

BIOS 477/877 Bioinformatics and Molecular Evolution

Instructor: Etsuko Moriyama
(School of Biological Sciences)

Spring 2026

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

Amino acid substitution matrices based on empirical data

- **PAM matrices**
Dayhoff, Schwartz, and Orcutt (1978)
- **BLOSUM matrices**
Henikoff and Henikoff (1992)
- Also see Eddy (2004)

BIOS477/877 L9 - 3

3

PAM matrices (Dayhoff et al. 1978)

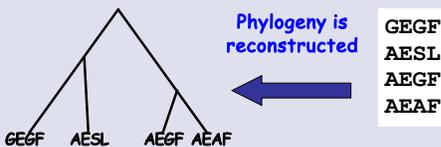
- **Accepted point mutations** (point accepted mutations, percent accepted mutations)
 - ➔ **accepted by selection**: no (or very weak) deleterious effect, maintaining the function
- Based on 1,572 changes in 71 groups of **closely related proteins** (34 protein families)
 - at least 85% identical
 - 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

Accepted point mutations: $f(a,b)$

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

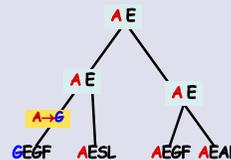


BIOS477/877 L9 - 5

5

Accepted point mutations: $f(a,b)$

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

Accepted point mutations: $f(a,b)$

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

Accepted point mutations: $f(a,b)$

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

Accepted point mutations: $f(a,b)$

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

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

BIOS477/877 L9 - 9

9

Accepted point mutations: $f(a,b)$

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

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

Ala Arg Asn Asp Cys Gln Glu Gly His Ile Leu Lys Met Phe Pro Ser Thr Trp Tyr Val
 Ala Arg Asn Asp Cys Gln Glu Gly His Ile Leu Lys Met Phe Pro Ser Thr Trp Tyr Val

Figure 80. Numbers of accepted point mutations (X10) accumulated from closely related sequences. Fifteen hundred and seventy-two exchanges are shown. Fractional exchanges result when ancestral sequences are ambiguous.
Dayhoff et al. (1978)

BIOS477/877 L9 - 10

10

Relative mutability: $m(a)$

- **Relative mutability**: $m(a)$
Probability that the amino acid a will change in a given small evolutionary interval

| Amino acid: a | A | E | F | G |
|---|---|---|---|-----|
| Changes: $\sum f(i,a)$ | 1 | 0 | 0 | 1 |
| Freq. of occurrence: $f(a)$ (Total # of the residue) | 1 | 2 | 2 | 3 |
| Relative mutability: $m(a)$ | 1 | 0 | 0 | 0.3 |

$m(a) = \frac{\sum f(i,a)}{f(a)}$
 [combined from multiple trees] $m(a) = \frac{\text{Number of times amino acid } a \text{ is substituted by any other amino acid}}{\sum_{\text{branch}} \{(\text{Freq. of amino acid } a) \times (\text{Number of total substitutions}) \times 100\}}$
 This denominator is called "the total exposure of the amino acid to mutation"

Substitutions are collected from trees with different lengths

BIOS477/877 L9 - 11

11

Relative mutability: $m(a)$

Table 21
Relative Mutabilities of the Amino Acids^a

| | | | | |
|--------------|-----|-----|-----|----|
| More 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 |
| Less mutable | | | | |

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

Dayhoff et al. (1978)

BIOS477/877 L9 - 12

12

Mutation probability: $M(a,b)$

- 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

Mutation probability matrix: M

Probabilities of AA_j replaced by AA_i

| | | ORIGINAL AMINO ACID | | | | | | | | | | | | | | | | | | | | | | |
|---|-------|---------------------|--------|--------|------|--------|--------|--------|------|--------|------|--------|--------|--------|------|--------|------|--------|--------|------|--------|------|--------|-----|
| | | N | I | A | R | K | H | D | C | Q | E | G | S | T | L | V | M | F | P | S | Y | W | V | Val |
| ACCEPTED POINT MUTATION PER 100 AMINO ACIDS | A Ala | 0.9967 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | R Arg | 0 | 0.9910 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | M Asn | 0 | 0 | 0.9822 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | D Asp | 0 | 0 | 0 | 0.42 | 0.9859 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | C Cys | 0 | 0 | 0 | 0 | 0 | 0.9970 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | G Glu | 0 | 0 | 0 | 0 | 0 | 0 | 0.9874 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | E Glu | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.35 | 0.9865 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | G Gly | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.37 | 0.9935 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | H His | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.9912 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | I Ile | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.9832 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | L Leu | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.22 | 0.9947 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | K Lys | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.19 | 0.9920 | 0 | 0 | 0 | 0 | 0 | 0 |
| | M Met | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.9874 | 0 | 0 | 0 | 0 | 0 |
| | F Phe | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.49 | 0.9946 | 0 | 0 | 0 |
| | P Pro | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.19 | 0.9926 | |
| | S Ser | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | Y Thr | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | W Trp | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | V Tyr | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | V Val | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Figure B2: 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 j will be replaced by the amino acid in row i after a given evolutionary interval, in this case 1 accepted point mutation per 100 amino acids. Thus, there is a 0.986% probability that Asp will be replaced by Glu. To simplify the appearance, the elements are shown multiplied by 10,000.

Dayhoff et al. (1978) BIOS477/877 L9 - 14

14

Mutation probability matrix: M

Probabilities of AA_j replaced by AA_i

| | | ORIGINAL AMINO ACID | | | | | | | | | | | | | | | | | | | | | | |
|---|-------|---------------------|--------|--------|------|--------|--------|--------|------|--------|--------|--------|--------|------|--------|---|------|--------|--------|------|--------|------|--------|-----|
| | | N | I | A | R | K | H | D | C | Q | E | G | S | T | L | V | M | F | P | S | Y | W | V | Val |
| ACCEPTED POINT MUTATION PER 100 AMINO ACIDS | A Ala | 0.9967 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | R Arg | 0 | 0.9910 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | M Asn | 0 | 0 | 0.9822 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | D Asp | 0 | 0 | 0 | 0.42 | 0.9859 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | C Cys | 0 | 0 | 0 | 0 | 0 | 0.9970 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | G Glu | 0 | 0 | 0 | 0 | 0 | 0 | 0.9874 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | E Glu | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.35 | 0.9865 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | G Gly | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.37 | 0.9935 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | H His | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.9912 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | I Ile | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.9832 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | L Leu | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.22 | 0.9947 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | K Lys | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.19 | 0.9920 | 0 | 0 | 0 | 0 | 0 | 0 |
| | M Met | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.9874 | 0 | 0 | 0 | 0 | 0 |
| | F Phe | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.49 | 0.9946 | 0 | 0 | 0 |
| | P Pro | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.19 | 0.9926 | |
| | S Ser | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | Y Thr | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | W Trp | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | V Tyr | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | V Val | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Ala will not be changed with 98.67% probability

Gly will be replaced by Ala with 0.21% probability

Figure B2: 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 j will be replaced by the amino acid in row i after a given evolutionary interval, in this case 1 accepted point mutation per 100 amino acids. Thus, there is a 0.986% probability that Asp will be replaced by Glu. To simplify the appearance, the elements are shown multiplied by 10,000.

Dayhoff et al. (1978) BIOS477/877 L9 - 15

15

Mutation probability matrix: M

Probabilities of AA_j replaced by AA_i

| | | ORIGINAL AMINO ACID | | | | | | | | | | | | | | | | | | | | | | |
|---|-------|---------------------|--------|--------|------|--------|--------|--------|------|--------|--------|--------|--------|------|--------|---|------|--------|--------|------|--------|------|--------|-----|
| | | N | I | A | R | K | H | D | C | Q | E | G | S | T | L | V | M | F | P | S | Y | W | V | Val |
| ACCEPTED POINT MUTATION PER 100 AMINO ACIDS | A Ala | 0.9967 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | R Arg | 0 | 0.9910 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | M Asn | 0 | 0 | 0.9822 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | D Asp | 0 | 0 | 0 | 0.42 | 0.9859 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | C Cys | 0 | 0 | 0 | 0 | 0 | 0.9970 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | G Glu | 0 | 0 | 0 | 0 | 0 | 0 | 0.9874 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | E Glu | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.35 | 0.9865 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | G Gly | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.37 | 0.9935 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | H His | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.9912 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | I Ile | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.9832 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | L Leu | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.22 | 0.9947 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | K Lys | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.19 | 0.9920 | 0 | 0 | 0 | 0 | 0 | 0 |
| | M Met | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.9874 | 0 | 0 | 0 | 0 | 0 |
| | F Phe | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.49 | 0.9946 | 0 | 0 | 0 |
| | P Pro | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.19 | 0.9926 | |
| | S Ser | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | Y Thr | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | W Trp | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | V Tyr | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | V Val | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Mutation probability matrix is not symmetrical

Figure B2: 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 j will be replaced by the amino acid in row i after a given evolutionary interval, in this case 1 accepted point mutation per 100 amino acids. Thus, there is a 0.986% probability that Asp will be replaced by Glu. To simplify the appearance, the elements are shown multiplied by 10,000.

Dayhoff et al. (1978) BIOS477/877 L9 - 16

16

PAM scores: $R(a,b)$

- 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)

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

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

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

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

$R(a,b) = R(b,a)$
 → Relatedness odds score matrix is symmetrical

BIOS477/877 L9 - 17

17

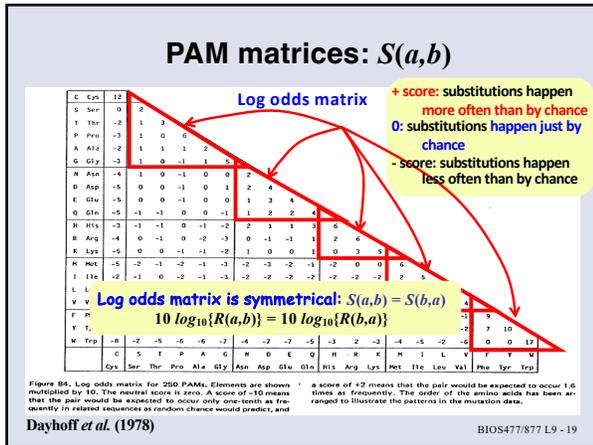
PAM scores: $S(a,b)$

- Relatedness odds score**

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

Odds ratio ↓

$\frac{P(a \text{ was derived from } b)}{P$



19

PAM1 matrix

- **PAM1 matrix**
 - normalized 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)
- **M₁: 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)
- **PAM1: from $S(a,b) = 10 \log_{10}\{M_1(a,b) / f(a)\}$,** where $M(a,b) = \lambda_m(b) \times f(a,b) / \sum_a f(a,b)$

Normalizing (scaling) factor

BIOS477/877 L9 - 20

20

PAM1 and PAMn

- **PAM1 matrix**
 - normalized 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)
- **M₁: 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 PAMn**
 - $M_n = (M_1)^n$ e.g., M_{250} (for PAM250) = M_1^{250}
 - Probability matrix after evolutionary interval of PAM250 (after 250 changes are produced in 100 aa)

BIOS477/877 L9 - 21

21

PAM distance

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 |
| 45 | 65 |
| 50 | 75 |
| 55 | 85 |
| 60 | 95 |
| 65 | 105 |
| 70 | 115 |
| 75 | 125 |
| 80 | 135 |
| 85 | 145 |
| 90 | 155 |
| 95 | 165 |
| 100 | 175 |

Annotations: A blue arrow points from the text $100[1 - \sum_b \{f(b)M_n(b,b)\}]$ for PAMn to the 'Observed Percent Difference' column. A red arrow points from 'PAM1' to the value '1' in the 'Evolutionary Distance' column. A yellow box highlights the values '20' and '30' in the 'Evolutionary Distance' column with the text 'PAM250 ≈ 20% similarity'. A red box highlights the values '80' and '246' in the 'Evolutionary Distance' column with the text '~PAM250'.

Dayhoff et al. (1978)

BIOS477/877 L9 - 22

22

PAM matrices updated

- **JTT matrices** (Jones, Taylor, and Thornton 1992)
 - Based on 59,190 accepted point mutations for 16,300 proteins
- **Gonnet matrices** (Gonnet, Cohen, and 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   GPCRRHODOPSN; 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_LOCM1|Q25321 ( 53) VTAVSLSLILITIVGNVLVLSVF 3
OAR2_LOCM1|Q25322 ( 53) VTAVSLSLILITIVGNVLVLSVF 3
OAL110 ( 22) ISLAVLEFINVLVGGNCLVIAAVF 22
OAR_DROME|P22270 ( 111) LITLVLSVIVLTIIGNILVLSVF 3
DOP2_DROME|Q24563 ( 110) GLLAFLELFSFATVFGNSLVILAVI 5
OAR_HELVI|Q25188 ( 54) CTAIVLTLIIISTIVGNILVLSVF 6
O31124 ( 35) ISLLALAPLNLMVAGNLLVMVAVF 9
A1AB_MESAU|F18841 ( 47) SVGLVLGAFILFAIVGNILVLSVA 5
A1AB_RAT|F15823 ( 47) SVGLVLGAFILFAIVGNILVLSVA 5
A1AB_HUMAN|P35368 ( 47) SVGLVLGAFILFAIVGNILVLSVA 5
O31127 ( 34) ANVALLLAILITIVGNLSVLSVF 3
A1AA_ORYLA|Q91175 ( 28) YLGMVLGIFLFGVGNILVLSVV 5
D2D1_XENLA|P24628 ( 30) YYAMLLTLLVVFVIFGNVLCIAVS 6
D2DR_CERAR|P52702 ( 36) YYATLTLTLLIIVIFGNVLCIAVS 5
O31119 ( 33) YYAVLTLTLLIIVIFGNVLCIAVS 5
OAR_ROMMO|Q17232 ( 56) CTAIILTMIIISTVVGNIIVLSVF 7
SRR5_MOUSE|Q08858 ( 39) LVPVLLVLTCTVGLGNTLVIVVVL 4
    
```

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

BIOS477/877 L9 - 25

25

BLOSUM matrices: how to count

Seq1 MCL
Seq2 GCY
Seq3 ICY
Seq4 MAI
Seq5 TCL

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

BIOS477/877 L9 - 26

26

BLOSUM matrices: how to count

Seq1 MCL
Seq2 GCY
Seq3 ICY
Seq4 MAI
Seq5 TCL

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

BIOS477/877 L9 - 27

27

BLOSUM matrices: how to count

Seq1 MCL
Seq2 GCY
Seq3 ICY
Seq4 MAI
Seq5 TCL

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

BIOS477/877 L9 - 28

28

BLOSUM matrices: how to count

Seq1 MCL
Seq2 GCY
Seq3 ICY
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: how to count

Seq1 MCL
Seq2 GCY
Seq3 ICY
Seq4 MAI
Seq5 TCL

Observed amino acid pairs:
CC + CC + CA + CC
+ CC + CA + CC
+ CA + CC
+ AC [3CA+AC]

Each column has 10 pairs
→ total 30 pairs
for 3 columns
= 6CC + 4CA

BIOS477/877 L9 - 30

30

Observed vs. estimated AA pairs

- Observed amino acid pairs in the alignment: **6CC, 4CA** (from 30 pairs total)
- Observed frequency of pairs
 $q_{CC} = 6/30 = 0.2$, $q_{CA} = 4/30 = 0.133$ (There are 20 other pairs, too)
- Observed frequency of each amino acid:
 $p_i = q_{ii} + \sum_{j \neq i} q_{ij}/2$
 $p_C = 0.2 + 0.133/2 = 0.267$ and $p_A = 0.13/2 = 0.067$
- Expected frequencies of amino acid pairs:
 $e_{ii} = p_i p_i = p_i^2$ and $e_{ij} = p_i p_j + p_j p_i = 2p_i p_j$
 $e_{CC} = 0.267^2 = 0.071$, $e_{AA} = 0.067^2 = 0.004$,
 $e_{CA} = 2 \times 0.267 \times 0.067 = 0.036$

Seq1 MCL
Seq2 GCV
Seq3 TCV
Seq4 MAL
Seq5 TCI

BIOS477/877 L9 - 31

31

BLOSUM scores

- Frequencies of amino acid pairs are cumulatively counted from all columns of the BLOCKs alignment
- Log odds score is calculated for each amino acid pair:
 - $S_{ij} = \log(q_{ij}/e_{ij})$
 - 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)

Seq1 MCL
Seq2 GCV
Seq3 TCV
Seq4 MAL
Seq5 TCI

[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 - 32

32

BLOSUM62 matrix

| | | | | | | | | | | | | | | | | | | | |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| C | S | T | P | A | G | N | D | E | Q | H | R | K | M | I | L | V | F | Y | W |
| S | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| T | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| P | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| A | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| G | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| N | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| D | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| E | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| Q | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| H | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| R | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| K | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| M | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 | -1 |
| I | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 | -1 |
| L | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 | -1 |
| V | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 | -1 |
| F | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 | -1 |
| Y | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 | -1 |
| W | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 4 |

+ score: substitutions happen more often than by chance
0: substitutions happen just by chance
- score: substitutions happen less often than by chance

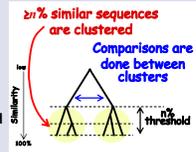
Note: Usually BLOSUM matrices are shown in half-bit units.

BIOS477/877 L9 - 33

33

BLOSUM matrices

- BLOSUMn: 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., For BLOSUM62:
 - Sequences with $\geq 62\%$ 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 - 34

34

Log odds matrix

- Log odds (Lod) score: **PAM matrix**
 $S(i,j) = 10\log_{10}[M(i,j)/f(i)] = 10\log_{10}\{f(i) \times f(j)\}$
 $M(i,j) = m(j) \times f(i,j) / \sum_j f(i,j)$: 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)
 $\sum_j f(i,j)$: Frequency of AA_i substituted by any other AA
 $m(j) = \sum_i f(i,j) / f(j)$: Relative mutability of AA_j
 $S(i,j) = 10\log_{10}\{m(j) \times f(i,j) / \sum_j f(i,j) / f(i)\}$
 $= 10\log_{10}\{[\sum_i f(i,j) \times f(i,j)] / [\sum_i f(i,j) \times f(i) \times f(j)]\}$
 $= 10\log_{10}\{f(i,j) / [f(i) \times f(j)]\}$
 $S(j,i) = 10\log_{10}\{f(j,i) / [f(j) \times f(i)]\} = S(i,j)$
 → PAM matrix is symmetrical: $S(i,j) = S(j,i)$

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

BIOS477/877 L9 - 35

35

Log odds matrix (PAM250)

| | | | | | | | | | | | | | | | | | | | |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| C | S | T | P | A | G | N | D | E | Q | H | R | K | M | I | L | V | F | Y | W |
| S | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| T | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| P | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| G | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| N | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| D | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| E | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Q | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| H | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| R | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| K | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| M | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 |
| I | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 |
| L | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 0 |
| V | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 |
| F | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 |
| Y | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 |
| W | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 |

Dayhoff et al. (1978)

BIOS477/877 L9 - 36

36

Relation to hypothesis testing

- Hypothesis testing
 - Test a hypothesis (H_1) against a null hypothesis (H_0)
 H_1 : also called as "alternative hypothesis"
 - Likelihood ratio test: likelihood (probability of an event given a hypothesis) of each hypothesis is compared

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

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

$$[- < \log(\text{LR}) < +]$$

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

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

↖ Target frequency (q_{ij})
↖ 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 - 44

44

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

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

λ can be estimated (matrix specific)

BIOS477/877 L9 - 45

45