

Spring 2024
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)

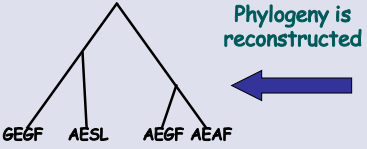
- **Accepted point mutations** (point accepted mutations, percent accepted mutations)
 - **accepted by selection**: no (or very weak) deleterious effect, maintaining the function
- Based on 1,572 changes in 71 groups of **closely related proteins** (34 protein families)
 - at least 85% identical
 - no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

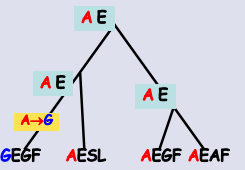
GEGF
 AESL
 AEGF
 AEF

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.*, $G \rightarrow A = A \rightarrow G$)



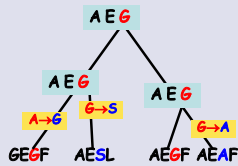
Using the maximum parsimony principle, ancestral sequences can be inferred

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., $G \rightarrow A = A \rightarrow G$)



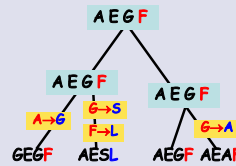
Using the maximum parsimony principle, ancestral sequences can be inferred

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., $G \rightarrow A = A \rightarrow G$)



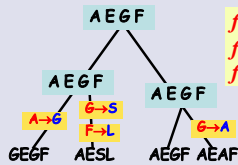
Substitutions can be identified along the phylogeny

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., $G \rightarrow A = A \rightarrow G$)



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

BIOS477/877 L9 - 9

9

PAM matrices

Numbers of accepted point mutations: $f(a,b)$
 Dayhoff et al. (1978)

Table showing the number of accepted point mutations between amino acids. The matrix is upper triangular. A yellow box notes: "Based on 1,572 changes, but still missing 35 types of substitutions".

Figure 80. Numbers of accepted point mutations. Ex 100 occurrences are shown. Fractional exchanges result when ancestral sequences are ambiguous.

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	G
Changes: $\sum f(i,a)$	1	0	1	2
Freq. of occurrence: $f(a)$ (Total # of the residue)	1	2	1	4
Relative mutability: $m(a)$	1	0	1	0.5

[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

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 to 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)

		ORIGINAL AMINO ACID																			
		I	A	R	N	D	C	Q	E	G	H	L	K	M	F	P	S	T	W	Y	V
REPLACEMENT AMINO ACID	I	Ala	9869	2	9	30	2	8	37	23	2	6	4	2	2	22	35	32	6	2	18
	A	Arg	1	9913	1	0	1	10	0	10	3	1	19	4	1	4	4	1	8	0	1
	R	Asn	4	1	9822	36	0	4	6	21	3	1	13	0	1	2	20	9	1	4	1
	N	Asp	6	0	42	9859	0	6	52	8	4	1	0	0	0	1	5	3	0	0	1
	D	Cys	1	1	0	0	9932	0	0	1	1	0	0	0	0	0	0	1	1	0	3
	C	Glu	3	9	4	5	0	9876	27	1	23	1	3	6	4	0	6	2	0	0	1
	Q	Gly	10	0	2	56	0	35	9865	4	2	3	1	4	1	0	3	4	2	0	1
	E	Gly	21	1	12	13	1	3	7	9955	1	0	1	1	0	1	0	21	3	0	5
	G	His	1	6	10	3	1	20	1	0	9932	0	1	3	0	2	3	1	1	1	4
	H	His	2	2	3	1	2	1	2	0	0	9872	9	2	12	7	0	1	7	0	1
	L	Leu	3	1	3	0	0	4	1	1	4	22	9947	2	45	13	3	1	3	4	-2
	K	Lys	2	37	25	6	0	12	7	2	2	4	1	9926	20	0	2	8	11	0	1
	M	Met	1	1	0	0	0	2	0	0	0	5	8	0	9874	1	0	1	2	0	4
	F	Phe	1	1	1	0	0	0	0	1	2	8	6	0	4	9946	0	2	1	3	20
	P	Pro	13	5	2	1	1	0	3	2	5	1	2	2	1	1	9926	12	4	0	2
	S	Ser	28	11	34	7	11	4	6	16	2	2	1	7	4	3	10	9840	34	5	2
	T	Thr	22	2	13	4	1	3	2	0	1	11	2	8	6	1	1	32	9871	0	2
	W	Tyr	0	2	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	9976	1
	Y	Tyr	1	0	3	0	3	0	1	0	4	1	1	0	0	21	0	1	1	1	9945
	V	Val	13	2	1	1	3	2	3	3	5	11	1	1	1	3	2	10	0	2	9903

Figure B2. Mutation probability matrix for the evolutionary 1 accepted point mutation per 100 amino acids. Thus, there is a 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

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)

		ORIGINAL AMINO ACID																			
		I	A	R	N	D	C	Q	E	G	H	L	K	M	F	P	S	T	W	Y	V
REPLACEMENT AMINO ACID	I	Ala	9869	2	9	30	2	8	37	23	2	6	4	2	2	22	35	32	6	2	18
	A	Arg	1	9913	1	0	1	10	0	10	3	1	19	4	1	4	4	1	8	0	1
	R	Asn	4	1	9822	36	0	4	6	21	3	1	13	0	1	2	20	9	1	4	1
	N	Asp	6	0	42	9859	0	6	52	8	4	1	0	0	0	1	5	3	0	0	1
	D	Cys	1	1	0	0	9932	0	0	1	1	0	0	0	0	0	0	1	1	0	3
	C	Glu	3	9	4	5	0	9876	27	1	23	1	3	6	4	0	6	2	0	0	1
	Q	Gly	10	0	2	56	0	35	9865	4	2	3	1	4	1	0	3	4	2	0	1
	E	Gly	21	1	12	13	1	3	7	9955	1	0	1	1	0	1	0	21	3	0	5
	G	His	1	6	10	3	1	20	1	0	9932	0	1	3	0	2	3	1	1	1	4
	H	His	2	2	3	1	2	1	2	0	0	9872	9	2	12	7	0	1	7	0	1
	L	Leu	3	1	3	0	0	4	1	1	4	22	9947	2	45	13	3	1	3	4	-2
	K	Lys	2	37	25	6	0	12	7	2	2	4	1	9926	20	0	2	8	11	0	1
	M	Met	1	1	0	0	0	2	0	0	0	5	8	0	9874	1	0	1	2	0	4
	F	Phe	1	1	1	0	0	0	0	1	2	8	6	0	4	9946	0	2	1	3	20
	P	Pro	13	5	2	1	1	0	3	2	5	1	2	2	1	1	9926	12	4	0	2
	S	Ser	28	11	34	7	11	4	6	16	2	2	1	7	4	3	10	9840	34	5	2
	T	Thr	22	2	13	4	1	3	2	0	1	11	2	8	6	1	1	32	9871	0	2
	W	Tyr	0	2	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	9976	1
	Y	Tyr	1	0	3	0	3	0	1	0	4	1	1	0	0	21	0	1	1	1	9945
	V	Val	13	2	1	1	3	2	3	3	5	11	1	1	1	3	2	10	0	2	9903

Figure B2. Mutation probability matrix for the evolutionary 1 accepted point mutation per 100 amino acids. Thus, there is a 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

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)

		ORIGINAL AMINO ACID																			
		I	A	R	N	D	C	Q	E	G	H	L	K	M	F	P	S	T	W	Y	V
REPLACEMENT AMINO ACID	I	Ala	9869	2	9	30	2	8	37	23	2	6	4	2	2	22	35	32	6	2	18
	A	Arg	1	9913	1	0	1	10	0	10	3	1	19	4	1	4	4	1	8	0	1
	R	Asn	4	1	9822	36	0	4	6	21	3	1	13	0	1	2	20	9	1	4	1
	N	Asp	6	0	42	9859	0	6	52	8	4	1	0	0	0	1	5	3	0	0	1
	D	Cys	1	1	0	0	9932	0	0	1	1	0	0	0	0	0	0	1	1	0	3
	C	Glu	3	9	4	5	0	9876	27	1	23	1	3	6	4	0	6	2	0	0	1
	Q	Gly	10	0	2	56	0	35	9865	4	2	3	1	4	1	0	3	4	2	0	1
	E	Gly	21	1	12	13	1	3	7	9955	1	0	1	1	0	1	0	21	3	0	5
	G	His	1	6	10	3	1	20	1	0	9932	0	1	3	0	2	3	1	1	1	4
	H	His	2	2	3	1	2	1	2	0	0	9872	9	2	12	7	0	1	7	0	1
	L	Leu	3	1	3	0	0	4	1	1	4	22	9947	2	45	13	3	1	3	4	-2
	K	Lys	2	37	25	6	0	12	7	2	2	4	1	9926	20	0	2	8	11	0	1
	M	Met	1	1	0	0	0	2	0	0	0	5	8	0	9874	1	0	1	2	0	4
	F	Phe	1	1	1	0	0	0	0	1	2	8	6	0	4	9946	0	2	1	3	20
	P	Pro	13	5	2	1	1	0	3	2	5	1	2	2	1	1	9926	12	4	0	2
	S	Ser	28	11	34	7	11	4	6	16	2	2	1	7	4	3	10	9840	34	5	2
	T	Thr	22	2	13	4	1	3	2	0	1	11	2	8	6	1	1	32	9871	0	2
	W	Tyr	0	2	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	9976	1
	Y	Tyr	1	0	3	0	3	0	1	0	4	1	1	0	0	21	0	1	1	1	9945
	V	Val	13	2	1	1	3	2	3	3	5	11	1	1	1	3	2	10	0	2	9903

Figure B2. Mutation probability matrix for the evolutionary 1 accepted point mutation per 100 amino acids. Thus, there is a 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

Mutation probability matrix is not symmetrical

Probabilities of AA_j replaced by AA_i

BIOS477/877 L9 - 16

16

PAM matrices

• Relatedness odds score

Odds ratio

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

$P(a \text{ was derived from } b) \text{ [nonrandom]}$
 $P(\text{random occurrence of } a)$

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

the probability

PAM matrices

Log odds matrix: $S(a,b)$ Dayhoff et al. (1978)

Figure B4. Log odds matrix for 250 PAMs. Elements are shown multiplied by 10. The neutral score is zero. A score of +10 means that the pair would be expected to occur only one tenth as frequently in related sequences as random chance would predict, and a score of +2 means that the pair would be expected to occur 1.6 times as frequently. The order of the amino acids has been arranged to illustrate the patterns in the mutation data.

Log odds matrix is symmetrical: $S(a,b) = S(b,a)$
 $10 \log_{10}\{R(a,b)\} = 10 \log_{10}\{R(b,a)\}$

BIOS477/877 L9 - 19

19

PAM matrices

➤ **PAM1 matrix** $M(a,b) = \lambda m(b) \times f(a,b) / \sum_a f(a,b)$

→ 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)

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)

e.g., M₂₅₀: 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)

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 PAM_n
 $M_n = (M_1)^n$ (e.g., PAM250 or M₂₅₀ = M₁²⁵⁰)

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	326

100[1 - S_b{f(b)M_n(b,b)}] ← PAM1

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

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

```

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   LOCNI|O25321 ( 53) VTAVLSLIIITIVGNLVLSVF 3
OAR2   LOCNI|O25322 ( 53) VTAVLSLIIITIVGNLVLSVF 3
OQ120  ( 22) ISLAVLEPIINVLVGGNCLVIKVF 22
OAR_DROME|P22270 ( 111) LTALVLSVIVLTIIGNILVLSVF 3
DOP2_DROME|O24563 ( 110) GLLAFLEFSPFATVFGNSLVILAVI 5
OAR_HBLVT|O25188 ( 54) CTAVLTLIIISTIVGNILVLSVF 6
OQ3128 ( 35) ISLLALFPLNLMVAGNLLVNVAVF 9
OQ3126 ( 35) ISLLALFPLNLMVAGNLLVNVAVF 9
A1AB_MESAU|P18841 ( 47) SVGLVLFAPILFAIVGNILVLSVA 5
A1AB_RAT|P15823 ( 47) SVGLVLFAPILFAIVGNILVLSVA 5
A1AB_HUMAN|P35368 ( 47) SVGLVLFAPILFAIVGNILVLSVA 5
OQ3127 ( 34) AATALLLAIILVTIIGNSLVIISVF 3
A1AA_ORYLA|O91175 ( 28) VLGMLFIFLFGVIGNILVLSVV 5
D2DI_XENLA|P24628 ( 30) YYAMLLTLLVVFVFGNVLVLSVA 6
D2DR_CRRAR|P52702 ( 36) YYATLLTLLIIVFVFGNVLVLSVA 5
OQ3810 ( 33) YYAVLTLLEFVIFVFGNVLVLSVA 5
OAR_BOHMO|O17232 ( 56) CTAILTMIISTVVGNIIVLSVF 7
SSR5_MOUSE|O08858 ( 39) LVFVLLVCTVGLGGNTLVIVVVL 4
    
```

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

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 ($p_i p_j$)

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

H₁: Homologous hypothesis (residues *i* and *j* are related)
H₀: 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 < +inf] = $\frac{\text{Prob}(\text{an event}|H_1)}{\text{Prob}(\text{an event}|H_0)}$

H₁: Hypothesis to be tested, H₀: 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 = $\log \left\{ \frac{\text{Likelihood of } H_1}{\text{Likelihood of } H_0} \right\}$

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

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

H₁: Hypothesis to be tested, H₀: 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