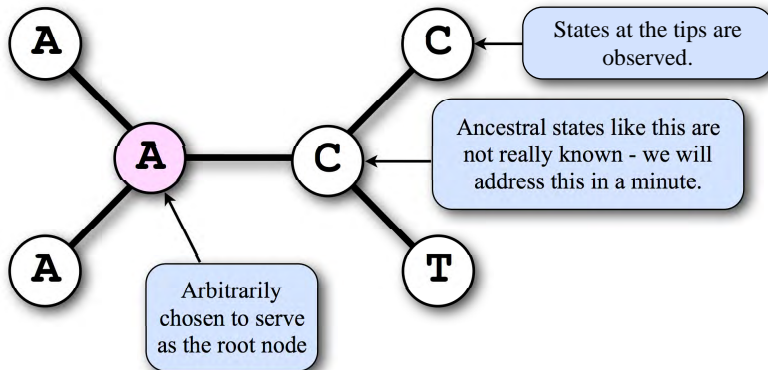
$$v = 3\alpha t$$

Rate and time are confounded

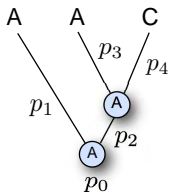


Likelihood of an unrooted tree

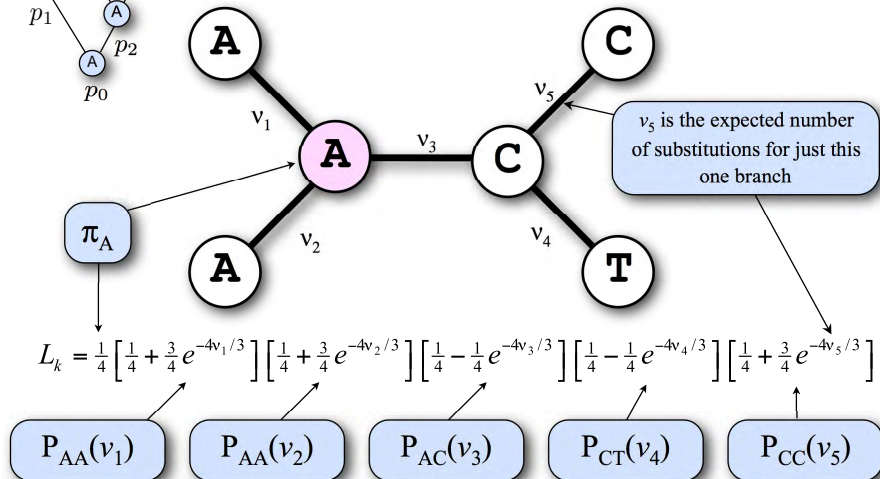
(data shown for only one site)



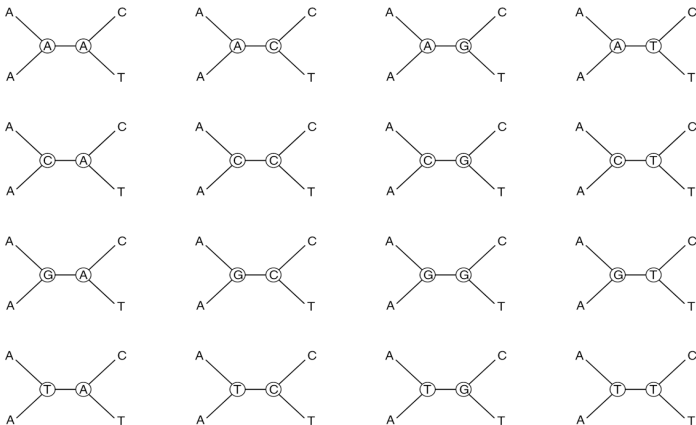
From slide 6



Likelihood for site k



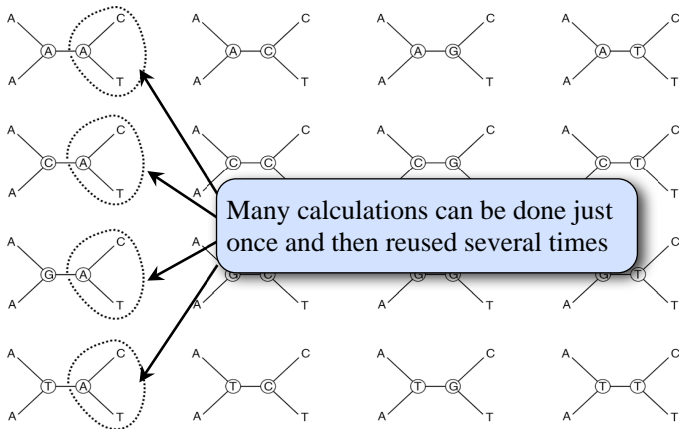
Brute force approach would be to calculate L_k for all 16 combinations of ancestral states and sum them



Note use of the OR probability rule

Pruning algorithm

(same result, less time)

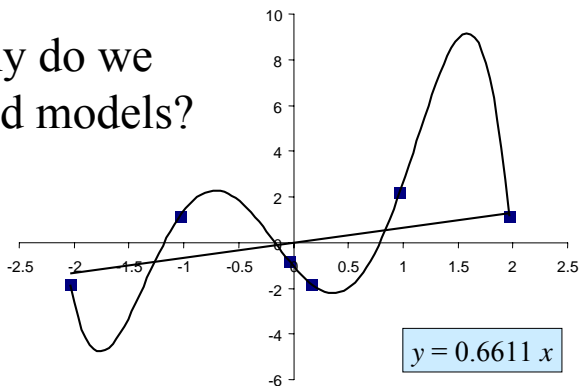


Felsenstein, J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach. *Journal of Molecular Evolution* 17:368-376

Substitution Models

$$y = -1.5972 x^5 + 23.167 x^4 - 126.18 x^3 + 319.17 x^2 - 369.22 x + 155.67$$

Why do we
need models?



Models

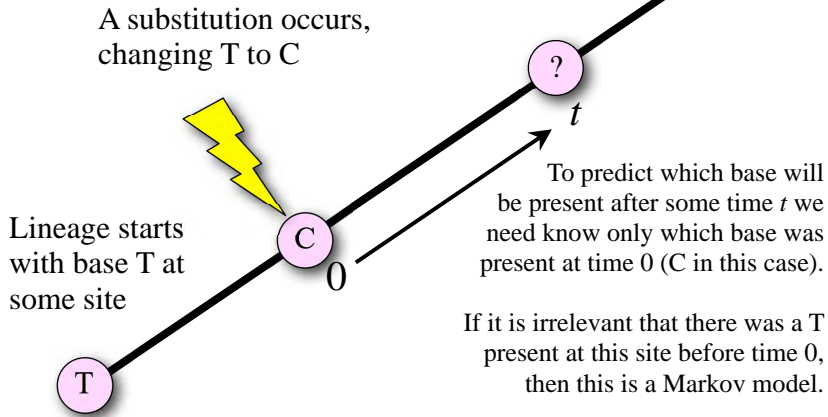
- Models help us intelligently **interpolate between our observations** for purposes of **making predictions**
- **Adding parameters** to a model generally increases its fit to the data
- **Underparameterized** models lead to poor fit to observed data points
- **Overparameterized** models lead to poor prediction of future observations
- Criteria for choosing models include likelihood ratio tests, AIC, BIC, Bayes Factors, etc.
 - all provide a way to choose a model that is neither underparameterized nor overparameterized

Jukes-Cantor (JC69) model

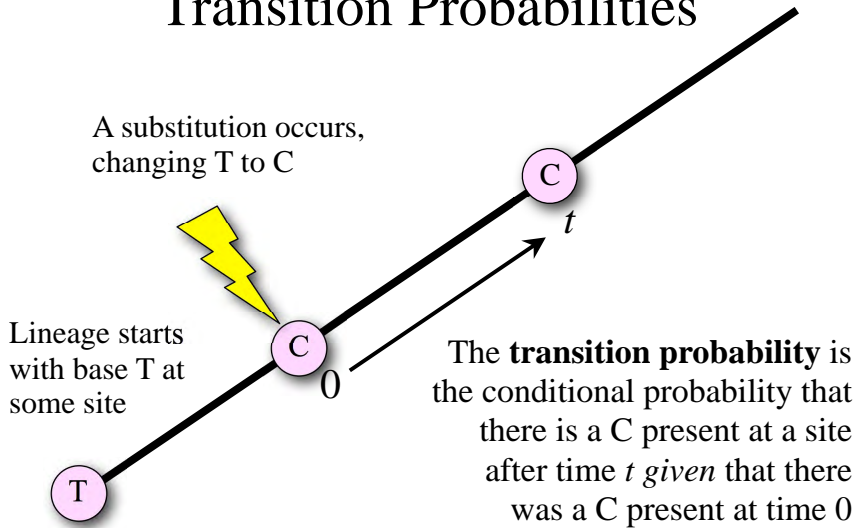
- The four bases (A, C, G, T) are expected to be **equally frequent** in sequences ($\pi_A = \pi_C = \pi_G = \pi_T = 0.25$)
- Assumes **same rate** for all types of substitution
($r_{A \leftrightarrow C} = r_{A \leftrightarrow G} = r_{A \leftrightarrow T} = r_{C \leftrightarrow G} = r_{C \leftrightarrow T} = r_{G \leftrightarrow T} = \alpha$)
- Usually described as a **1-parameter** model (the parameter being the edge length)
 - Remember, however, that each edge in a tree can have its own length, so there are really as many parameters in the model as there are edges in the tree!
- Assumes substitution is a **Markov** process...

Jukes, T. H., and C. R. Cantor. 1969. Evolution of protein molecules. Pages 21-132 in H. N. Munro (ed.), *Mammalian Protein Metabolism*. Academic Press, New York.

What is a Markov process?



Transition Probabilities



Jukes-Cantor transition probabilities

Here is the probability that a site starting in state T will end up in state G after time t when the individual substitution rates are all α :

$$P_{TG}(t) = \frac{1}{4} (1 - e^{-4\alpha t}) = \Pr(G|T, \alpha t)$$

The JC69 model has only one unknown quantity: αt

(The symbol e represents the base of the natural logarithms: its value is 2.718281828459045...)

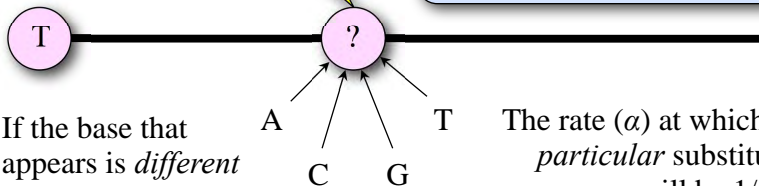
Where does a transition probability formula such as this come from?

"ACHNyons" vs. substitutions

ACHN =
"Anything
Can Happen
Now"

When an *achnyon* occurs, any base can appear in a sequence.

Note: *achnyon* is *my term* for this make-believe event. You will not see this term in the literature.



If the base that appears is *different* from the base that was already there, then a **substitution** event has occurred.

The rate (α) at which any *particular* substitution occurs will be 1/4 the achnyon rate (μ).
That is, $\alpha = \mu/4$
(or $\mu = 4\alpha$)

The Poisson distribution

Probability distribution on the number of events when:

1. events are assumed to be independent,
2. the *rate* of events some constant, μ , and
3. the process continues for some duration of time, t .

The expectation of the number of events is $\nu = \mu t$.

Note that ν can be any non-negative number, but the Poisson is a discrete distribution – it gives the probabilities of the number of events (and this number will always be a non-negative integer).

Poisson distribution can be used to explain statistical regularities of rare events

$$P(k \text{ events in interval}) = \frac{e^{-\nu} \nu^k}{k!}$$

- ▶ ν is the average number of events per interval (rate times time)
- ▶ e is the number 2.71828... (Euler's number) the base of the natural logarithms
- ▶ k takes values 0, 1, 2, ...
- ▶ $k! = k * (k - 1) * (k - 2) * \dots * 2 * 1$ is the factorial of k .

from wikipedia

$$P(\text{k events in interval}) = \frac{e^{-\nu} \nu^k}{k!}$$
$$P(0 \text{ events}) = \frac{e^{-\nu} \nu^0}{0!} = e^{-\nu} = e^{-\mu t}$$
$$P(\geq 1 \text{ events}) = 1 - e^{-\mu t}$$

Deriving a transition probability

Calculate the probability that a site currently T will change to G over time t when the rate of this particular substitution is α :

$$\Pr(\text{zero achnyons}) = e^{-\mu t} \quad (\text{Poisson probability of zero events})$$

$$\Pr(\text{at least 1 achnyon}) = 1 - e^{-\mu t}$$

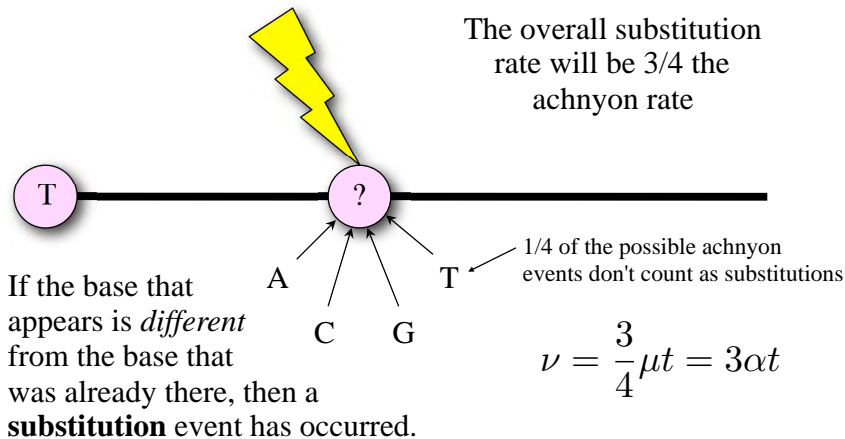
$$\Pr(\text{last achnyon results in base G}) = \frac{1}{4}$$

$$\Pr(\text{end in G} \mid \text{start in T}) = \frac{1}{4} (1 - e^{-\mu t})$$

Remember that the rate (α) of any particular substitution is one fourth the achnyon rate (μ):

$$P_{GT}(t) = \frac{1}{4} (1 - e^{-4\alpha t})$$

Expected number of substitutions



Transition Probabilities: Remarks

$$P_{TA}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TC}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TG}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TT}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

These should add to 1.0 because T *must* change to something!

$$1 - e^{-4\alpha t}$$

Doh! Something must be wrong here...

Transition Probabilities: Remarks

$$P_{TA}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TC}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TG}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TT}(t) = \frac{1}{4}(1 - e^{-4\alpha t}) + e^{-4\alpha t}$$

Forgot to account for the possibility of *no* acnyons over time t

Equilibrium frequencies

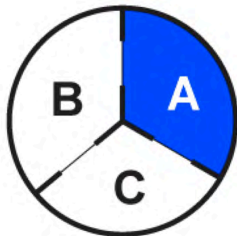
- The JC69 model assumes that the frequencies of the four bases (A, C, G, T) are equal
- The equilibrium relative frequency of each base is thus 0.25
- Why are they called *equilibrium* frequencies?

Equilibrium Frequencies

Imagine a bottle of perfume has been spilled in room A.

The doors to the other rooms are closed, so the perfume has, thus far, not been able to spread.

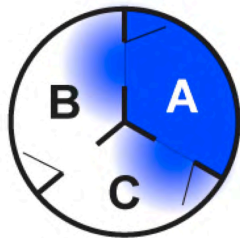
What would happen if we opened all the doors?



Equilibrium Frequencies

If the doors are suddenly opened, the perfume would begin diffusing from the area of highest concentration to lowest.

Molecules of perfume go both ways through open doors, but more pass one way than another, leading to a net flow from room A to rooms B and C.

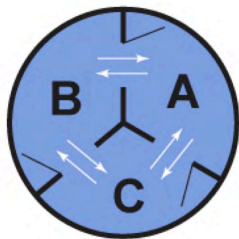


In the instant that the doors are opened, A is losing perfume molecules at *twice the rate* each of the other rooms is gaining molecules. As diffusion progresses, however, the rate of loss from A drops, approaching an equilibrium.

Equilibrium Frequencies

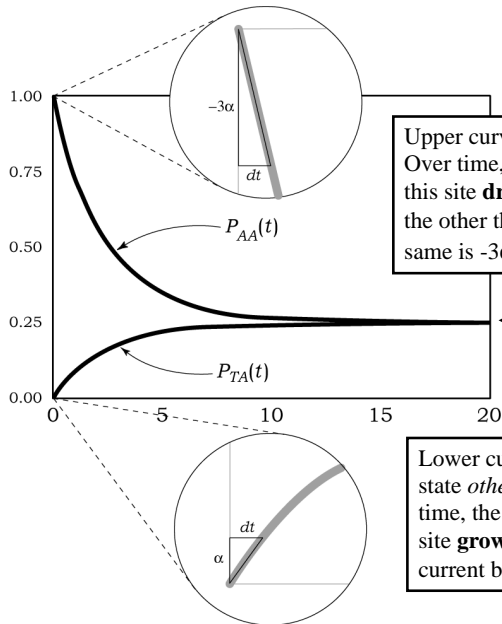
Eventually, all 3 rooms have essentially the same concentration of perfume.

Molecules still move through open doors, but now the rates are the same in all directions.



Back to sequence evolution: assume a sequence began with only A nucleotides (a poly-A sequence). Over time, substitution would begin converting some of these As to Cs, Gs, and Ts, just as the perfume diffused into adjacent rooms.

Pr(A|A) and Pr(A|T) as a function of time

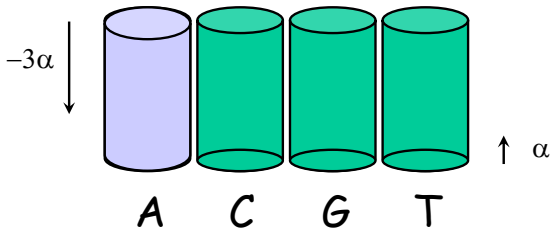


Upper curve assumes we started with A at time 0. Over time, the probability of still seeing an A at this site **drops** because rate of changing to one of the other three bases is 3α (so rate of staying the same is -3α).

The equilibrium relative frequency of A is 0.25

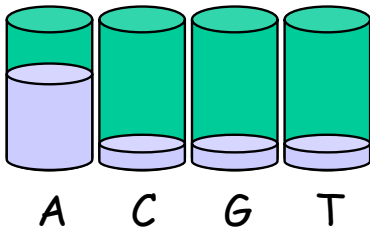
Lower curve assumes we started with some state *other* than A (T is used here). Over time, the probability of seeing an A at this site **grows** because the rate at which the current base will change into an A is α .

Water analogy (time 0)



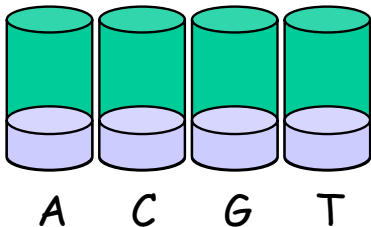
- Start with container A completely full and others empty
- Imagine that all containers are connected by tubes that allow same rate of flow between any two
- Initially, A will be losing water at 3 times the rate that C (or G or T) gains water

Water analogy (after some time)



A's level is not dropping as fast now because it is now also *receiving* water from C, G and T

Water analogy (after a very long time)



Eventually, all containers are one fourth full and there is zero *net* volume change – **stationarity** (equilibrium) has been achieved

(Thanks to Kent Holsinger for this analogy)

JC69 rate matrix

1 parameter:
 α

		To			
		A	C	G	T
From	A	-3α	α	α	α
	C	α	-3α	α	α
	G	α	α	-3α	α
	T	α	α	α	-3α

Jukes, T. H., and C. R. Cantor. 1969. Evolution of protein molecules. Pages 21-132 in H. N. Munro (ed.), *Mammalian Protein Metabolism*. Academic Press, New York.

K80 (or K2P) rate matrix

2 parameters:

α
 β

		To			
		A	C	G	T
From	A	$-\alpha - 2\beta$	β	α	β
	C	β	$-\alpha - 2\beta$	β	α
	G	α	β	$-\alpha - 2\beta$	β
	T	β	α	β	$-\alpha - 2\beta$

↑ transition rate ↑ transversion rate

Kimura, M. 1980. A simple method for estimating evolutionary rate of base substitutions through comparative studies of nucleotide sequences. *Journal of Molecular Evolution* 16:111-120.

K80 rate matrix

(looks different, but actually the same)

2 parameters:

κ

β

	A	C	G	T
A	$-\beta(\kappa + 2)$	β	$\kappa\beta$	β
C	β	$-\beta(\kappa + 2)$	β	$\kappa\beta$
G	$\kappa\beta$	β	$-\beta(\kappa + 2)$	β
T	β	$\kappa\beta$	β	$-\beta(\kappa + 2)$

All I've done is re-parameterize the rate matrix,
letting κ equal the *transition/transition rate ratio*

$$\longrightarrow \kappa = \frac{\alpha}{\beta}$$

Note: the K80 model is identical to the JC69 model if $\kappa = 1$ ($\alpha = \beta$)

Transition/transversion ratio (ratio) versus Transition/transversion *rate* ratio (kappa)



Cobbler analogy:

- 4 cobblers in a factory make loafers
- 8 cobblers in the factory make work boots
- all cobblers produce the same number of shoes per unit time, regardless of shoe type
- what is the loafer/boot *rate ratio* and how does that compare to the loafer/boot *ratio*?

The loafer/boot *rate ratio* is 1.0 because each cobbler cranks out shoes at the same rate.

The loafer/boot *ratio*, however, is 0.5 because there are twice as many cobblers making boots as there are cobblers making loafers.

There are 8 possible transversion-type substitutions and only 4 possible transition-type substitutions: the transition/transversion ratio is thus 0.5 when the transition/transversion rate ratio is 1.

F81 rate matrix

4 parameters:

μ

π_A

π_C

π_G

	A	C	G	T
A	$-\mu(1 - \pi_A)$	$\pi_C\mu$	$\pi_G\mu$	$\pi_T\mu$
C	$\pi_A\mu$	$-\mu(1 - \pi_C)$	$\pi_G\mu$	$\pi_T\mu$
G	$\pi_A\mu$	$\pi_C\mu$	$-\mu(1 - \pi_G)$	$\pi_T\mu$
T	$\pi_A\mu$	$\pi_C\mu$	$\pi_G\mu$	$-\mu(1 - \pi_T)$

Note: the F81 model is identical to the JC69 model if all base frequencies are equal

HKY85 rate matrix

5 parameters:

κ
 β
 π_A
 π_C
 π_G

	A	C	G	T
A	—	$\pi_C \beta$	$\pi_G \beta \kappa$	$\pi_T \beta$
C	$\pi_A \beta$	—	$\pi_G \beta$	$\pi_T \beta \kappa$
G	$\pi_A \beta \kappa$	$\pi_C \beta$	—	$\pi_T \beta$
T	$\pi_A \beta$	$\pi_C \beta \kappa$	$\pi_G \beta$	—

A dash means equal to negative sum of other elements on the same row

Note: the HKY85 model is identical to the F81 model if $\kappa = 1$. If, in addition, all base frequencies are equal, it is identical to JC69.

Hasegawa, M., H. Kishino, and T. Yano. 1985. Dating of the human-ape splitting by a molecular clock of mitochondrial DNA. *Journal of Molecular Evolution* 21:160-174.

F84 vs. HKY85

F84 model:

μ rate of process generating *all types of substitutions*

$k\mu$ rate of process generating *only transitions*

Becomes F81 model if $k = 0$

HKY85 model:

β rate of process generating *only transversions*

$\kappa\beta$ rate of process generating *only transitions*

Becomes F81 model if $\kappa = 1$

F84 first used in Felsenstein's PHYLIP package in 1984, first published by: Kishino, H., and M. Hasegawa. 1989. Evaluation of the maximum likelihood estimate of the evolutionary tree topologies from DNA sequence data, and the branching order in hominoidea. *Journal of Molecular Evolution* 29: 170-179.

GTR rate matrix

	A	C	G	T
A	—	$\pi_C a \mu$	$\pi_G b \mu$	$\pi_T c \mu$
C	$\pi_A a \mu$	—	$\pi_G d \mu$	$\pi_T e \mu$
G	$\pi_A b \mu$	$\pi_C d \mu$	—	$\pi_T f \mu$
T	$\pi_A c \mu$	$\pi_C e \mu$	$\pi_G f \mu$	—

9 parameters:

π_A
 π_C
 π_G
 a
 b
 c
 d
 e
 f

Identical to the F81 model if $a = b = c = d = e = f = 1$. If, in addition, all the base frequencies are equal, GTR is identical to JC69. If $a = c = d = f = \beta$ and $b = e = \kappa\beta$, GTR becomes the HKY85 model.

Lanave, C., G. Preparata, C. Saccone, and G. Serio. 1984.
A new method for calculating evolutionary substitution
rates. *Journal of Molecular Evolution* 20:86-93.

Rate Heterogeneity

Green Plant *rbcL*

First 88 amino acids (translation is for *Zea mays*)

M--S--P--Q--T--E--T--K--A--S--V--G--F--K--A--G--V--K--D--Y--K--L--T--Y--Y--T--P--E--Y--E--T--K--D--T--D--I--L--A--A--F--R--V--T--P--	
Chara (green alga; land plant lineage)	AAAGATTACAGATTAACTTACTATACTCTGAGTATAAACTAAAGATACTGACATTTTAGCTGCAITTCGTGTAACCTCCA
Chlorella (green alga)C.....C.T.....T.....CC.C.A.....C.....T.....C.T.....G.C.....A.G.....T
Volvox (green alga)TC.T.....A.....C.....A.....C.....GT.GTA.....C.....C.....A.....A.G.....
Conocephalum (liverwort)TC.....T.....T.....G.T.....G.....G.....T.....G.....A.....A.A.G.....T
Bazzania (moss)T.....C.....T.....G.....A.....G.G.C.....G.....A.....T.....G.....A.....A.G.....C
Anthoceros (hornwort)T.....CC.T.....C.....T.....CG.G.C.....G.....T.....G.....A.....G.C.T.AA.G.....T
Osmunda (fern)TC.....G.....C.....C.....C.....T.....G.G.C.....G.....T.....G.....A.....C.....AA.G.....C
Lycopodium (club "moss")GG.....C.....C.T.....C.....T.....G.C.....A.....C.....T.....C.G.....A.....AA.G.....T
Ginkgo (gymnosperm; Ginkgo biloba)G.....T.....A.....C.....C.....C.....T.....C.G.....A.....C.....A.....T.....T
Picea (gymnosperm; spruce)T.....T.....A.....A.....C.G.C.....C.....G.....T.....G.....A.....C.....A.....T.....T
Iris (flowering plant)C.....G.....T.....CG.C.....C.....C.....T.....C.G.....A.....C.....A.....T.....T
Asplenium (fern; spleenwort)TC.....C.G.....T.....C.....C.....C.....A.....C.....G.C.....C.....T.....C.G.....A.....T.....C.....GA.G.....C.....
Nicotiana (flowering plant; tobacco)G.....A.....G.....T.....T.....CC.....C.....G.....T.....A.G.....A.....C.....A.....T.....T

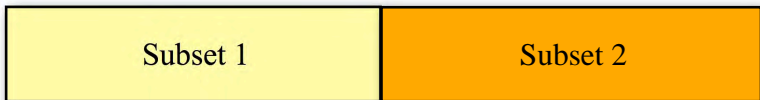
Q--L--G--V--P--P--E--E--A--G--A--A--V--A--A--E--S--S--T--G--T--W--T--T--V--W--T--D--G--L--T--S--L--D--R--Y--K--G--R--C--Y--H--I--E--	
CAACCTGSCGTTCCACCTGAAGAAGCAGGGGCTGCAGTAGCTGCAGAAATCTTCTACTGGTACATGGACTACTGTTTGGAC	GTGACGGGATTAAGTAGTTTGGACCGATACAAAGGAAGATGCTACGATATTGAA
.....A.....T.....A.....G.....T.....G.....A.....A.....A.....T.....A.....A.....T.....T.....TC.T.....T.....C.....C.....G.....T.....T.....A.....C.....T.....T.....TC.T.....T.....C.....C.....G.....
.....A.....T.....TGT.....T.....T.....T.....A.....A.....T.....A.....A.....T.....T.....TC.T.....T.....C.....C.....G.....T.....TC.T.....ACC.T.....T.....T.....TC.T.....T.....G.....C.....
.....G.....A.....G.A.....A.....A.....T.....T.....A.....A.....G.....C.....G.....C.....G.....C.....T.....GC.T.....A.....C.....C.....T.....TC.....T.....C.....C.....C.....T.....C.....C.....C.....C.....T.....TC.....C.....C.....
.....T.....A.....A.....C.....G.....G.A.....C.....T.....C.....C.....C.....C.....C.....G.....C.....T.....C.....C.....C.....T.....TC.....G.....T.....A.....C.....T.....C.....C.....C.....T.....T.....G.....C.....T.....C.....C.....G.....
.....A.....G.....G.....G.....G.....A.....C.....C.....C.....C.....C.....C.....C.....C.....T.....C.....C.....C.....T.....T.....G.....GC.....T.....C.....C.....G.....C.....T.....C.....C.....C.....T.....T.....G.....GC.....T.....C.....C.....G.....
.....C.....A.....TG.....G.....C.....G.....C.....C.....C.....A.....A.....G.....T.....C.....C.....C.....C.....T.....T.....G.....CC.....C.....C.....C.....G.....C.....A.....G.....C.....A.....C.....C.....G.....C.....A.....G.....G.....C.....C.....T.....T.....G.....CC.....C.....C.....G.....
.....A.....A.....G.....C.....G.....C.....C.....C.....C.....A.....A.....A.....C.....T.....C.....T.....C.....CC.T.....T.....T.....GC.....C.....G.....C.....G.....	

All four bases are observed at some sites...

...while at other sites, only one base is observed

Site-specific rates

Each defined subset (e.g. gene, codon position) has its own relative rate



r_1 applies to subset 1
(e.g. sites 1 - 1000)

r_2 applies to subset 2
(e.g. sites 1001-2000)

Relative rates have mean 1:

$$\frac{r_1 + r_2}{2} = 1$$

More generally:

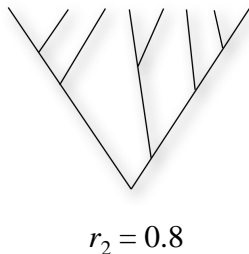
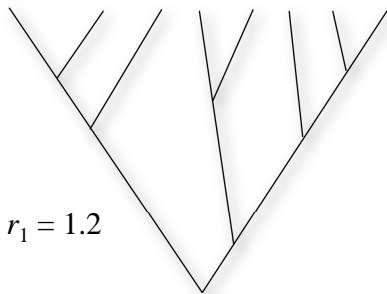
$$r_1 p(r_1) + r_2 p(r_2) = 1$$

Site-specific rates

$$L = \underbrace{\Pr(D_1|r_1) \cdots \Pr(D_{1000}|r_1)}_{\text{Gene 1}} \underbrace{\Pr(D_{1001}|r_2) \cdots \Pr(D_{2000}|r_2)}_{\text{Gene 2}}$$

Gene 1

Gene 2



Site-specific rates

JC69 transition probabilities that would be used for every site if rate *homogeneity* were assumed:

$$P_{ii}(t) = \frac{1}{4} + \frac{3}{4}e^{-4\alpha t}$$

$$P_{ij}(t) = \frac{1}{4} - \frac{1}{4}e^{-4\alpha t}$$

Site specific rates

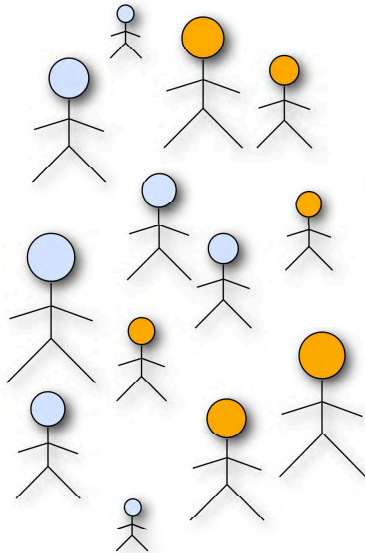
JC69 transition probabilities that would be used for sites in **gene 1**:

$$\begin{aligned}P_{ii}(t) &= \frac{1}{4} + \frac{3}{4}e^{-4r_1\alpha t} \\P_{ij}(t) &= \frac{1}{4} - \frac{1}{4}e^{-4r_1\alpha t}\end{aligned}$$

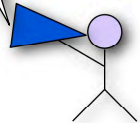
JC69 transition probabilities that would be used for sites in **gene 2**:

$$\begin{aligned}P_{ii}(t) &= \frac{1}{4} + \frac{3}{4}e^{-4r_2\alpha t} \\P_{ij}(t) &= \frac{1}{4} - \frac{1}{4}e^{-4r_2\alpha t}\end{aligned}$$

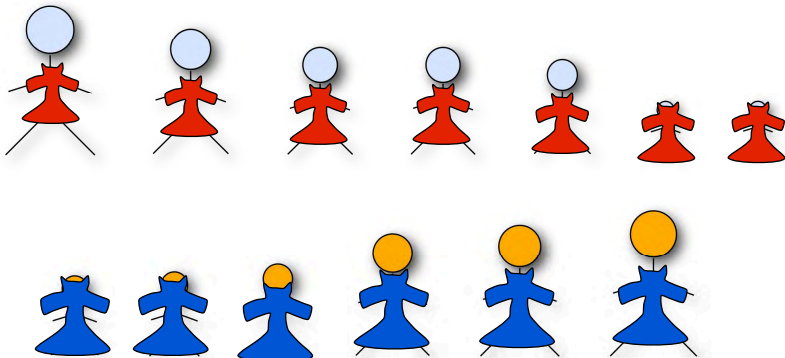
Site-specific Approach



Ok, I am going to divide you into 2 groups based on the color of your head, and everyone in each group will get a coat of the average size for their group. Very sorry if this does not work well for some people who are unusually large or small compared to their group.



Site-specific Approach

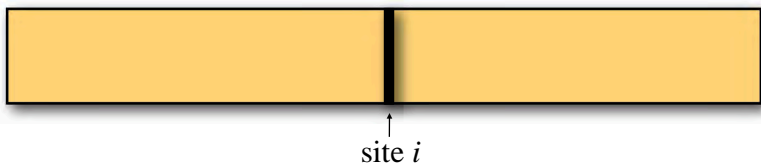


Good: costs less: need to buy just one coat for every person

Bad: every person in a group has to wear the same size coat, so the fit will be poor for some people if they are much bigger or smaller than the average size for the group in which they have been placed

Mixture Models

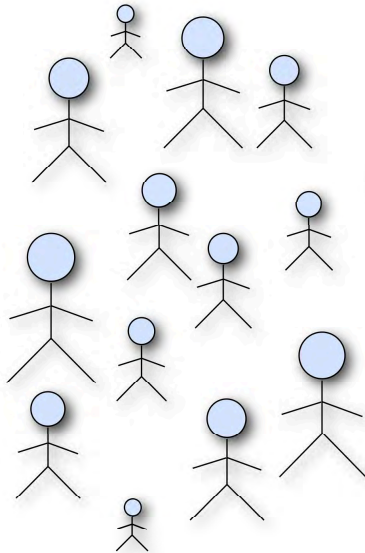
All relative rates applied to every site



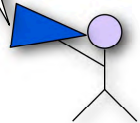
$$L_i = \Pr(D_i|r_1) \Pr(r_1) + \Pr(D_i|r_2) \Pr(r_2)$$

Common examples $\left\{ \begin{array}{l} \text{Invariable sites (I) model} \\ \text{Discrete Gamma (G) model} \end{array} \right.$

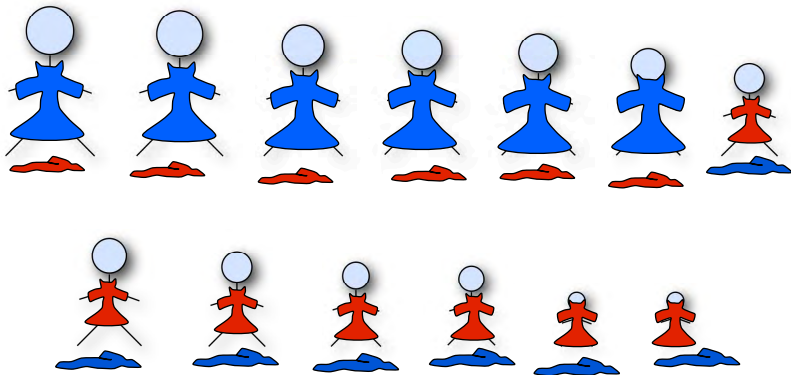
Mixture Model Approach



Ok, I am going to give each of you 2 coats: use the one that fits you best and throw away the other one. This costs twice as much for me, but on average leads to better fit for you. I have determined the two sizes of coats based on the distribution of your sizes.



Mixture Model Approach

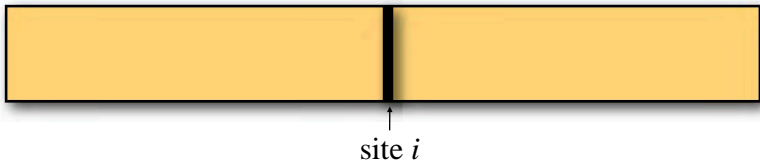


Good: every person experiences better fit because they can choose the size coat that fits best

Bad: costs more because two coats much be provided for each person

Invariable Sites Model

A fraction p_{invar} of sites are assumed to be invariable (i.e. rate = 0.0)



$$L_i = \Pr(D_i | r_1) p_{\text{invar}} + \Pr(D_i | r_2) (1 - p_{\text{invar}})$$

$$r_1 = 0.0$$

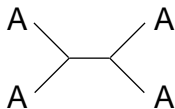
$$r_2 = \frac{1}{1 - p_{\text{invar}}}$$

Allows for the possibility that any given site could be variable or invariable

Reeves, J. H. 1992. Heterogeneity in the substitution process of amino acid sites of proteins coded for by mitochondrial DNA. *Journal of Molecular Evolution* 35:17-31.

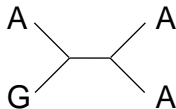
Invariable sites model

If site i is a *constant* site, both terms will contribute to the site likelihood:



$$L_i = \Pr(D_i|0.0)p_{\text{invar}} + \Pr(D_i|r_2)(1 - p_{\text{invar}})$$

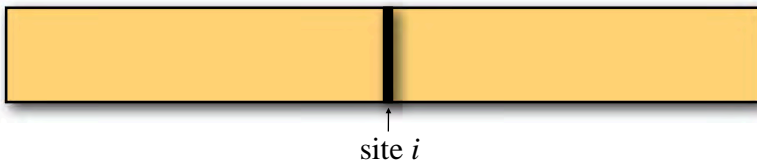
If site i is a *variable* site, there is no way to explain the data with a zero rate, so the first term is zero:



$$L_i = \cancel{\Pr(D_i|0.0)p_{\text{invar}}} + \Pr(D_i|r_2)(1 - p_{\text{invar}})$$

Discrete Gamma Model

No relative rate is exactly 0.0, and all are equally probable



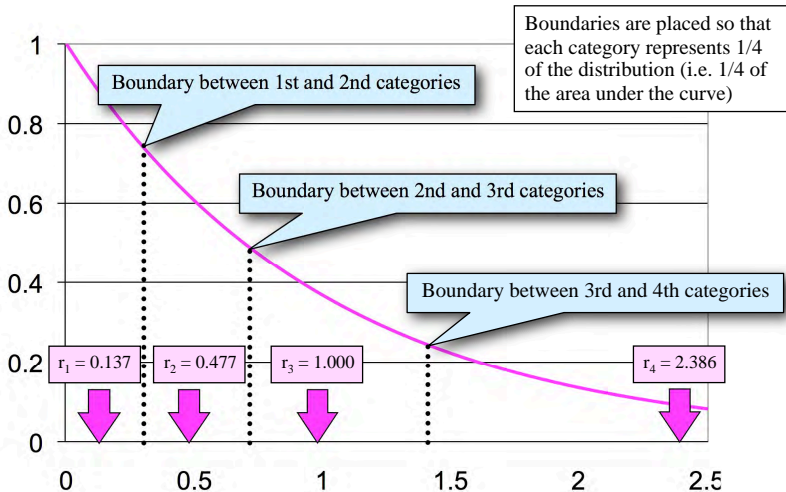
$$L = \left(\frac{1}{4}\right) \Pr(D_i|r_1) + \left(\frac{1}{4}\right) \Pr(D_i|r_2) + \left(\frac{1}{4}\right) \Pr(D_i|r_3) + \left(\frac{1}{4}\right) \Pr(D_i|r_4)$$

Relative rates are constrained to a discrete gamma distribution
Number of rate categories can vary (4 used here)

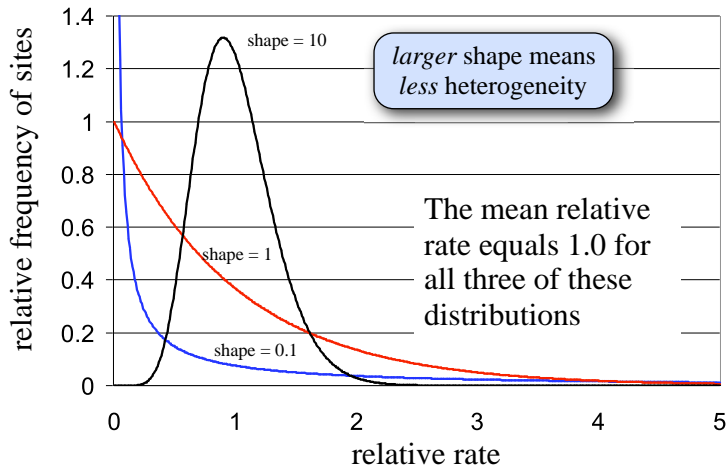
Yang, Z. 1993. Maximum-likelihood estimation of phylogeny from DNA sequences when substitution rates differ over sites. *Molecular Biology and Evolution* 10:1396-1401.

Yang, Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. *Journal of Molecular Evolution* 39:306-314.

Relative rates in 4-category case



Gamma distributions



Next class - lab on ML tree searching.