

Logistics

- Problem set 3 available later today
- 7Eleven-1 due tomorrow
- Midterm Exam, October 10th
 - Covers sequence alignment, models of sequence substitution
 - [See Study Guide on Syllabus page](#)
 - Lectures through Sept. 19 (and review on Sept 21)
 - Closed book
 - Two pages of notes

Two widely used families of Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

- PAM matrices, Dayhoff *et al*, 1978
- BLOSUM (Block Sum) matrices, Hennikoff & Hennikoff, 1991

Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

Overall strategy for both PAM and BLOSUM

1. Trusted amino acid alignments
2. Obtain amino acid pair counts (A_{xy}^N) with corrections for
 - Evolutionary divergence
 - Sample biases
3. Estimate substitution frequencies, q_{xy}^N , from pair counts, A_{xy}^N
4. Log odds substitution matrix: $S^N[x,y] = c \log \frac{q_{xy}^N}{p_x p_y}$

Log odds substitution matrices

Two sequences have N PAMs divergence, if, on average, N amino acid replacements per 100 residues occurred since their separation

$$S^N[x,y] = c \log \frac{q_{xy}^N}{p_x p_y}$$

Scaling constant

Frequency of x aligned with y in sequences with divergence N

Frequency of x aligned with y in "random" sequences

Two widely used families of Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

- PAM matrices, Dayhoff *et al*, 1978
- BLOSUM (Block Sum) matrices, Hennikoff & Hennikoff, 1991

PAM Matrices

Atlas of Protein Sequence & Structure
1965 - 1978



Examined 1572 changes in 71 groups
of closely related proteins



Margaret Dayhoff
PhD in Chemistry, 47
Watson Computing Lab
Fellow 47 - 48

1. Trusted multiple sequence alignments

Examined 1572 changes in 71 groups
of closely related proteins

At least 85% identical

```
DVLVAHQHILRAFQRHWGHTALNPSILLEAGGVOTLSVEHHLLDEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIADNPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIMDQPTFLLEAGGVOTLSVEHKSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIDQPTFLLEAGGVOTLSVEHKSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
DVLVAHQHILRAFQRHWGHTIENPSILLEAGGVOTLSVEHHSLEPAIILGQAFV
```

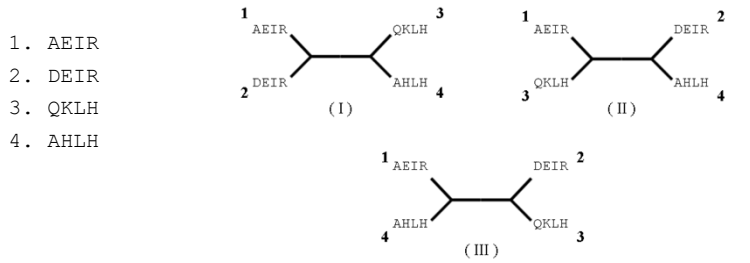
2. Obtain amino acid pair counts (A_{xy}) with corrections for evolutionary divergence and sample biases

Counting amino acid pairs on a tree:

For each unrooted tree with k leaves

- Select the tree(s) that require the fewest substitutions to explain the data
- Count amino acid pairs on the branches of the tree

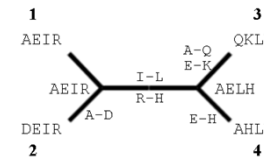
Suppose we have an alignment of four sequences. There are 3 hypotheses (i.e., 3 unrooted trees) for their evolutionary relationships



How to select the tree(s) that require the fewest substitutions to explain the data...

For a given a tree, assign labels to internal nodes that minimize the number of changes required to explain the data

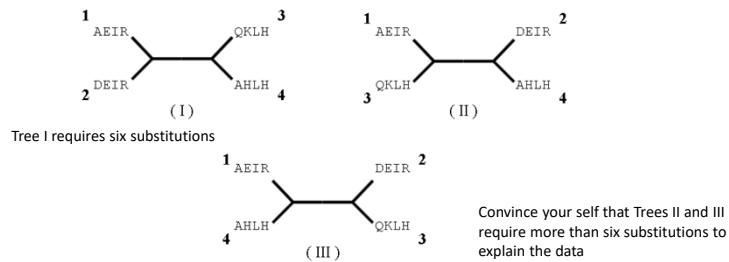
- 1. AEIR
- 2. DEIR
- 3. QKLH
- 4. AHLH



Tree I requires six substitutions

There may be more than one set of labels that satisfies this criterion

Select the most parsimonious tree; i.e., the tree that requires the fewest substitutions to explain the data.



2. Obtain amino acid pair counts (A_{xy}) with corrections for evolutionary divergence and sample biases

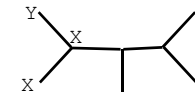
Counting amino acid pairs on a tree:

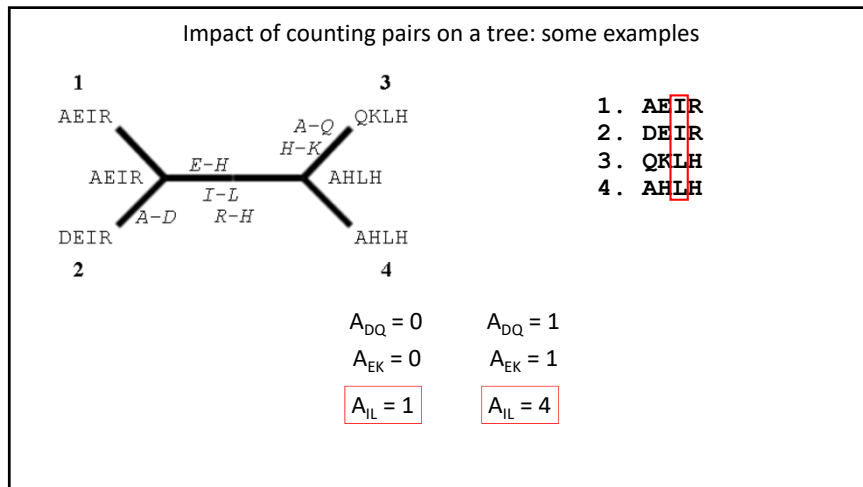
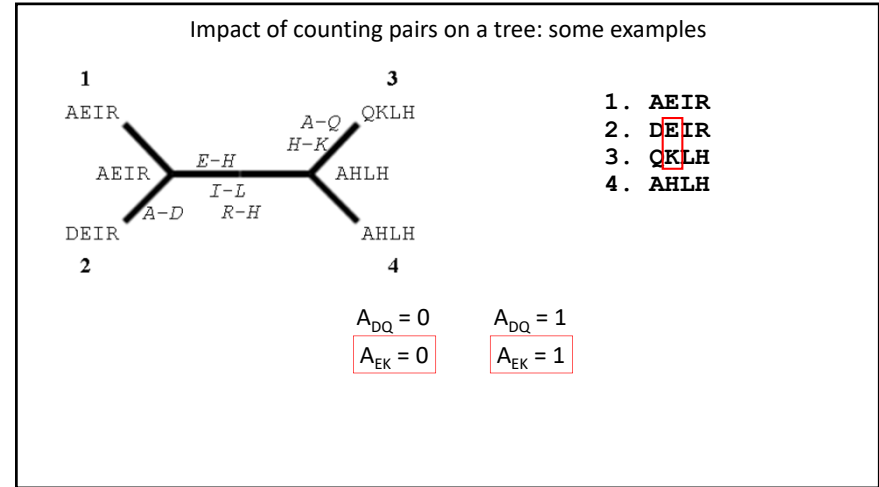
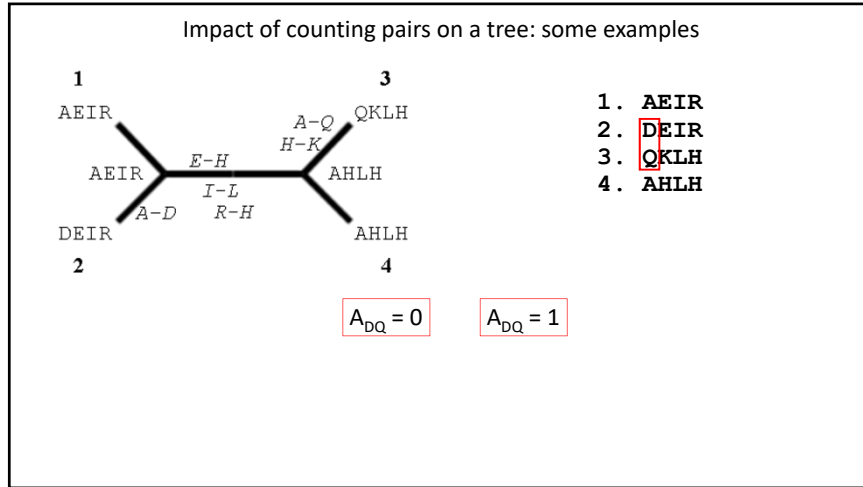
For each unrooted tree with k leaves

- Select the tree(s) that require the fewest substitutions to explain the data
- Count amino acid pairs on the branches of the tree

• For each branch,

- if labeled x — y , $A_{xy}^N = A_{xy}^N + 1$ and $A_{yx}^N = A_{yx}^N + 1$
- if labeled x — x , $A_{xx}^N = A_{xx}^N + 2$





3. Estimate substitution frequencies, q_{xy}^N , from pair counts, A_{xy}

- Markov model with 20 states (A, C, D, E ... Y)
- Estimate 1 PAM transition matrix P^1 from A_{xy}
- N-PAM transition matrix: $P^1 = (P^1)^N$
- $q_{xy}^N = p_x P_{xy}^N$
- $S^N[x,y] = c \log \frac{q_{xy}^N}{p_x p_y}$

Is P_{xy}^N a symmetric matrix? No. (Check this algebraically).
Is $S^N[x,y]$ a symmetric matrix? Yes (Check this algebraically).

Two widely used families of Amino Acid Substitution Matrices Parameterized for evolutionary divergence (*N*)

- PAM matrices, Dayhoff *et al*, 1978
- BLOSUM (Block Sum) matrices, Hennikoff & Hennikoff, 1991

BLOSUM Matrices

- Trusted data
 - 2000 blocks of conserved regions in ~500 groups of proteins
- Count amino acid pairs: A_{xy}^N
 - Parameterize by evolutionary distance, *N*
 - Correct for sample bias
- Calculate amino acid frequencies:
 - Related pairs: q_{xy}^N
 - Background pair frequencies calculated from blocks: E_{xy}
- Log likelihood scoring matrix
 - $S^N = 2 \log_2 \frac{q_{xy}^N}{E_{xy}}$

1. Trusted multiple sequence alignments

BLOCKS from MOTIF
9 sequences are included in 4 blocks

```

    DbBari.1 ( 21) 218  YFFPYPNQVPCV
    DMB1 ( 24) 216  DGGIFVNGCF
    DMI1 ( 22) 224  FFAITVNGCQ
    TC1A_CARB ( 23) 222  DGGIFVNGCF
    TC1A_CARB ( 23) 222  DGGIFVNGCF
    TC1A_CARB ( 23) 222  DGGIFVNGCF
    TC1A_CARB ( 20) 216  DGGIFVNGCF
    Bcfw1 ( 22) 214  DGGIFVNGCF
    DfB1 ( 22) 214  DGGIFVNGCF
    DfB1 ( 22) 214  DGGIFVNGCF

    TC1 B, width = 44
    DbBari.1 ( 21) 218  YFFPYPNQVPCV
    DMB1 ( 24) 216  DGGIFVNGCF
    DMI1 ( 22) 224  FFAITVNGCQ
    TC1A_CARB ( 23) 222  DGGIFVNGCF
    TC1A_CARB ( 23) 222  DGGIFVNGCF
    TC1A_CARB ( 23) 222  DGGIFVNGCF
    TC1A_CARB ( 20) 216  DGGIFVNGCF
    Bcfw1 ( 22) 214  DGGIFVNGCF
    DfB1 ( 22) 214  DGGIFVNGCF
    DfB1 ( 22) 214  DGGIFVNGCF

    TC1 C, width = 11
    DbBari.1 ( 21) 218  YFFPYPNQVPCV
    DMB1 ( 24) 216  DGGIFVNGCF
    DMI1 ( 22) 224  FFAITVNGCQ
    TC1A_CARB ( 23) 222  DGGIFVNGCF
    TC1A_CARB ( 23) 222  DGGIFVNGCF
    TC1A_CARB ( 23) 222  DGGIFVNGCF
    TC1A_CARB ( 20) 216  DGGIFVNGCF
    Bcfw1 ( 22) 214  DGGIFVNGCF
    DfB1 ( 22) 214  DGGIFVNGCF
    DfB1 ( 22) 214  DGGIFVNGCF

    TC1 D, width = 15
    DbBari.1 ( 18) 217  MFIQKQKVFYV
    DMB1 ( 18) 247  WQAPFSELPFVSEV
    DMI1 ( 18) 204  MFIQKQKVFYV
    TC1A_CARB ( 18) 201  MFIQKQKVFYV
    TC1A_CARB ( 18) 201  MFIQKQKVFYV
    TC1A_CARB ( 18) 201  MFIQKQKVFYV
    TC1A_CARB ( 18) 201  MFIQKQKVFYV
    Bcfw1 ( 22) 297  MFIQKQKVFYV
    DfB1 ( 18) 229  MFIQKQKVFYV
  
```

~2000 blocks representing
500+ groups of proteins

Automated construction and graphical presentation of protein blocks from unaligned sequences
Steven Henikoff^{1,2*}, Jojo G. Henikoff^{1,2}, William J. Alford^{1,2}, Shmel Pietrowski^{1,2}
¹Academy of Natural Sciences, Penn. State Univ., University Park, PA 16802, USA
²Department of Biology, Penn. State Univ., University Park, PA 16802, USA
Received 30 May 1995, revised 23 Jun 1995, accepted 1 July 1995, published 1 August 1995

2. Count amino acid pairs: A_{xy}^N

Parameterize by evolutionary distance, *N*
Correct for sample bias

- Cluster sequences such that if *s1* and *s2* are in different clusters, then *identity(s1, s2) < N%*
- Count amino acid pairs in *s1* aligned with *s2* only if *s1* and *s2* are in different clusters
- Normalize for cluster size

As an example...

BLOSUM clustering example

- 1: KKRK
- 2: KKKK
- 3: KNRN
- 4: NRNR
- 5: KNKN
- 6: KRNR

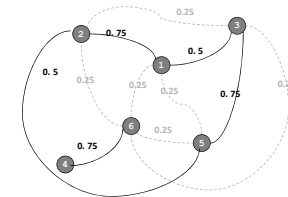
Percent Sequence Identity						
	[2]	[3]	[4]	[5]	[6]	
[1]	0.75	0.5	0	0.25	0.25	
[2]		0.25	0	0.5	0.25	
[3]			0	0.75	0.25	
[4]				0	0.75	
[5]					0.25	

Unclassified sequences: Every sequence is at least 25% identical

Percent Sequence Identity						
	[2]	[3]	[4]	[5]	[6]	
[1]	0.75	0.5	0	0.25	0.25	
[2]		0.25	0	0.5	0.25	
[3]			0	0.75	0.25	
[4]				0	0.75	
[5]					0.25	

- 1: KKRK
- 2: KKKK
- 3: KNRN
- 5: KNKN
- 4: NRNR
- 6: KRNR

< 45% identical



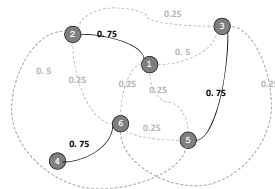
- 1: KKRK
- 2: KKKK
- 3: KNRN
- 5: KNKN
- 4: NRNR
- 6: KRNR

< 65% id

< 65% id

< 65% id

Percent Sequence Identity						
	[2]	[3]	[4]	[5]	[6]	
[1]	0.75	0.5	0	0.25	0.25	
[2]		0.25	0	0.5	0.25	
[3]			0	0.75	0.25	
[4]				0	0.75	
[5]					0.25	



Two widely used families of Amino Acid Substitution Matrices Parameterized for evolutionary divergence (*N*)

- PAM matrices, Dayhoff *et al*, 1978
- BLOSUM (Block Sum) matrices, Hennikoff & Hennikoff, 1991

Similarities and differences between PAM and BLOSUM

	PAM	BLOSUM
Evolutionary model	Explicit evolutionary model	None
Data	Full length MSAs of closely related sequences.	Conserved blocks. i.e., ungapped local MSAs
Bias correction	Trees	Clustering
Multiple substitutions	Markov model: $P^n = (P^1)^n$	Implicitly represented in data (clustering)
Evolutionary distance	Markov model: $P^n = (P^1)^n$	Clustering
Matrices	Transition and log odds scoring matrices	Log odds scoring matrix only.
Parameter n	Distance increases with n	Distance decreases with n
Biophysical properties	Derived indirectly from data	Derived indirectly from data

