

# Molecular Evolution

Malay K Basu  
Department of Pathology  
[malay@uab.edu](mailto:malay@uab.edu)

# *Ad Verecundiam*


Because Dobzhansky says so.




*“Nothing in Biology Makes Sense Except  
in the Light of Evolution”*

- Theodosius Dobzhansky 1973

# Classification of all living organisms



Kingdom  
Phylum  
Class  
Order  
Family  
Genus  
Species




Carl Linnaeus (1707-1778)


# Quiz

Do you "know thyself"?

I think



Thus between A + B. various  
sorts of relation. C + B. The  
first gradation. B + D  
rather greater distinction  
thus genera would be  
formed. - binary relation





**I think...>,**

- Charles Darwin  
8B Notebook p. 37<sup>9</sup>

Unity of Descent  
Last Universal Common Ancestor  
(LUCA)

"Common descent with modification"




Earnst Haeckel (1834-1919)

# Tree kingdoms of Life




# Carl Richard Woese (1928-2012)



# Horizontal Gene Transfer

Transfer of genetic material from surroundings to genome

# Phylogenetic forest (many trees)



# 3 pillars of evolution

- 1 .Mutation - Random error in DNA replication
- 2 .Selection - Increase/decrease fitness
- 3 .Drift - Random fluctuation in allele frequency

# Bone in classical evolutionary theories



Most changes in DNA are  
"neutral".

Genetic drift is the major cause of evolution

Earnst Haeckel (1924-1994)

# "Molecular" evolution

Evolutionary changes in molecules: DNA and protein sequences

# Mutational changes in DNA

(A) Substitution

Thr	Tyr	Leu	Leu
ACC	<b>TAT</b>	TTG	CTG
	↓		
ACC	<b>TCT</b>	TTG	CTG

Thr Ser Leu Leu

(C) Insertion

Thr	Tyr	Leu	Leu
ACC	<b>TAT</b>	TTG	CTG
	↓		
ACC	<b>TAC</b>	TTT	GCT G--

Thr Tyr Phe Ala

(B) Deletion

Thr	Tyr	Leu	Leu
ACC	<b>TAT</b>	TTG	CTG
	↓		
ACC	<b>TAT</b>	<b>TGC</b>	TG-

Thr Tyr Cys

(D) Inversion

Thr	Tyr	Leu	Leu
ACC	<b>TAT</b>	TTG	CTG
	↓		
ACC	<b>TTT</b>	<b>ATG</b>	CTG

Thr Phe Met Leu

FIGURE 1.2. Four basic types of mutation at the nucleotide level. Nucleotide sequences are presented in units of codons or nucleotide triplets in order to show how the amino acids encoded are affected by the nucleotide changes. The nucleotides affected by the mutational changes are shown in boldface.

# Nucleotide substitution




FIGURE 1.3. Transitional ( $A \leftrightarrow G$  and  $T \leftrightarrow C$ ) and transversional (others) nucleotide substitutions.  $\alpha$  and  $\beta$  are the rates of transitional and transversional substitutions, respectively.

# Standard Genetic Code

TTT F Phe	TCT S Ser	TAT Y Tyr	TGT C Cys
TTC F Phe	TCC S Ser	TAC Y Tyr	TGC C Cys
TTA L Leu	TCA S Ser	TAA * Ter	TGA * Ter
TTG L Leu i	TCG S Ser	TAG * Ter	TGG W Trp
CTT L Leu	CCT P Pro	CAT H His	CGT R Arg
CTC L Leu	CCC P Pro	CAC H His	CGC R Arg
CTA L Leu	CCA P Pro	CAA Q Gln	CGA R Arg
CTG L Leu i	CCG P Pro	CAG Q Gln	CGG R Arg
ATT I Ile	ACT T Thr	AAT N Asn	AGT S Ser
ATC I Ile	ACC T Thr	AAC N Asn	AGC S Ser
ATA I Ile	ACA T Thr	AAA K Lys	AGA R Arg
ATG M Met i	ACG T Thr	AAG K Lys	AGG R Arg
GTT V Val	GCT A Ala	GAT D Asp	GGT G Gly
GTC V Val	GCC A Ala	GAC D Asp	GGC G Gly
GTA V Val	GCA A Ala	GAA E Glu	GGA G Gly
GTG V Val	GCG A Ala	GAG E Glu	GGG G Gly

# Codon bias

Phe UUU	<b>15 (0.51)</b>	Ser UCU	<b>32 (1.86)</b>	Tyr UAU	<b>18 (0.64)</b>	Cys UGU	<b>5 (1.00)</b>
<b>UUC</b>	<b>44 (1.49)</b>	UCC	<b>38 (2.21)</b>	<b>UAC</b>	<b>38 (1.36)</b>	UGC	<b>5 (1.00)</b>
Leu UUA	<b>2 (0.07)</b>	UCA	<b>2 (0.12)</b>	Ter UAA		Ter UGA	
UUG	<b>8 (0.27)</b>	UCG	<b>5 (0.29)</b>	Ter UAG		Trp UGG	<b>8 (1.00)</b>
Leu CUU	<b>11 (0.36)</b>	Pro CCU	<b>9 (0.48)</b>	His CAU	<b>5 (0.36)</b>	Arg CGU	<b>89 (3.93)</b>
CUC	<b>18 (0.60)</b>	CCC	<b>0 (0.00)</b>	CAC	<b>23 (1.64)</b>	CGC	<b>46 (2.03)</b>
CUA	<b>1 (0.03)</b>	CCA	<b>11 (0.59)</b>	Gln CAA	<b>15 (0.34)</b>	CGA	<b>1 (0.04)</b>
<b>CUG</b>	<b>141 (4.67)</b>	CCG	<b>55 (2.93)</b>	<b>CAG</b>	<b>73 (1.66)</b>	CGG	<b>0 (0.00)</b>
Ile AUU	<b>29 (0.69)</b>	Thr ACU	<b>19 (0.78)</b>	Asn AAU	<b>4 (0.11)</b>	Ser AGU	<b>3 (0.17)</b>
<b>AUC</b>	<b>98 (2.31)</b>	ACC	<b>63 (2.57)</b>	AAC	<b>66 (1.89)</b>	AGC	<b>23 (1.34)</b>
AUA	<b>0 (0.00)</b>	ACA	<b>3 (0.12)</b>	Lys AAA	<b>77 (1.35)</b>	Arg AGA	<b>0 (0.00)</b>
Met AUG	<b>60 (1.00)</b>	ACG	<b>13 (0.53)</b>	AAG	<b>37 (0.65)</b>	AGG	<b>0 (0.00)</b>
Val GUU	<b>55 (1.53)</b>	Ala GCU	<b>30 (0.94)</b>	Asp GAU	<b>60 (0.83)</b>	Gly GGU	<b>78 (2.40)</b>
GUC	<b>21 (0.58)</b>	GCC	<b>19 (0.59)</b>	GAC	<b>85 (1.17)</b>	GGC	<b>47 (1.45)</b>
<b>GUA</b>	<b>34 (0.94)</b>	<b>GCA</b>	<b>30 (0.94)</b>	Glu GAA	<b>147 (1.52)</b>	GGA	<b>0 (0.00)</b>
<b>GUG</b>	<b>34 (0.94)</b>	<b>GCG</b>	<b>49 (1.53)</b>	GAG	<b>46 (0.48)</b>	GGG	<b>5 (0.15)</b>

FIGURE 1.4. Codon frequencies observed in the RNA polymerase genes (*rpo B* and *D* genes) of the bacterium *Escherichia coli*. The codons optimal for the translational system are shown in boldface. Relative synonymous codon usages (RSCU) given in the parentheses were computed by Equation (1.1). Data from Ikemura (1985).


# Result of substitution

- Synonymous - Does not change the AA
- Non-synonymous (missense) - Changes the AA
- Nonsense - Creates a stop codon

# P Distance

No of times two sequences differ

# Sequence divergence with time



# Models of nucleotide substitutions

Table 3.2 Models of nucleotide substitution.

	A	T	C	G		A	T	C	G
(A) Jukes-Cantor model					(E) HKY model				
A	-	$\alpha$	$\alpha$	$\alpha$	A	-	$\beta g_T$	$\beta g_C$	$\alpha g_G$
T	$\alpha$	-	$\alpha$	$\alpha$	$\beta g_A$	-	-	$\alpha g_C$	$\beta g_G$
C	$\alpha$	$\alpha$	-	$\alpha$	$\beta g_A$	$\alpha g_T$	-	-	$\beta g_G$
G	$\alpha$	$\alpha$	$\alpha$	-	$\alpha g_A$	$\beta g_T$	$\beta g_C$	-	-
(B) Kimura model					(F) Tamura-Nei model				
A	-	$\beta$	$\beta$	$\alpha$	A	-	$\beta g_T$	$\beta g_C$	$\alpha_1 g_G$
T	$\beta$	-	$\alpha$	$\beta$	$\beta g_A$	-	-	$\alpha_2 g_C$	$\beta g_G$
C	$\beta$	$\alpha$	-	$\beta$	$\beta g_A$	$\alpha_2 g_T$	-	-	$\beta g_G$
G	$\alpha$	$\beta$	$\beta$	-	$\alpha_1 g_A$	$\beta g_T$	$\beta g_C$	-	-
(C) Equal-input model					(G) General reversible model				
A	-	$\alpha g_T$	$\alpha g_C$	$\alpha g_G$	A	-	$ag_T$	$bg_C$	$cg_G$
T	$\alpha g_A$	-	$\alpha g_C$	$\alpha g_G$	$ag_A$	-	$dg_C$	$eg_G$	-
C	$\alpha g_A$	$\alpha g_T$	-	$\alpha g_G$	$bg_A$	$dg_T$	-	$fg_G$	-
G	$\alpha g_A$	$\alpha g_T$	$\alpha g_C$	-	$cg_A$	$eg_T$	$fg_C$	-	-
(D) Tamura model					(H) Unrestricted model				
A	-	$\beta\theta_2$	$\beta\theta_1$	$\alpha\theta_1$	A	-	$a_{12}$	$a_{13}$	$a_{14}$
T	$\beta\theta_2$	-	$\alpha\theta_1$	$\beta\theta_1$	$a_{21}$	-	-	$a_{23}$	$a_{24}$
C	$\beta\theta_2$	$\alpha\theta_2$	-	$\beta\theta_1$	$a_{31}$	$a_{32}$	-	-	$a_{34}$
G	$\alpha\theta_2$	$\beta\theta_2$	$\beta\theta_1$	-	$a_{41}$	$a_{42}$	$a_{43}$	-	-

Note: An element  $(e_{ij})$  of the above substitution matrices stands for the substitution rate from the nucleotide in the  $i$ -th row to the nucleotide in the  $j$ -th column.  $g_A$ ,  $g_T$ ,  $g_C$ , and  $g_G$  are the nucleotide frequencies.  $\theta_1 = g_G + g_C$ .  $\theta_2 = g_A + g_T$ .

# Empirical AA substitution table

PAM (Point Accepted Mutation)

- Created by Margaret Dayhoff
- Different matrix for different evolutionary distance

BLOSUM

# BLOSUM62

# How BLOSUM is calculated

$$S_{ij} = \frac{1}{\lambda} \log \left( \frac{f_{ij}}{f_i \times f_j} \right)$$

$\lambda$  = a scaling parameter

$f_{ij}$  = frequency of number of times one AA changes to another

$f_i, f_j$  = frequency of each AA

# Quiz


Why the self substitution scores are different for different amino acids?

Find the score of PQG  
matching PQG using  
**BLOSUM62**

# Homologs

Genes related by evolution.

# Homologs



All 4 : **homologs**

Species1 (a1,a2) and  
Species 2 (a1, a2):  
**orthologs**

a1 and a2: **paralogs**



Fitch W. (1970). "Distinguishing homologous from analogous proteins".  
*Syst Zool* 19 (2): 99–113.

## DISTINGUISHING HOMOLOGOUS FROM ANALOGOUS PROTEINS

WALTER M. FITCH

### Abstract

Fitch, W. M. (Dept. Physiological Chem., U. Wisconsin, Madison 53706) 1970. Distinguishing homologous from analogous proteins. *Syst. Zool.*, 19:99–113.—This work provides a means by which it is possible to determine whether two groups of related proteins have a common ancestor or are of independent origin. A set of 16 random amino acid sequences were shown to be unrelated by this method. A set of 16 real but presumably unrelated proteins gave a similar result. A set of 24 model proteins which was composed of two independently evolving groups, converging toward the same chemical goal, was correctly shown to be convergently related, with the probability that the result was due to chance being  $<10^{-n}$ . A set of 24 cytochromes composed of 5 fungi and 19 metazoans was shown to be divergently related, with the probability that the result was due to chance being  $< 10^{-n}$ . A process was described which leads to the absolute minimum of nucleotide replacements required to account for the divergent descent of a set of genes given a particular topology for the tree depicting their ancestral relations. It was also shown that the convergent processes could realistically lead to amino acid sequences which would produce positive tests for relatedness, not only by a chemical criterion, but by a genetic (nucleotide sequence) criterion as well. Finally, a realistic case is indicated where truly homologous traits, behaving in a perfectly expectable way, may nevertheless lead to a ludicrous phylogeny.

The demonstration that two proteins are related has been attempted using two different criteria. One criterion is to show that their chemical structures are very similar. An early example of this approach was the observation of the relatedness of the oxygen carrying proteins, myoglobin and hemoglobin (Watson and Kendrew, 1961). More recent is the relatedness of two enzymes in carbohydrate metabolism, lysozyme and alpha-lactalbumin (Brew, Vanaman and Hill, 1967). The other criterion is to show that underlying genetic structures of the proteins are more alike than one would expect by chance. This is now possible because our knowledge of the genetic code permits us to determine how many nucleotide positions, at the minimum, must differ in the genes encoding the two presumptively homologous proteins. One then compares the answer obtained to the number of differences one would expect for unrelated proteins. An example of this approach is the observation of the relatedness of plant and bacterial ferredoxins (Matsubara,

Jukes and Cantor, 1969) for which added evidence has been produced (Fitch, 1970a). But regardless of the approach, the impulse, too powerful to resist, is to conclude that a particular pair of proteins had a common genic ancestor if they meet whichever criterion the observer uses.

Now two proteins may appear similar because they descend with divergence from a common ancestral gene (i.e., are homologous in a time-honoured meaning dating back at the least to Darwin's *Origin of Species*) or because they descend with convergence from separate ancestral genes (i.e., are analogous). And, if a common genic ancestor is to be the conclusion, a genetic criterion should be superior to a chemical criterion. This is because analogous gene products, although they have no common ancestor, do serve similar functions and may well be expected to have similar chemical structures and thereby be confused with homologous gene products. This danger can only be increased by using a chemical, as opposed to a genetic, criterion.

**Sequence similarity is not  
homology**

# Homology vs Homoplasy



# Detecting selection

# Types of selection

Purifying /Negative selection - Does not allow change

Positive/Adaptive selection - Faster change

Neutral

# How to measure selection

$d_n = \text{Non-synonymous substitutions} / \text{non-synonymous site}$

$d_s = \text{synonymous substitutions} / \text{synonymous site}$

$d_n/d_s > 1$  = Positive selection

$d_n/d_s < 1$  = Negative selection

$d_n/d_s = 1$  = Neutral selection

1/3 synonymous  
2/3 nonsynonymous  
nucleotide site

1 non-synonymous  
nucleotide site


1 synonymous  
nucleotide site

ATA    GTA    TTA  
(Ile)    (Val)    (Leu)

CGA    CCA    CAA  
(Arg)    (Pro)    (Gln)

CTC    CTG    CTT  
(Leu)    (Leu)    (Leu)

... GGT AGG CCA CTA AAT CGA TTA ...  
(Leu)



T	P	N	G	A	L	E	L	K	P	V	R
ACT	CCG	AAC	GGGGCG	TTAGAGTT	GAAACCCG	TTAGA					
*	*	*	*	*	*	*	*	*	*	**	
ACG	CCG	ATC	GGCGCG	ATAGGGTT	CAAGCTCGT	ACGA					
T	P	I	G	A	I	G	F	K	L	V	R

syn    00100100 $\frac{1}{2}$ 001001 $\frac{1}{4}$ 0 $\frac{1}{2}$ 00 $\frac{1}{3}\frac{1}{3}$ 0 $\frac{1}{3}$ 00 $\frac{1}{3}$ 001001 $\frac{1}{3}$ 0 $\frac{2}{3}$       sum = 7.5833

non    11011011 $\frac{1}{2}$ 110110 $\frac{3}{4}$ 1 $\frac{1}{2}$ 11 $\frac{2}{3}\frac{2}{3}$ 1 $\frac{2}{3}$ 11 $\frac{2}{3}$ 110100 $\frac{2}{3}$ 1 $\frac{1}{3}$       sum = 28.4167

$$dN = \frac{\text{No. non-synonymous substitutions}}{\text{No. non-synonymous sites}} = \frac{5}{28.417} = 0.176$$

$$dS = \frac{\text{No. synonymous substitutions}}{\text{No. synonymous sites}} = \frac{5}{7.583} = 0.659$$

The ratio is then


$$\frac{dN}{dS} = \frac{0.176}{0.659} = 0.269$$

# Phylogeny


# Evolutionary Tree

A graph structure showing the relationship amongst species or in case of genes, relationship amongst gene.

### A. Rooted tree



### B. Unrooted tree



**Figure 6.5.** Structure of evolutionary trees.

# Tree features

Taxon (plural taxa) are atomic units of the tree.

Branch length represent the estimate of the sequence change.

Each internal node represent a speciation even.

# Tree features


The branch length may differ due to “accelerated evolution” after speciation.

May phylogenetic techniques assume that the branch lengths are same “molecular clock”. Such assumption is only valid for closely related species.

# Rooted trees are hard to make

**Table 6.2.** *Number of possible evolutionary trees to consider as a function of number of sequences*

Taxa or sequence no.	No. of rooted trees	No. of unrooted trees
3	3	1
4	15	3
5	105	15
—	—	—
7	10,395	954



# Rooted tree

Root represent common ancestor of all nodes.

In general, root is fixed by a taxon that branched off earlier than the others “outgroup”.

Root can also be predicted provided molecular clock assumption holds true.

# $3^{1/2}$ Methods

Parsimony

Distance method

Maximum Likelihood

Bayesian

# List of phylogenetic software

[http://evolution.gs.washington.edu/phylip/  
software.html](http://evolution.gs.washington.edu/phylip/software.html)

# Phylip

[http://evolution.genetics.washington.edu/phylip/  
getme.html](http://evolution.genetics.washington.edu/phylip/getme.html)

# Parsimony

Smallest number of evolutionary changes  
that explain the observed sequences.

Usually used for ancestral reconstruction  
using binary characters.

# Occam's razor






# William of Ockham

## 14th Century

<http://upload.wikimedia.org/>

# Ancestral reconstruction using Parsimony

DOLLO



# Main Parsimony programs in phylip

DNAPARS for DNA

PROTPARS for protein

# Parsimony




	1	2	3	4
Seq1	A	G	G	A
Seq2	A	G	G	G
Seq3	A	A	C	A
Seq4	A	A	C	G

To be informative at least one change is required  
Position 1: uninformative  
Positions 2-4: informative

# Parsimony

3 possible unrooted trees for position

2



Tree 1 is parsimonious tree with just one change

Best tree is the one that explains all the position with least number of changes.

# Distance method

Step 1: Calculate distance between all pairs  
of sequence in a multiple alignment

Step 2: Create a phylogenetic tree from this  
distance matrix

# Creating tree from distance matrix

FITCH: Fitch Margoliash method. No molecular clock.

KITSCH: Fitch Margoliash but under assumption of molecular clock.

NEIGHBOR: Neighbor joining or UPGMA.

NJ trees are unrooted and no assumption of molecular clock.

Align each pair of sequences and calculate distance as (number of mismatches/ number of matches) and create a distance matrix

	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>
<b>A</b>	-	$D_{AB} = 20$	$D_{AC} = 25$	$D_{AD} = 37$
<b>B</b>	-	-	$D_{BC} = 45$	$D_{BD} = 42$
<b>C</b>	-	-	-	$D_{CD} = 15$
<b>D</b>	-	-	-	-

# Programs to calculate distance matrix in PHYLIP

DNADIST for DNA. Uses various models for  
DNA

PROTDIST for protein. Uses various models  
including PAMs.

# Creating tree using PHYLIP

## Step 1

Create a multiple alignment

`muscle -in infile -phyout outfile`

# Creating tree using PHYLIP

## Step 2

Run a distance program

protdist

# Creating tree using PHYLIP

## Step 3

Run a distance program  
fitch

**TABLE 27.11.** Neighbor-joining example

	Cycle 1	Cycle 2	Cycle 3	Cycle 4	Cycle 5																																																																																										
Distance matrix	<table border="1"> <thead> <tr> <th></th><th>A</th><th>B</th><th>C</th><th>D</th><th>E</th></tr> </thead> <tbody> <tr> <td>B</td><td>5</td><td></td><td></td><td></td><td></td></tr> <tr> <td>C</td><td>4</td><td>7</td><td></td><td></td><td></td></tr> <tr> <td>D</td><td>7</td><td>10</td><td>7</td><td></td><td></td></tr> <tr> <td>E</td><td>6</td><td>9</td><td>6</td><td>5</td><td></td></tr> <tr> <td>F</td><td>8</td><td>11</td><td>8</td><td>9</td><td>8</td></tr> </tbody> </table>		A	B	C	D	E	B	5					C	4	7				D	7	10	7			E	6	9	6	5		F	8	11	8	9	8	<table border="1"> <thead> <tr> <th></th><th>U<sub>1</sub></th><th>C</th><th>D</th><th>E</th></tr> </thead> <tbody> <tr> <td>C</td><td>3</td><td></td><td></td><td></td></tr> <tr> <td>D</td><td>6</td><td>7</td><td></td><td></td></tr> <tr> <td>E</td><td>5</td><td>6</td><td>5</td><td></td></tr> <tr> <td>F</td><td>7</td><td>8</td><td>9</td><td>8</td></tr> </tbody> </table>		U <sub>1</sub>	C	D	E	C	3				D	6	7			E	5	6	5		F	7	8	9	8	<table border="1"> <thead> <tr> <th></th><th>U<sub>1</sub></th><th>C</th><th>U<sub>2</sub></th></tr> </thead> <tbody> <tr> <td>C</td><td>3</td><td></td><td></td></tr> <tr> <td>U<sub>2</sub></td><td>3</td><td>4</td><td></td></tr> <tr> <td>F</td><td>7</td><td>8</td><td>6</td></tr> </tbody> </table>		U <sub>1</sub>	C	U <sub>2</sub>	C	3			U <sub>2</sub>	3	4		F	7	8	6	<table border="1"> <thead> <tr> <th></th><th>U<sub>2</sub></th><th>U<sub>3</sub></th></tr> </thead> <tbody> <tr> <td>U<sub>3</sub></td><td>2</td><td></td></tr> <tr> <td>F</td><td>6</td><td>6</td></tr> </tbody> </table>		U <sub>2</sub>	U <sub>3</sub>	U <sub>3</sub>	2		F	6	6	<table border="1"> <thead> <tr> <th></th><th>U<sub>4</sub></th></tr> </thead> <tbody> <tr> <td>F</td><td>5</td></tr> </tbody> </table>		U <sub>4</sub>	F	5
	A	B	C	D	E																																																																																										
B	5																																																																																														
C	4	7																																																																																													
D	7	10	7																																																																																												
E	6	9	6	5																																																																																											
F	8	11	8	9	8																																																																																										
	U <sub>1</sub>	C	D	E																																																																																											
C	3																																																																																														
D	6	7																																																																																													
E	5	6	5																																																																																												
F	7	8	9	8																																																																																											
	U <sub>1</sub>	C	U <sub>2</sub>																																																																																												
C	3																																																																																														
U <sub>2</sub>	3	4																																																																																													
F	7	8	6																																																																																												
	U <sub>2</sub>	U <sub>3</sub>																																																																																													
U <sub>3</sub>	2																																																																																														
F	6	6																																																																																													
	U <sub>4</sub>																																																																																														
F	5																																																																																														
Step 1																																																																																															
S calculations	$S_A = (5+4+7+6+8)/4 = 7.5$	$S_{U_1} = (3+6+5+7)/3 = 7$	$S_{U_1} = (3+3+7)/2 = 6.5$	$S_{U_2} = (2+6)/1 = 8$	Because $N - 2 = 0$ , we cannot do this calculation.																																																																																										
$S_x = (\text{sum all } D_x)/(N - 2)$ , where $N$ is the # of OTUs in the set.	$S_B = (5+7+10+9+11)/4 = 10.5$ $S_C = (4+7+7+6+8)/4 = 8$ $S_D = (7+10+7+5+9)/4 = 9.5$ $S_E = (6+9+6+5+8)/4 = 8.5$ $S_F = (8+11+8+9+8)/4 = 11$	$S_C = (3+7+6=8)/3 = 8$ $S_D = (6+7+5+9)/3 = 9$ $S_E = (5+6+5+8)/3 = 8$ $S_F = (7+8+9+8)/3 = 10.6$	$S_C = (3+4+8)/2 = 7.5$ $S_{U_2} = (3+4+6)/2 = 6.5$ $S_F = (7+8+6)/2 = 10.5$	$S_{U_3} = (2+6)/1 = 8$ $S_F = (6+6)/1 = 12$																																																																																											
Step 2																																																																																															
Calculate pair with smallest ( $M$ ), where $M_{ij} = D_{ij} - S_i - S_j$ .	Smallest are $M_{AB} = 5 - 7.5 - 10.5 = -13$ $M_{DE} = 5 - 9.5 - 8.5 = -13$ Choose one of these (AB here).	Smallest is $M_{CU_1} = 3 - 7 - 8 = -12$ $M_{DE} = 5 - 9 - 8 = -12$ Choose one of these (DE here).	Smallest is $M_{CU_1} = 3 - 6.5 - 7.5 = -11$	Smallest is $M_{U_2F} = 6 - 8 - 12 = -14$ $M_{U_3F} = 6 - 8 - 12 = -14$ $M_{U_2U_3} = 2 - 8 - 8 = -14$ Choose one of these ( $M_{U_2U_3}$ here).																																																																																											
Step 3																																																																																															
Create a node (U) that joins pair with lowest $M_{ij}$ such that $S_{IU} = D_{ij}/2 + (S_i - S_j)/2$ .	U <sub>1</sub> joins A and B: $S_{AU_1} = D_{AB}/2 + (S_A - S_B)/2 = 1$ $S_{BU_1} = D_{AB}/2 + (S_B - S_A)/2 = 4$	U <sub>2</sub> joins D and E: $S_{DU_2} = D_{DE}/2 + (S_D - S_E)/2 = 3$ $S_{EU_2} = D_{DE}/2 + (S_E - S_D)/2 = 2$	U <sub>3</sub> joins C and U <sub>1</sub> : $S_{CU_3} = D_{CU_1}/2 + (S_C - S_{U_1})/2 = 2$ $S_{U_1U_3} = D_{CU_1}/2 + (S_{U_1} - S_C)/2 = 1$	U <sub>4</sub> joins U <sub>2</sub> and U <sub>3</sub> : $S_{U_2U_4} = D_{U_2U_3}/2 + (S_{U_2} - S_{U_3})/2 = 1$ $S_{U_3U_4} = D_{U_2U_3}/2 + (S_{U_3} - S_{U_2})/2 = 1$	For last pair, connect U <sub>4</sub> and F with branch length = 5.																																																																																										
Step 4																																																																																															
Join $i$ and $j$ according to $S$ above and make all other taxa in form of a star. Branches in black are of unknown length. Branches in red are of known length.																																																																																															
Step 5																																																																																															
Calculate new distance matrix of all other taxa to U with $D_{xU} = D_{ix} + D_{jx} - D_{ij}$ , where $i$ and $j$ are those selected from above.					Comments Note this is the same tree we started with (drawn in unrooted form here).																																																																																										

# Output tree format

## Newick

```
(P73_HUMAN/:0.16068,((P53_XENLA/:0.18610,((P53_ONCMY/:0.12081,  
P53_DANRE/:0.12111):0.02394,P53_HUMAN/:0.22849):0.03528):0.04183,  
P53_ORYLA/:0.20291):0.11899,Q27937_LOL:0.48924);
```

# Reliability Bootstrapping

Randomly sample the original alignment

Create many alignments

Create many trees

Create a consensus tree

# Bootstrapping in phylip

seqboot

protdist

fitch

consense

Don't forget to use "multiple" parameters

# Maximum Likelihood

# Conditional probability

Likelihood

$$Prob(H|D) = \frac{Prob(D|H)Prob(H)}{Prob(D)}$$

$H$  = Hypothesis

D= data

# Calculating likelihood

Given a dataset

$$D = D_1, D_2, \dots, D_n$$

Likelihood

$$L = \text{Prob}(D_1|H) \text{Prob}(D_2|H) \dots \text{Prob}(D_n|H)$$

# Maximum likelihood

Given the likelihood

$$L = \text{Prob}(D_1|H)\text{Prob}(D_2|H) \dots \text{Prob}(D_n|H)$$

We calculate the likelihood for a set of probabilities of H

The probability of H is “most probably” where the likelihood is maximum.

# Let's calculate the probability of heads

HHTTH


$$\begin{aligned}L &= pp(1-p)(1-p)p \\&= p^3(1-p)^2\end{aligned}$$

$$\ln L = \ln p^3 + \ln(1-p)^2$$

$$\frac{d(\ln L)}{dp} = \frac{3}{p} - \frac{2}{(1-p)} = 0$$

$$p = \frac{3}{5}$$

# Probability of a tree



For nucleotide sequence  
 $x = (A, T, G, \text{ or } C)$

$$L = \sum_x Prob(x) Prob(A|x) Prob(B|x)$$

# RAX-ML

[http://icwww.epfl.ch/~stamatak/index-Dateien/  
Page443.htm](http://icwww.epfl.ch/~stamatak/index-Dateien/Page443.htm)

# Text books

