

Spring 2025

BIOS 477/877

Bioinformatics and Molecular Evolution

Lecture 9

BIOS477/877 L9 - 1

1

TODAY'S TOPICS

- Amino Acid Substitution Matrix
 - Dayhoff's PAM Matrix
 - BLOSUM Matrix

BIOS477/877 L9 - 2

2

Substitution matrices based on empirical data

- **PAM matrices**
 - Dayhoff, Schwartz, and Orcutt (1978)
- **BLOSUM matrices**
 - Henikoff and Henikoff (1992)

Also see Eddy (2004) Nature Biotechnology 22: 1035-36

BIOS477/877 L9 - 3

3

PAM matrices (Dayhoff *et al.* 1978)

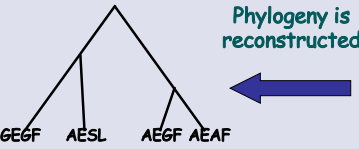
- **Accepted point mutations** (point accepted mutations, percent accepted mutations)
 - **accepted by selection**: no (or very weak) deleterious effect, maintaining the function
- Based on 1,572 changes in 71 groups of **closely related proteins** (34 protein families)
 - at least 85% identical
 - no ambiguity in alignments, no gap
 - most likely observed substitutions do not affect protein functions (accepted by selection, close to neutral)
 - successive (multiple) substitutions at one site are minimal (no hidden substitution)

BIOS477/877 L9 - 4

4

PAM matrices

- Numbers of **accepted point mutations**: $f(a,b)$ are counted based on phylogenies
- Assumption: substitutions are equally likely in each direction (*e.g.*, $G \rightarrow A = A \rightarrow G$)



Phylogeny is reconstructed

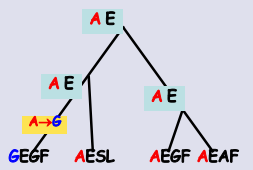
GEGF
AESL
AEGF
AEAf

BIOS477/877 L9 - 5

5

PAM matrices

- Numbers of **accepted point mutations**: $f(a,b)$ are counted based on phylogenies
- Assumption: substitutions are equally likely in each direction (*e.g.*, $G \rightarrow A = A \rightarrow G$)



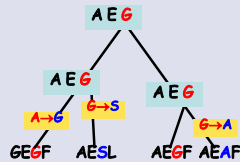
Using the maximum parsimony principle, ancestral sequences can be inferred

BIOS477/877 L9 - 6

6

PAM matrices

- Numbers of **accepted point mutations**: $f(a,b)$ are counted based on phylogenies
 - Assumption: substitutions are equally likely in each direction (e.g., $G \rightarrow A = A \rightarrow G$)



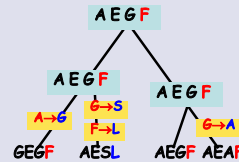
Using the maximum parsimony principle, ancestral sequences can be inferred

BIOS477/877 L9 - 7

7

PAM matrices

- Numbers of **accepted point mutations**: $f(a,b)$ are counted based on phylogenies
 - Assumption: substitutions are equally likely in each direction (e.g., $G \rightarrow A = A \rightarrow G$)



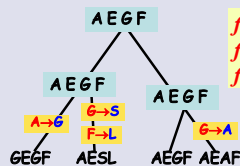
Substitutions can be identified along the phylogeny

BIOS477/877 L9 - 8

8

PAM matrices

- Numbers of **accepted point mutations**: $f(a,b)$ are counted based on phylogenies
 - Assumption: substitutions are equally likely in each direction (e.g., $G \rightarrow A = A \rightarrow G$)


$$\begin{aligned} f(G,A) &= \text{Freq}(G \rightarrow A) + \text{Freq}(A \rightarrow G) = 2 \\ f(G,S) &= 1 \\ f(F,L) &= 1 \end{aligned}$$

BIOS477/877 L9 - 9

9

PAM matrices

Numbers of accepted point mutations: $f(a,b)$

Dayhoff *et al.* (1978)

Based on 1,572 changes, but still missing 35 types of substitutions

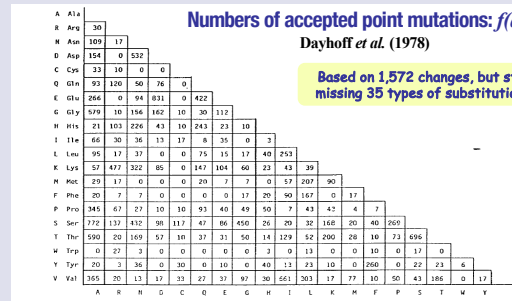


Figure 80. Numbers of accepted point mutations ($\times 10$) accumulated from closely related sequences. Fifteen hundred and seventy-

BIOS477/877 L9 - 10

10

PAM matrices

- **Relative mutability:** $m(a)$

Probability that the amino acid a will change in a given small evolutionary interval
[from a pair]

Amino acid: a	A	E	F
-----------------	---	---	---

GEGF
 AEGG

$$m(a) = \frac{\sum_i f(i,a)}{f(a)}$$

[combined from multiple trees]

Substitutions are collected from trees with different lengths

$$m(a) = \frac{\text{Number of times amino acid } a \text{ is substituted by any other amino acid}}{\sum_{\text{branch}} \{ \text{Freq. of amino acid } a \} \times (\text{Number of total substitutions}) \times 100}$$

(Number of occurrence of amino acid a) / (Total number of residues)

This denominator is called "the total exposure of the amino acid to mutation"

BIOS477/877 L9 - 11

11

PAM matrices

Relative mutability: $m(a)$

Dayhoff *et al.* (1978)

Table 21
Relative Mutabilities of the Amino Acids^a

More mutable		Less mutable	
Asn	134	His	66
Ser	120	Arg	65
Asp	106	Lys	56
Glu	102	Pro	56
Ala	100	Gly	49
Thr	97	Tyr	41
Ile	96	Phe	41
Met	94	Leu	40
Gln	93	Cys	20
Val	74	Trp	18

^aThe value for Ala has been arbitrarily set at 100.

BIOS477/877 L9 - 12

12

PAM matrices

• Mutation probability

$$M(a,b) = \lambda m(b) \times f(a,b) / \sum_a f(a,b), \text{ where } a \neq b$$

$m(b)$: relative mutability of amino acid b

$f(a,b)$: frequency of accepted point mutations between amino acids a and b

$\sum_a f(a,b)$: number of times the amino acid b is substituted by any other amino acid

λ : proportionality constant (normalization factor)

→ The probability of the amino acid b being replaced by the amino acid a after a given evolutionary time

$$M(b,b) = 1 - \lambda m(b)$$

unchange probability (the diagonal elements)

BIOS477/877 L9 - 13

13

PAM matrices

Mutation probability matrix: $M(a,b)$ Dayhoff et al. (1978)

EVOLUTIVE AMINO ACID	ORIGINAL AMINO ACID																			
	N	A	R	D	C	Q	E	G	H	I	L	K	M	F	S	T	W	Y	V	
A	Ala	1000	7	9	10	2	8	17	23	2	6	4	2	2	22	35	32	6	18	
R	Arg	1	9913	1	0	1	10	0	0	10	3	1	19	4	1	4	6	1	8	0
D	Asp	4	1	9822	36	0	4	6	0	21	3	1	13	0	1	7	20	9	1	4
E	Asp	6	0	42	9859	0	4	53	6	4	1	0	3	0	0	1	5	3	0	1
C	Cys	1	0	0	9973	6	0	0	0	1	1	0	0	0	1	5	1	0	3	2
Q	Gln	3	9	4	5	0	9874	27	1	23	1	3	6	4	0	6	2	2	0	1
G	Glu	10	0	7	56	0	35	9865	4	2	3	1	4	1	0	3	4	2	0	1
H	Gly	21	1	12	11	1	7	9555	1	0	1	1	0	1	1	3	21	3	0	5
I	Ile	1	8	18	3	1	20	1	0	9912	0	1	1	0	2	3	1	1	4	1
L	Ile	2	2	3	1	2	1	2	0	0	9872	9	2	12	7	0	1	7	0	1
K	Lys	3	1	3	0	0	6	1	1	4	22	9847	2	45	13	3	1	5	4	2
M	Met	2	37	25	6	0	12	7	2	2	4	1	9555	20	0	3	8	11	0	1
F	Phe	1	1	0	0	0	2	0	0	0	0	5	8	4	9874	1	0	1	2	0
P	Phe	1	1	1	0	0	0	0	1	2	8	6	0	4	9946	0	2	1	7	28
S	Ser	13	5	2	1	1	8	1	2	5	1	2	2	1	1	9926	12	4	0	2
T	Thr	20	13	34	7	11	4	5	16	2	2	3	7	4	3	17	9840	38	1	2
W	Tyr	22	2	12	4	1	3	2	2	1	11	2	8	6	1	5	32	9871	0	2
Y	Tyr	0	2	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	9976	1
V	Tyr	1	0	3	0	0	1	0	4	1	1	0	0	21	0	1	1	0	9861	1
V	Val	13	2	1	1	3	2	2	3	3	5	11	1	3	7	10	0	2	9901	

Figure 82. Mutation probability matrix for the evolutionary distance of 1 PAM. An element of this matrix, M_{ij} , gives the probability that the amino acid in column i will be replaced by the amino acid in row j after a given evolutionary interval, in this case 1 accepted point mutation per 100 amino acids. Thus, there is a 0.66% probability that Asp will be replaced by Glu. To simplify the appearance, the elements are shown multiplied by 10,000.

Probabilities of AA_j replaced by AA_i

BIOS477/877 L9 - 14

14

PAM matrices

Mutation probability matrix: $M(a,b)$ Dayhoff et al. (1978)

EVOLUTIVE AMINO ACID	ORIGINAL AMINO ACID																			
	N	A	R	D	C	Q	E	G	H	I	L	K	M	F	S	T	W	Y	V	
A	Ala	1000	7	9	10	2	8	17	23	2	6	4	2	2	22	35	32	6	18	
R	Arg	1	9913	1	0	1	10	0	0	10	3	1	19	4	1	4	6	1	8	0
D	Asp	4	1	9822	36	0	4	6	0	21	3	1	13	0	1	7	20	9	1	4
E	Asp	6	0	42	9859	0	4	53	6	4	1	0	3	0	0	1	5	3	0	1
C	Cys	1	0	0	9973	6	0	0	0	1	1	0	0	0	1	5	1	0	3	2
Q	Gln	3	9	4	5	0	9874	27	1	23	1	3	6	4	0	6	2	2	0	1
G	Glu	10	0	7	56	0	35	9865	4	2	3	1	4	1	0	3	4	2	0	1
H	Gly	21	1	12	11	1	7	9555	1	0	1	1	0	1	1	3	21	3	0	5
I	Ile	1	8	18	3	1	20	1	0	9912	0	1	1	0	2	3	1	1	4	1
L	Ile	2	2	3	1	2	1	2	0	0	9872	9	2	12	7	0	1	7	0	1
K	Lys	3	1	3	0	0	6	1	1	4	22	9847	2	45	13	3	1	5	4	2
M	Met	2	37	25	6	0	12	7	2	2	4	1	9555	20	0	3	8	11	0	1
F	Phe	1	1	1	0	0	2	0	0	0	0	5	8	4	9874	1	0	1	2	0
P	Phe	1	1	1	0	0	0	0	1	2	8	6	0	4	9946	0	2	1	7	28
S	Ser	13	5	2	1	1	8	1	2	5	1	2	2	1	1	9926	12	4	0	2
T	Thr	20	13	34	7	11	4	5	16	2	2	3	7	4	3	17	9840	38	1	2
W	Tyr	22	2	12	4	1	3	2	2	1	11	2	8	6	1	5	32	9871	0	2
Y	Tyr	0	2	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	9976	1
V	Tyr	1	0	3	0	0	1	0	4	1	1	0	0	21	0	1	1	0	9861	1
V	Val	13	2	1	1	3	2	2	3	3	5	11	1	3	7	10	0	2	9901	

Figure 82. Mutation probability matrix for the evolutionary distance of 1 PAM. An element of this matrix, M_{ij} , gives the probability that the amino acid in column i will be replaced by the amino acid in row j after a given evolutionary interval, in this case 1 accepted point mutation per 100 amino acids. Thus, there is a 0.66% probability that Asp will be replaced by Glu. To simplify the appearance, the elements are shown multiplied by 10,000.

Probabilities of AA_j replaced by AA_i

BIOS477/877 L9 - 15

15

PAM matrices

Mutation probability matrix: $M(a,b)$ Dayhoff et al. (1978)

EVOLUTIVE AMINO ACID	ORIGINAL AMINO ACID																			
	N	A	R	D	C	Q	E	G	H	I	L	K	M	F	S	T	W	Y	V	
A	Ala	1000	7	9	10	2	8	17	23	2	6	4	2	2	22	35	32	6	18	
R	Arg	1	9913	1	0	1	10	0	0	10	3	1	19	4	1	4	6	1	8	0
D	Asp	4	1	9822	36	0	4	6	0	21	3	1	13	0	1	7	20	9	1	4
E	Asp	6	0	42	9859	0	4	53	6	4	1	0	3	0	0	1	5	3	0	1
C	Cys	1	0	0	9973	6	0	0	0	1	1	0	0	0	1	5	1	0	3	2
Q	Gln	3	9	4	5	0	9874	27	1	23	1	3	6	4	0	6	2	2	0	1
G	Glu	10	0	7	56	0	35	9865	4	2	3	1	4	1	0	3	4	2	0	1
H	Gly	21	1	12	11	1	7	9555	1	0	1	1	0	1	1	3	21	3	0	5
I	Ile	1	8	18	3	1	20	1	0	9912	0	1	1	0	2	3	1	1	4	1
L	Ile	2	2	3	1	2	1	2	0	0	9872	9	2	12	7	0	1	7	0	1
K	Lys	3	1	3	0	0	6	1	1	4	22	9847	2	45	13	3	1	5	4	2
M	Met	2	37	25	6	0	12	7	2	2	4	1	9555	20	0	3	8	11	0	1
F	Phe	1	1	1	0	0	2	0	0	0	0	5	8	4	9874	1	0	1	2	0
P	Phe	1	1	1	0	0	0	0	1	2	8	6	0	4	9946	0	2	1	7	28
S	Ser	13	5	2	1	1	8	1	2	5	1	2	2	1	1	9926	12	4	0	2
T	Thr	20	13	34	7	11	4	5	16	2	2	3	7	4	3	17	9840	38	1	2
W	Tyr	22	2	12	4	1	3	2	2	1	11	2	8	6	1	5	32	9871	0	2
Y	Tyr	0	2	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	9976	1
V	Tyr	1	0	3	0	0	1	0	4	1	1	0	0	21	0	1	1	0	9861	1
V	Val	13	2	1	1	3	2	2	3	3	5	11	1	3	7	10	0	2	9901	

Figure 82. Mutation probability matrix for the evolutionary distance of 1 PAM. An element of this matrix, M_{ij} , gives the probability that the amino acid in column i will be replaced by the amino acid in row j after a given evolutionary interval, in this case 1 accepted point mutation per 100 amino acids. Thus, there is a 0.66% probability that Asp will be replaced by Glu. To simplify the appearance, the elements are shown multiplied by 10,000.

Mutation probability matrix is not symmetrical

Probabilities of AA_j replaced by AA_i

BIOS477/877 L9 - 16

16

PAM matrices

• Relatedness odds score

$$R(a,b) = M(a,b) / f(a)$$

$$M(a,b) = \lambda m(b) \times f(a,b) / \sum_a f(a,b)$$

the probability that amino acid b will change to a in a related sequence in a given interval

$f(a)$: the chance of a random occurrence of amino acid a

[frequency of occurrence of amino acid a]

= (number of occurrences of the residue a) / (total number of residues)

$$R(a,b) = M(a,b) / f(a) = \{\lambda m(b) \times f(a,b) / \sum_a f(a,b)\} / f(a)$$

$$= \{\lambda \sum_a f(a,b) \times f(a,b)\} / \{\sum_a f(a,b) \times f(a)\} = \lambda f(a,b) / \{f(a) \times f(b)\}$$

$$R(b,a) = M(b,a) / f(b) = \{\lambda m(a) \times f(b,a) / \sum_b f(b,a)\} / f(b)$$

$$= \lambda f(b,a) / \{f(b) \times f(a)\} = R(a,b)$$

$$R(a,b) = R(b,a) \rightarrow \text{Relatedness odds score matrix is symmetrical}$$

BIOS477/877 L9 - 17

17

PAM matrices

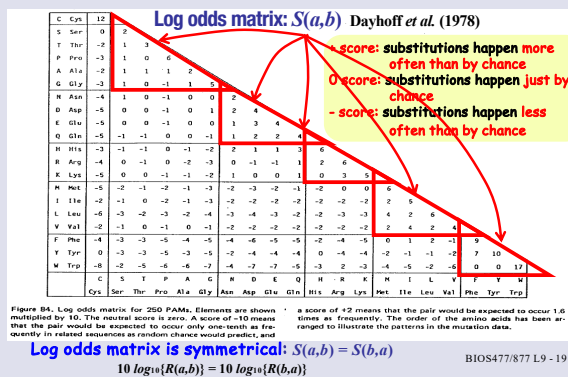
• Relatedness odds score

$$R(a,b) = M(a,b) / f(a)$$

$$M(a,b) = \lambda m(b) \times f(a,b) / \sum_a f(a,b)$$

the probability that amino acid b will change to a in a related

PAM matrices



19

PAM matrices

PAM1 matrix

→ normalized (using λ) to represent an amount of evolution producing an average of one mutation per hundred amino acids [Evolutionary interval of PAM1]

$100 \times \sum_b \{f(b)M(b,b)\} = 99$ where $M(b,b) = 1 - \lambda m(b)$
within 100 amino acids 99 are unchanged (or only 1 changed)

M_1 : PAM1 mutation probability matrix

→ shows the probability of AA_i replaced by AA_j after the evolutionary interval of PAM1 (when one mutation per 100 aa is found)

e.g., M_{250} : Probability matrix after evolutionary interval of PAM250 (after 250 changes are produced in 100 aa)

BIOS477/877 L9 - 20

20

PAM matrices

PAM1 matrix

→ normalized (using λ) to represent an amount of evolution producing an average of one mutation per hundred amino acids [Evolutionary interval of PAM1]

$100 \times \sum_b \{f(b)M(b,b)\} = 99$ where $M(b,b) = 1 - \lambda m(b)$
within 100 amino acids 99 are unchanged (or only 1 changed)

M_1 : PAM1 mutation probability matrix

→ shows the probability of AA_i replaced by AA_j after the evolutionary interval of PAM1 (when one mutation per 100 aa is found)

M_n : mutation probability matrix for PAM n

$M_n = (M_1)^n$ (e.g., PAM250 or $M_{250} = M_1^{250}$)

BIOS477/877 L9 - 21

21

PAM matrices

Table 23
Correspondence between Observed Differences and the Evolutionary Distance

Observed Percent Difference	Evolutionary Distance in PAMs
1	1
5	5
10	11
15	17
20	23
25	30
30	38
35	47
40	56
60	112
65	133
70	159
75	190
80	246
85	328

← PAM1

$100[1 - S_b \{f(b)M_n(b,b)\}]$

PAM250 ≈ 20% similarity

~PAM250

BIOS477/877 L9 - 22

22

PAM matrices updated

JTT matrices

by Jones, Taylor, and Thornton (1992)
→ 59,190 accepted point mutations for 16,300 proteins

Gonnet matrices

by Gonnet, Cohen, Benner (1992)
→ Based on exhaustive pairwise alignment from the protein database (~8,344,353 amino acids).

BIOS477/877 L9 - 23

23

BLOSUM matrices (Henikoff and Henikoff 1992)

Blocks substitution matrix

→ Based on ~2,000 conserved amino acid patterns (or ungapped blocks), representing more than 500 families.

→ Based on local, multiple alignment of all commonly-occurring motifs (blocks) in the protein sequence database.

The Blocks Database

(no longer available, but used to generate BLOSUM matrices)

BIOS477/877 L9 - 24

24

BLOCK entry example

Block PR00237A

```
ID GPCRHHODPSN; BLOCK
AC PR00237A; distance from previous block=(5,490)
DE Rhodopsin-like GPCR superfamily signature
BL adapted; width=25; seqs=739; 99.5%=1613; strength=1138
OAR1_LOCHI|Q25121 ( 53) VTAVSLSLIITIVGNVLVLSVF 3
OAR2_LOCHI|Q25122 ( 53) VTAVSLSLIITIVGNVLVLSVF 3
O61710 ( 22) ISLAVLEPINVLVIGGNCLVIAAVF 22
OAR_DROME|P22270 (111) LTALVLSVIVLTIGNIVLSVF 3
DOP2_DROME|Q24563 (110) GLLAFLFLFSFATVFGNSLVILAVI 5
OAR_HELVI|Q25188 ( 54) CTAIVLTLLIISTVGNILVLSVF 6
Q23128 ( 35) ISLLALAFNLHMVAGNLLVMMAVF 9
Q23126 ( 35) ISLLALAFNLHMVAGNLLVMMAVF 9
ALAB_MESAU|P18841 ( 47) SVGLVLGAFILFAIVGNILVLSVA 5
ALAB_RAT|P15823 ( 47) SVGLVLGAFILFAIVGNILVLSVA 5
ALAB_HUMAN|P35368 ( 47) SVGLVLGAFILFAIVGNILVLSVA 5
Q23127 ( 34) AATALLLAILVTIVGNISLVISVF 3
ALAA_ORYLA|Q91175 ( 28) VLGMVLGIFILFGVIGNILVLSVV 5
D2D1_XENLA|P24628 ( 30) YYAMLLTLLVVFVIFGNVLVCIASV 6
Q2DR_CERAR|P52702 ( 36) YYATLLTLLAVIVFGNVLVCHAVS 5
Q73810 ( 33) YVAVLLTLLFVIFGNVLVCHAVS 5
OAR_BOHMO|Q17212 ( 56) CTAIILTHIISTVVGNNLVLSVF 7
SR5_MOUSE|Q08858 ( 39) LVFVLYLLVCTVGLGGNTLVIVVL 4
```

Blocks are multiply aligned **ungapped** segments corresponding to the most highly conserved regions of proteins

BIOS477/877 L9 - 25

25

BLOSUM matrices

Seq1 MCL
Seq2 CCV
Seq3 ICV
Seq4 MAI
Seq5 TCL

Observed amino acid pairs:
 $CC + CC + CA + CC$

BIOS477/877 L9 - 26

26

BLOSUM matrices

Seq1 MCL
Seq2 CCV
Seq3 ICV
Seq4 MAI
Seq5 TCL

Observed amino acid pairs:
 $CC + CC + CA + CC$
 $+ CC + CA + CC$

BIOS477/877 L9 - 27

27

BLOSUM matrices

Seq1 MCL
Seq2 CCV
Seq3 ICV
Seq4 MAI
Seq5 TCL

Observed amino acid pairs:
 $CC + CC + CA + CC$
 $+ CC + CA + CC$
 $+ CA + CC$

BIOS477/877 L9 - 28

28

BLOSUM matrices

Seq1 MCL
Seq2 CCV
Seq3 ICV
Seq4 MAI
Seq5 TCL

Observed amino acid pairs:
 $CC + CC + CA + CC$
 $+ CC + CA + CC$
 $+ CA + CC$
 $+ AC$

BIOS477/877 L9 - 29

29

BLOSUM matrices

Seq1 MCL
Seq2 CCV
Seq3 ICV
Seq4 MAI
Seq5 TCL

Observed amino acid pairs:
 $CC + CC + CA + CC$
 $+ CC + CA + CC$
 $+ CA + CC$
 $+ AC$
 $[3CA + AC]$
 $= 6CC + 4CA$

Each column has 10 pairs
→ total 30 pairs
for 3 columns

BIOS477/877 L9 - 30

30

BLOSUM matrices

- **Observed amino acid pairs:** 6CC, 4CA (from 30 pairs total)

Seq3	TCV
Seq4	MAL
Seq5	TCI
 - **Observed frequency of pairs in the alignment:**

$$q_{CC} = 6/30 = 0.2, q_{CA} = 4/30 = 0.133$$

(There are 20 other pairs, too)
 - **Observed frequency of each amino acid in the alignment:**

$$p_i = q_{ii} + \sum_{j \neq i} q_{ij} / 2$$

$$\rightarrow p_C = 0.2 + 0.133/2 = 0.267 \text{ and } p_A = 0.133/2 = 0.067$$
 - **Expected frequencies of amino acid pairs:**

$$e_{ii} = p_i p_i = p_i^2 \text{ and } e_{ij} = p_i p_j + p_j p_i = 2 p_i p_j$$

$$\rightarrow e_{CC} = 0.267^2 = 0.071$$

$$e_{AA} = 0.067^2 = 0.004$$

$$e_{CA} = 2 \times 0.267 \times 0.067 = 0.036$$
- BIOS477/877 L9 - 31

BIOS477/877 L9 - 31

31

BLOSUM matrices

- **Observed and Expected frequencies** of amino acid pairs are cumulatively counted from all columns of the BLOCKS
- **Log odds scores** are calculated for each amino acid pairs:

$$S_{ij} = \log(q_{ij}/e_{ij})$$

In bit units: $S_{ij} = \log_2 (q_{ij}/e_{ij})$

Usually in **half-bit units**: $S_{ij} = 2\log_2(q_{ij}/e_{ij})$

***bit = binary digit (0 or 1)**

BIOS477/877 L9 - 32

32

BLOSUM matrices

- **Observed and Expected frequencies** of amino acid pairs are cumulatively counted from all columns of the BLOCKS
- **Log odds scores** are calculated for each amino acid pairs:

$$S_{ij} = 2\log_2 (q_{ij}/e_{ij})$$

From the example:

$$q_{CC} = 0.2, e_{CC} = 0.071 \rightarrow S_{CC} = 2\log_2(0.2/0.071) = 2.99$$

$$q_{CA} = 0.133, e_{AC} = 0.036 \rightarrow S_{CA} = 2\log_2(0.133/0.036) = 3.77$$

BIOS477/877 L9 - 33

33

BLOSUM62 matrix

Log odds matrix (Henikoff and Henikoff 1992)

[illegible]

Note: Usually BLOSUM matrices are shown in **half-bit units**.

BIOS477/877 L9 - 34

34

BLOSUM matrices

- **BLOSUM_{*n*}** represents the similarity threshold (*e.g.*, BLOSUM₆₂, BLOSUM₄₅, BLOSUM₈₀)
 - for any *n*, the corresponding BLOSUM matrix is generated mainly comparing sequences that are **less than *n*% identical**
 - e.g.*, BLOSUM₆₂: Sequences with **≥62%** identity are clustered and treated as a single sequence for counting.
 - 3 sequences in one cluster are counted as 1/3 each, instead of 1 each for counting.
 - **All BLOSUM matrices are based on observed alignments; they are not extrapolated from comparisons of closely related proteins**
-
- BIOS477/877 L9 - 35

BIOS477/877 L9 - 35

35

BLOSUM and PAM matrices

- **PAM matrices:** based on mutational model of evolution
 - a transition probability matrix for a Markov process
 - any M_n matrix can be extrapolated based on PAM1 (M_1) matrix ($e.g., M_{250} = M_1^{250}$)
 - assume more distant changes are a reflection of the short-term changes
 - designed to track the evolutionary origins of proteins
- **BLOSUM matrices:** not based on explicit evolutionary model
 - based on local similarity
 - derived from all changes observed in the conserved blocks regardless of the overall degree of similarity
 - generated based on different similarity levels (BLOSUM50, BLOSUM62, etc.)
 - all BLOSUM matrices are generated based on observed data
 - designed to find conserved domains

BIOS477/877 L9 - 36

36

Log Odds Matrix

	AA ₁	AA ₂
AA ₁	S ₁₁	S ₂₁
AA ₂	S ₁₂	S ₂₂

➤ Log odds (Lod) score: PAM matrix

$$S(i,j) = 10 \log_{10} \{M(i,j)/f(i)\} = 10 \log_{10} \{f(i,j)/\{f(i) \times f(j)\}\}$$

$M(i,j)$: Mutation probability $AA_i \rightarrow AA_j$

$f(i,j)$: Frequency of accepted point mutations $AA_i \leftrightarrow AA_j$

$f(i), f(j)$: (frequency of AA_i or AA_j) / (total no. residues)

$\sum f(i,j)$: Frequency of AA_j substituted by any other AA

$m(j) = \sum_i f(i,j) / f(j)$: Relative mutability of AA_j

$$\begin{aligned} S(i,j) &= 10 \log_{10} \{m(j) \times f(i,j) / \sum_i f(i,j) / f(i)\} \\ &= 10 \log_{10} \{[\sum_i f(i,j) \times f(i,j)] / [\sum_i f(i,j) \times f(i) \times f(j)]\} \\ &= 10 \log_{10} \{f(i,j) / \{f(i) \times f(j)\}\} \end{aligned}$$

$$S(j,i) = 10 \log_{10} \{f(j,i) / \{f(j) \times f(i)\}\} = S(i,j)$$

→ PAM matrix is symmetrical: $S(i,j) = S(j,i)$

BIOS477/877 L9 - 37

37

Log Odds Matrix (PAM250)

PAM matrix is symmetrical

$$S(j,i) = S(i,j)$$

C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
S	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
T	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
P	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
G	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
N	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
D	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Q	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
H	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
R	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
M	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
I	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
L	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
V	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
F	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Y	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
W	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Dayhoff et al. (1978)

BIOS477/877 L9 - 38

38

Log Odds Matrix

	AA ₁	AA ₂
AA ₁	S ₁₁	S ₂₁
AA ₂	S ₁₂	S ₂₂

➤ Log odds (Lod) score: BLOSUM matrix

$$S(i,j) = 2 \log_2 (q_{ij}/e_{ij})$$

q_{ij} : Observed frequency of (AA_i, AA_j) pairs

e_{ij} : Expected frequencies of (AA_i, AA_j) pairs

$$e_{ii} = p_i p_i = p_i^2 \text{ and } e_{ij} = p_i p_j + p_j p_i = 2p_i p_j$$

p_i : Observed frequency of AA_i in the pairs

$$p_i = q_{ii} + \sum_{j \neq i} q_{ij} / 2$$

→ BLOSUM matrix is symmetrical: $S(i,j) = S(j,i)$

BIOS477/877 L9 - 39

39

Log Odds Matrix (BLOSUM62)

BLOSUM matrix is symmetrical

$$S(j,i) = S(i,j)$$

C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
S	4	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
T	1	5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
P	0	0	7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
G	0	0	0	0	6	0	0	0	0	0	0	0	0	0	0	0	0	0	0
N	0	0	0	0	0	6	0	0	0	0	0	0	0	0	0	0	0	0	0
D	0	0	0	0	0	0	6	0	0	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	5	0	0	0	0	0	0	0	0	0	0	0
Q	0	0	0	0	0	0	0	0	5	0	0	0	0	0	0	0	0	0	0
H	0	0	0	0	0	0	0	0	0	8	0	0	0	0	0	0	0	0	0
R	0	0	0	0	0	0	0	0	0	0	5	0	0	0	0	0	0	0	0
K	0	0	0	0	0	0	0	0	0	0	0	5	0	0	0	0	0	0	0
M	0	0	0	0	0	0	0	0	0	0	0	0	5	0	0	0	0	0	0
I	0	0	0	0	0	0	0	0	0	0	0	0	0	5	0	0	0	0	0
L	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	0	0	0	0
V	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	0	0	0
F	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	0	0
Y	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	0
W	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4

Henikoff and Henikoff (1992)

BIOS477/877 L9 - 40

40

Log Odds Matrix

• PAM matrix

$$S(i,j) = 10 \log_{10} \{M(i,j)/f(i)\}$$

$M(i,j)$: Mutation probability from AA_i to AA_j

$f(i)$: Frequency of AA_i (number of AA_i / total number of residues)

Probability to find AA_i by chance

• BLOSUM matrix

$$S(i,j) = 2 \log_2 (q_{ij}/e_{ij})$$

q_{ij} : Observed frequency of AA_i, AA_j pairs

e_{ij} : Expected frequencies of AA_i, AA_j pairs

• General form

$$S(i,j) = 1/\lambda \log_2 (q_{ij}/p_i p_j) \text{ [in bit unit]}$$

$$S(i,j) = 1/\lambda \log_e (q_{ij}/p_i p_j) \text{ [in nat unit]}$$

BIOS477/877 L9 - 41

41

Log Odds Matrix

	AA ₁	AA ₂
AA ₁	S ₁₁	S ₂₁
AA ₂	S ₁₂	S ₂₂

➤ Log odds (Lod) score: general also called log odds ratio or log likelihood ratio

$$S(i,j) = 1/\lambda \log_2 (q_{ij}/p_i p_j) \text{ [in bit unit]}$$

$$S(i,j) = 1/\lambda \log_e (q_{ij}/p_i p_j) \text{ [in nat unit]}$$

q_{ij} : the frequency of the amino acid pair, AA_i and AA_j

p_i, p_j : the individual frequency of AA_i or AA_j

λ : a scaling factor

If $\lambda=1/2$,
 $S(i,j) = 2 \log_2 (q_{ij}/p_i p_j)$
[a half-bit unit]

(1/λ=2 is used with BLOSUM)

$$S(i,j) = 1/\lambda \log \{ \text{Observed freq. of amino acid pair } i \text{ and } j / \text{Expected freq. of amino acid pair } i \text{ and } j \}$$

* In the general format, substitutions does not have to be symmetrical.
 $S_{12} = S_{21}$ is not assumed.

BIOS477/877 L9 - 42

42

Log Odds Matrix

	AA_1	AA_2
AA_1	S_{11}	S_{21}
AA_2	S_{12}	S_{22}

➤ **Log odds (Lod) score:** general
also called **log odds ratio** or **log likelihood ratio**

$S(i,j) = 1/\lambda \log_2(q_{ij}/p_i p_j)$ [in bit unit]
 $S(i,j) = 1/\lambda \log_e(q_{ij}/p_i p_j)$ [in nat unit]

Target frequency (q_{ij})

$S(i,j) = 1/\lambda \log \left\{ \frac{\text{Observed freq. of amino acid pair } i,j}{\text{Expected freq. of amino acid pair } i,j} \right\}$

Background frequency ($p_i p_j$)

$[- < S(i,j) < +]$

H_1 : Homologous hypothesis (residues i and j are related)
 H_0 : Random hypothesis (residues i and j are unrelated)

BIOS477/877 L9 - 43

43

Log Odds Matrix

	AA_1	AA_2
AA_1	S_{11}	S_{21}
AA_2	S_{12}	S_{22}

➤ **Log odds (Lod) score:** general
also called **log odds ratio** or **log likelihood ratio**

$S(i,j) = 1/\lambda \log_2(q_{ij}/p_i p_j)$ [in bit unit]
 $S(i,j) = 1/\lambda \log_e(q_{ij}/p_i p_j)$ [in nat unit]

Likelihood ratio (LR) = $\frac{\text{Likelihood of } H_1}{\text{Likelihood of } H_0}$

$[0 < LR < +\infty]$ **= $\frac{\text{Prob}(\text{an event}|H_1)}{\text{Prob}(\text{an event}|H_0)}$**

H_1 : Hypothesis to be tested, H_0 : Null hypothesis

BIOS477/877 L9 - 44

44

Log Odds Matrix

	AA_1	AA_2
AA_1	S_{11}	S_{21}
AA_2	S_{12}	S_{22}

➤ **Log odds (Lod) score:** general
also called **log odds ratio** or **log likelihood ratio**

$S(i,j) = 1/\lambda \log_2(q_{ij}/p_i p_j)$ [in bit unit]
 $S(i,j) = 1/\lambda \log_e(q_{ij}/p_i p_j)$ [in nat unit]

Log likelihood ratio = $\log \left\{ \frac{\text{Likelihood of } H_1}{\text{Likelihood of } H_0} \right\}$

= $\log \{ \text{Prob}(\text{an event}|H_1) \} - \log \{ \text{Prob}(\text{an event}|H_0) \}$

[- < log(LR) < +]

H_1 : Hypothesis to be tested, H_0 : Null hypothesis

BIOS477/877 L9 - 45

45

Log Odds Score and Target Frequencies

$S(i,j) = 1/\lambda \log_e(q_{ij}/p_i p_j)$
[or $S(i,j) = 1/\lambda \log_2(q_{ij}/p_i p_j)$ for BLOSUM]

$\lambda S(i,j) = \log_e(q_{ij}/p_i p_j)$
 $e^{\lambda S(i,j)} = q_{ij}/p_i p_j$

$q_{ij} = p_i p_j e^{\lambda S(i,j)}$

Target frequency Expected (or background) frequency

$\sum_i \sum_j q_{ij} = \sum_i \sum_j p_i p_j e^{\lambda S(i,j)} = 1$
(i < j)

λ can be estimated (matrix specific)

BIOS477/877 L9 - 46

46