

Maximum Likelihood in Phylogenetics

26 January 2011

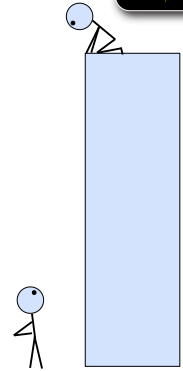
Workshop on Molecular Evolution
Český Krumlov, Česká republika

Paul O. Lewis
Department of Ecology & Evolutionary Biology
University of Connecticut, Storrs, CT

Goals

- Explain jargon
- Increase comfort level
- Provide background
- In other words...give a hand up

$$f(r) = \frac{r^{\alpha-1} e^{-r/\beta}}{\beta^\alpha \Gamma(\alpha)}$$



The Plan

- Probability review
- Likelihood
- Substitution models

- The AND and OR rules
- Independence of events

- What does it mean?
- Likelihood of a single sequence
- Maximum likelihood distances
- Likelihoods of trees

- Markov model basics
- Transition probabilities
- Survey of models
- Rate heterogeneity
- Codon models

Combining probabilities

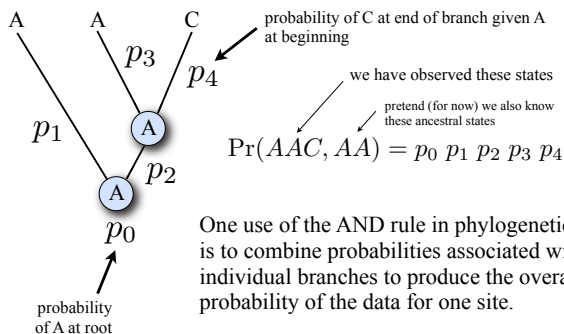
- Multiply* probabilities if the component events must happen **simultaneously** (i.e. where you would naturally use the word AND when describing the problem)

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)

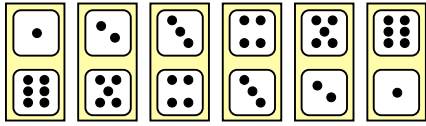
Using one die, what is the probability of



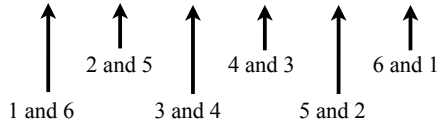
$$(1/6) + (1/6) = 1/3$$

Combining AND and OR

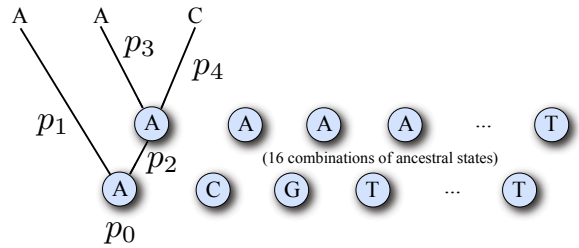
What is the probability that the sum of two dice is 7?



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



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...

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

joint probability
conditional probability

If we can say this...

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

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

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

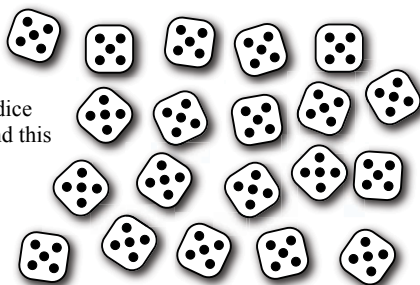
Likelihood

The Likelihood Criterion

The probability of the observations computed using a model tells us how surprised we should be.

The preferred model is the one that surprises us least.

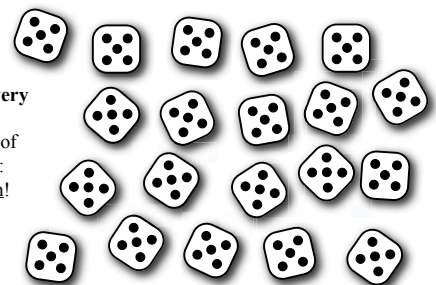
Suppose I threw 20 dice down on the table and this was the result...



The Fair Dice model

$$\Pr(\text{obs.}|\text{fair dice model}) = \left(\frac{1}{6}\right)^{20} = \frac{1}{3,656,158,440,062,976}$$

You should have been **very surprised** at this result because the probability of this event is **very small**: only 1 in 3.6 **quadrillion**!

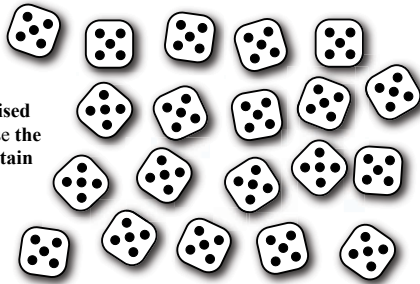


The Trick Dice model

(assumes dice each have 5 on every side)

$$\Pr(\text{obs.}|\text{trick dice model}) = 1^{20} = 1$$

You should **not be surprised at all** at this result because the **observed outcome is certain** under this model



Results

Model	Likelihood	Surprise level
Fair Dice	$\frac{1}{3,656,158,440,062,976}$	Very, very, very surprised
Trick Dice	1.0	Not surprised at all

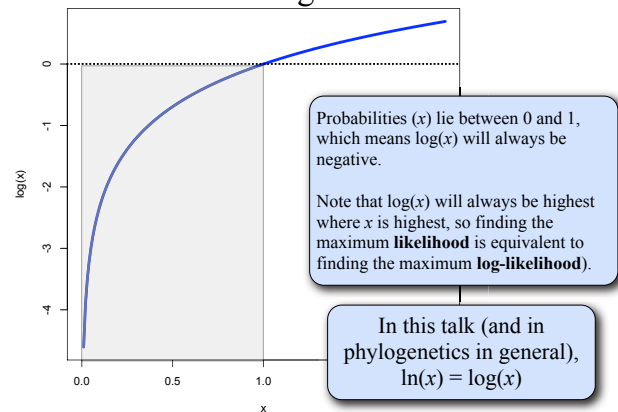
winning model maximizes likelihood (and thus minimizes surprise)

Likelihood and model comparison

- Analyses using likelihoods ultimately involve **model comparison**
- The models compared can be **discrete** (as in the fair vs. trick dice example)
- More often the models compared differ **continuously**:
 - Model 1: branch length is 0.05
 - Model 2: branch length is 0.06

Rather than having an infinity of models, we instead think of the branch length as a **parameter** within one model

Likelihoods vs. log-likelihoods



Likelihood of a single sequence

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

GAAGTCCTTGAGAAATAAACTGCACACACTGG

$$L = \pi_A^{12} \pi_C^7 \pi_G^7 \pi_T^6$$

Note that we are assuming independence between 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 (e.g. F81) allowing **unequal** base frequencies will **fit better** than a model that assumes **equal** base frequencies (e.g. JC69) because there are about twice as many As as there are Cs, Gs and Ts.

Likelihood ratio test

Likelihood ratio tests can be used to evaluate whether an **unconstrained** model fits the data significantly better than a **constrained** version of the same model.

Find *maximum* logL under F81 (unconstrained) model:

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

Find *maximum* logL under JC69 (constrained) model:

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

Likelihood ratio test

Calculate the likelihood ratio test statistic:

$$R = -2 \left[\log(L_{JC69}) - \log(L_{F81}) \right]$$

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

$$= -2 \left[-44.4 - (-43.1) \right] = 2.6$$

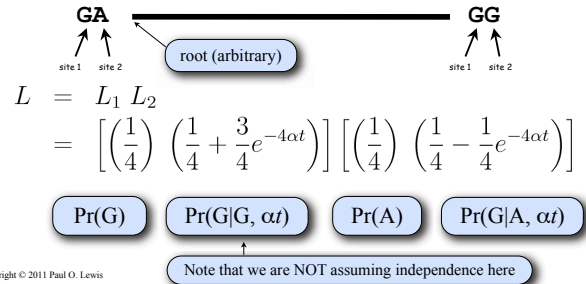
F81 does fit better (-43.1 > -44.4), but not significantly better (P = 0.457, chi-squared with 3 d.f.*)

*The number of degrees of freedom equals the difference between the two models in the number of free parameters. In this case, F81 has 3 parameters and JC69 has 0, so d.f. = 3 - 0 = 3

Likelihood of the simplest tree

sequence 1 ————— sequence 2

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

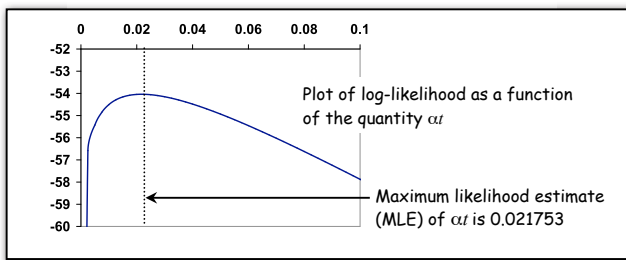


Maximum likelihood estimation

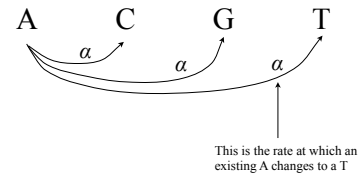
First 32 nucleotides of the $\psi\eta$ -globin gene of gorilla and orangutan:

gorilla **GAAGTCCTTGAGAAATAAACTGCACACACTGG**
 orangutan **GGACTCCTTGAGAAATAAACTGCACACACTGG**

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



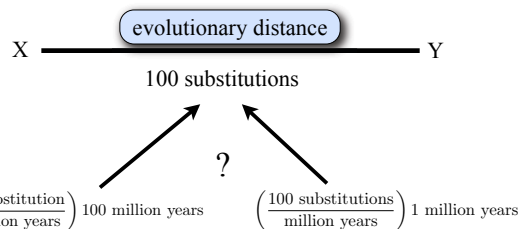
number of substitutions = rate \times time



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

$$v = 3\alpha t$$

Rate and time are confounded



Andrew will discuss (next week) models in which constraints on times can be used to infer rates (and vice versa), but without some extra information or constraints, sequence data allow only estimation of the **number** of substitutions.

A convenient convention

Because rate and time are confounded, it is convenient to arbitrarily standardize things by setting the rate to a value such that **one substitution** is expected to occur in **one unit of time** for each site.

This results in "time" (the length of a branch) being measured in units of **evolutionary distance (expected number of substitutions per site)** rather than years (or some other calendar unit).

evolutionary distance $v = 3\alpha t$

$$v = 3 \left(\frac{1}{3} \right) t$$

Setting $\alpha=1/3$ results in v equalling t

Evolutionary distances for several common models

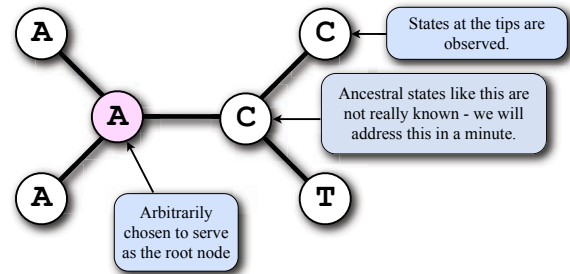
Model	Expected no. substitutions: $v = \{r\}t$
JC69	$v = \{3\alpha\}t$
F81	$v = \{2\mu(\pi_R\pi_Y + \pi_A\pi_G + \pi_C\pi_T)\}t$
K80	$v = \{\beta(\kappa + 2)\}t$
HKY85	$v = \{2\mu[\pi_R\pi_Y + \kappa(\pi_A\pi_G + \pi_C\pi_T)]\}t$

In the formulas above, the overall rate r (in curly brackets) is a function of all parameters in the substitution model.

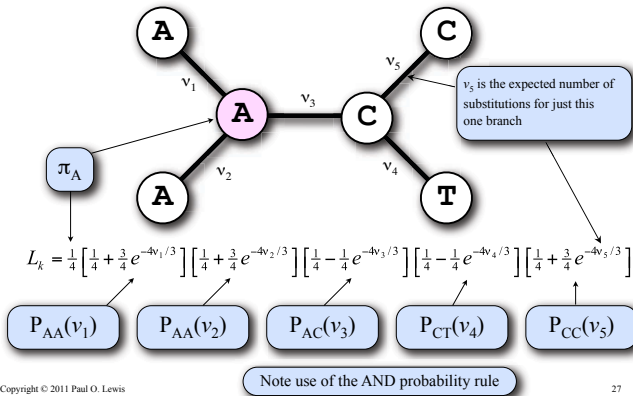
One of the parameters of the substitution model can always be determined from the branch length (using our convention that $v = t$).

Likelihood of an unrooted tree

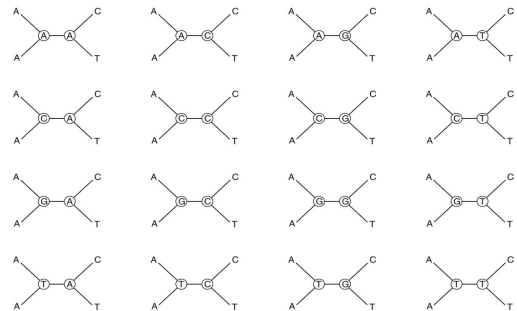
(data shown for only one site)



Likelihood for site k

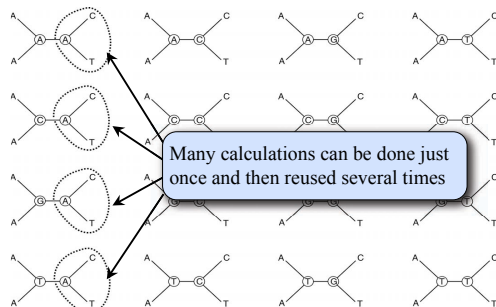


Brute force approach would be to calculate L_k for all 16 combinations of ancestral states and sum them



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

Jukes-Cantor (JC69) model

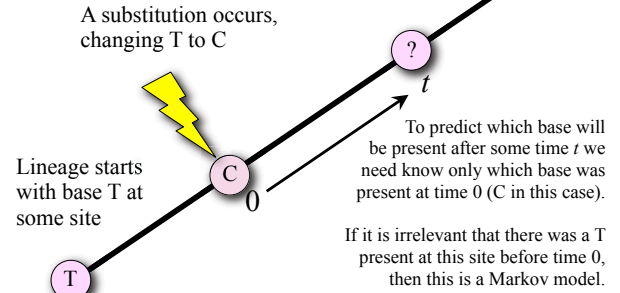
- The four bases (A, C, G, T) are expected to be **equally frequent** in sequences ($\pi_A = \pi_C = \pi_G = \pi_T = 0.25$)
- Assumes **same rate** for all types of substitution ($r_{A \rightarrow C} = r_{A \rightarrow G} = r_{A \rightarrow T} = r_{C \rightarrow G} = r_{C \rightarrow T} = r_{G \rightarrow T} = \alpha$)
- Usually described as a **1-parameter** model (the parameter being the branch length)
 - Remember, however, that each branch 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.

Copyright © 2011 Paul O. Lewis

31

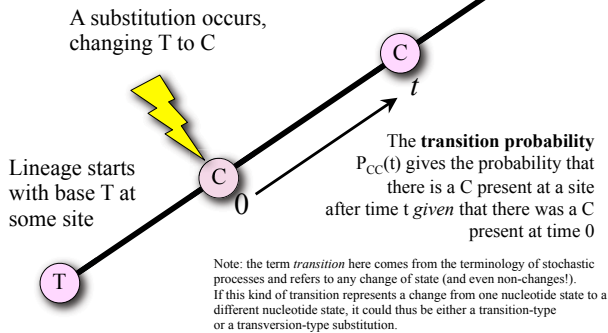
What is a Markov model?



Copyright © 2011 Paul O. Lewis

32

Transition Probabilities



Copyright © 2011 Paul O. Lewis

33

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 :

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

The JC69 model has only one unknown quantity: αt

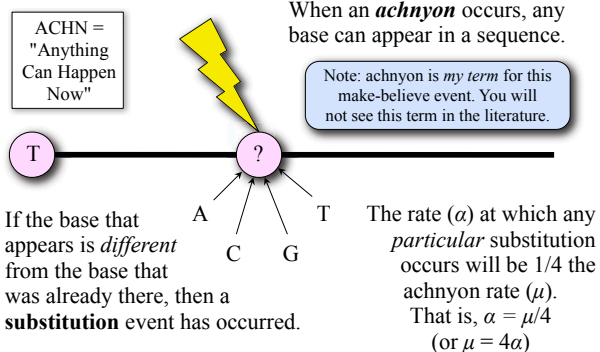
(The symbol e represents the base of the natural logarithms: its value is 2.718281828459045...)

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

Copyright © 2011 Paul O. Lewis

34

"ACHNyons" vs. substitutions



Copyright © 2011 Paul O. Lewis

35

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 :

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

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

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

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

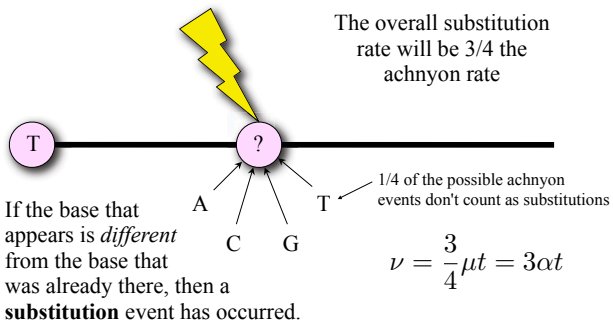
Remember that the rate (α) of any particular substitution is one fourth the achnyon rate (μ):

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

Copyright © 2011 Paul O. Lewis

36

Expected number of substitutions



Transition Probabilities: Remarks

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

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

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

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

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

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

Doh! Something must be wrong here...

Transition Probabilities: Remarks

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

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

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

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

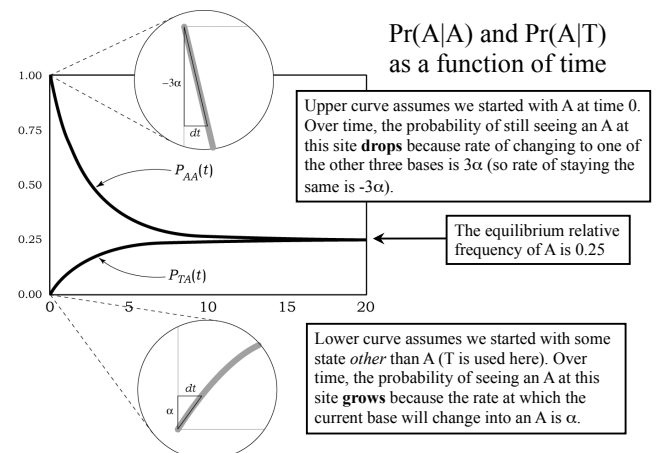
Forgot to account for the possibility of *no* achnyons over time *t*

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 Frequency Simulation

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



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.

Copyright © 2011 Paul O. Lewis

43

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.

Copyright © 2011 Paul O. Lewis

44

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/transversion rate ratio* $\rightarrow \kappa = \frac{\alpha}{\beta}$

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

Copyright © 2011 Paul O. Lewis

45

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

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

Copyright © 2011 Paul O. Lewis

46

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.

Copyright © 2011 Paul O. Lewis

47

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 hominoidae. *Journal of Molecular Evolution* 29: 170-179.

Copyright © 2011 Paul O. Lewis

48

Site specific rates

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

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

$$P_{ij}(t) = \frac{1}{4} - \frac{1}{4}e^{-4r_1\alpha t}$$

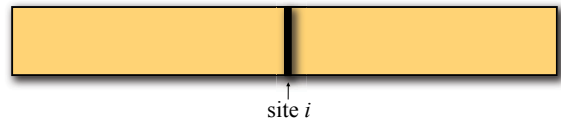
JC69 transition probabilities that would be used for sites in **gene 2**:

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

$$P_{ij}(t) = \frac{1}{4} - \frac{1}{4}e^{-4r_2\alpha t}$$

Mixture Models

All relative rates applied to every site

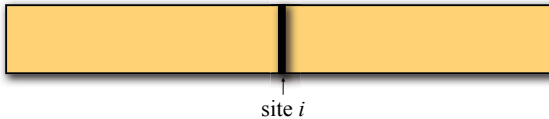


$$L_i = \Pr(D_i|r_1) \Pr(r_1) + \Pr(D_i|r_2) \Pr(r_2)$$

Common examples { Invariable sites (I) model
Discrete Gamma (G) model

Invariable Sites Model

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



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

$$r_1 = 0.0$$

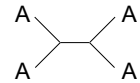
$$r_2 = \frac{1}{1 - p_{invar}}$$

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.

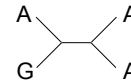
Invariable sites model

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



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

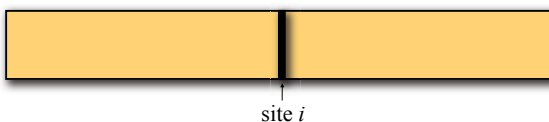
If site *i* is a *variable* site, there is no way to explain the data with a zero rate, so the first term is zero:



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

Discrete Gamma Model

No relative rates are exactly 0.0, and all are equally probable



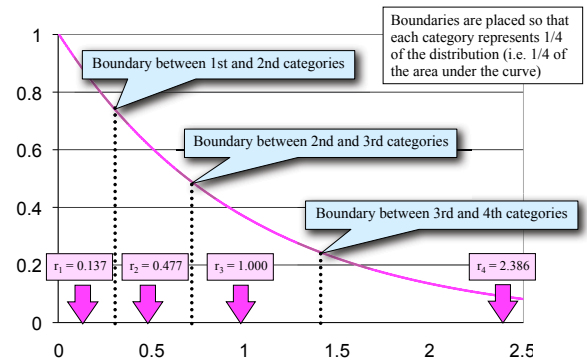
$$L = \left(\frac{1}{4}\right) \Pr(D_i|r_1) + \left(\frac{1}{4}\right) \Pr(D_i|r_2) + \left(\frac{1}{4}\right) \Pr(D_i|r_3) + \left(\frac{1}{4}\right) \Pr(D_i|r_4)$$

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

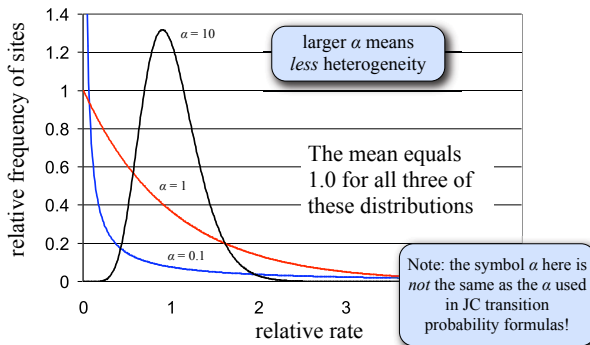
Yang, Z. 1993. Maximum-likelihood estimation of phylogeny from DNA sequences when substitution rates differ over sites. *Molecular Biology and Evolution* 10:1396-1401.

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



Copyright © 2011 Paul O. Lewis

61

Codon models

Copyright © 2011 Paul O. Lewis

62

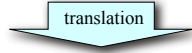
The Genetic Code

First 12 nucleotides at the 5' end of the *rbcl* gene in corn:

5' -ATG | TCA | CCA | CAA-3' coding strand } DNA double helix
 3' -TAC | AGT | GGT | GTT-5' template strand



5' -AUG | UCA | CCA | CAA-3' mRNA



N-Met | Ser | Pro | Gln-C polypeptide
 Codon models treat codons as the independent units, not individual nucleotide sites.

Genetic Code

	U	C	A	G	
U	UUU Phe UUC Phe UUA Leu UUG Leu	UCU Ser UCC Ser UCA Ser UCG Ser	UAU Tyr UAC Tyr UAA Stop UAG Stop	UGU Cys UGC Cys UGA Stop UGG Trp	U C A G
C	CUU Leu CUC Leu CUA Leu CUG Leu	CCU Pro CCC Pro CCA Pro CCG Pro	CAU His CAC His CAA Gln CAG Gln	CGU Arg CGC Arg CGA Arg CGG Arg	U C A G
A	AUU Ile AUC Ile AUA Ile AUG Met	ACU Thr ACC Thr ACA Thr ACG Thr	AUU Asn AAC Asn AAA Lys AAG Lys	AGU Ser AGC Ser AGA Arg AGG Arg	U C A G
G	GUU Val GUC Val GUA Val GUG Val	GCU Ala GCC Ala GCA Ala GCG Ala	GAU Asp GAC Asp GAA Glu GAG Glu	GGU Gly GGC Gly GGA Gly GGG Gly	U C A G

<http://www.ncbi.nlm.nih.gov/Taxonomy/NCBI/Taxonomy/Html/genbank/code.html>

Copyright © 2011 Paul O. Lewis

63

First codon models

- Muse and Gaut model (MG94) is simplest
 - α = synonymous substitution rate
 - β = nonsynonymous substitution rate
 - $\pi_A, \pi_C, \pi_G, \pi_T$ = base frequencies
- Goldman and Yang model (GY94) similar
 - accounts for synon./nonsynon. *and* trs/trv bias *and* amino acid properties (later simplified, see Yang et al. 1998)

Muse, S. V., and B. S. Gaut. 1994. A likelihood approach for comparing synonymous and nonsynonymous substitution rates, with application to the chloroplast genome. *Molecular Biology and Evolution* 11:715-724.

Goldman, N., and Z. Yang. 1994. A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Molecular Biology and Evolution* 11:725-736.

Yang, Z., Nielsen, R., and Hasegawa, M. 1998. Models of amino acid substitution and applications to mitochondrial protein evolution. *Molecular Biology and Evolution* 15:1600-1611.

Copyright © 2011 Paul O. Lewis

64

Table I. Part of Muse and Gaut's 61 × 61 instantaneous rate matrix^a

Codon before substitution (the 'from' state)	Codon after substitution (the 'to' state)							
	TTT (Phe)	TTC (Phe)	TTA (Leu)	TTG (Leu)	CTT (Leu)	CTC (Leu)	...	GGG (Gly)
TTT (Phe)	---	$\alpha\pi_C$	$\beta\pi_A$	$\beta\pi_G$	$\beta\pi_C$	0	...	0
TTC (Phe)	$\alpha\pi_T$	---	$\beta\pi_A$	$\beta\pi_G$	0	$\beta\pi_C$...	0
TTA (Leu)	$\beta\pi_T$	$\beta\pi_C$	---	$\alpha\pi_G$	0	0	...	0
TTG (Leu)	$\beta\pi_T$	$\beta\pi_C$	$\alpha\pi_A$	---	0	0	...	0
CTT (Leu)	0	0	0	---	$\alpha\pi_C$...	0	0
CTC (Leu)	0	$\beta\pi_T$	0	0	$\alpha\pi_T$	---	...	0
...
GGG (Gly)	0	0	0	0	0	0	...	---

Note that it is still easy for the change CTT → TTA to occur, it just requires more than one instant of time

Instantaneous rate is 0.0 if two or more nucleotides must change during the codon transition

Table I from: Lewis, P. O. 2001. Phylogenetic systematics turns over a new leaf. *Trends in Ecology and Evolution* 16:30-37.

Copyright © 2011 Paul O. Lewis

65

Interpreting codon model results

$\omega = \beta/\alpha$ is the nonsynonymous/synonymous rate ratio

omega	mode of selection	example(s)
$\omega < 1$	stabilizing selection (nucleotide substitutions rarely change the amino acid)	functional protein coding genes
$\omega = 1$	neutral evolution (synonymous and nonsynonymous substitutions occur at the same rate)	pseudogenes
$\omega > 1$	positive selection (nucleotide substitutions often change the amino acid)	envelope proteins in viruses under active positive selection

Copyright © 2011 Paul O. Lewis

66