

# BIOS 477/877 Bioinformatics and Molecular Evolution

Instructor: Etsuko Moriyama  
(School of Biological Sciences)

Spring 2026    Lecture 10

BIOS477/877 L10 - 1

1

## Today's topics

- Amino Acid Substitution Matrix
  - Information Theory
- Statistical Significance of Alignment Scores
- Assignment 5

BIOS477/877 L10 - 2

2

## 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 [pairing by chance]

### ▪ General form

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

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

BIOS477/877 L10 - 3

3

## Relative entropy ( $H$ )

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

### ➤ Expected score ( $E$ )

$$E = \sum_i \sum_j p_i p_j S(i,j) \quad [p_i, p_j: \text{expected freq. of } AA_i, AA_j]$$

### ➤ Relative entropy ( $H$ )

$$H = \sum_i \sum_j q_{ij} \lambda S(i,j)$$

[ $q_{ij}$ : observed freq. of  $AA_i, AA_j$  pair]

Entropy in Information Theory:  
completely unrelated to "entropy"  
in thermodynamics!

- The average information per residue pair
- Summarizes the behavior of the scoring matrix
- Measures the ability of the matrix to discriminate related from unrelated (nonrandom from random matching)
  - $H = 0$  when target distribution equals to background distribution  
 [ If  $q_{ij} = p_i p_j \rightarrow S(i,j) = 1/\lambda \log_e(q_{ij}/p_i p_j) = 1/\lambda \log_e(1) = 0$  ]
  - $H$  increases when the two distributions become more distinguishable
- Can be used to compare scoring matrices

BIOS477/877 L10 - 4

4

## Introduction to information theory

- Information: a decrease in uncertainty (unpredictability, a degree of surprise)
- If you are asking questions to somebody...
  - if you can guess every answer correctly
    - there is no surprise
    - you cannot gain any new information
  - but if you have no idea what answer you get
    - every answer is a surprise
    - you gain a lot of information

See "Information Theory Primer" by Tom Schneider (on CANVAS)    BIOS477/877 L10 - 5

5

## How to quantify information

- Information: a decrease in uncertainty (unpredictability, a degree of surprise)

Device ⇒ A, A, A, A, A, A, A, A, ...

Prob(A) = 1  
Only one possible symbol

No surprise  
No information

BIOS477/877 L10 - 6

6

### How to quantify information

➤ **Information: a decrease in uncertainty**  
(unpredictability, a degree of surprise)

Device  $\Rightarrow$  A, B!, A, B!, A, B, A, B, ...

Prob(A) = 0.5  
Prob(B) = 0.5

Two possible symbols  
A little surprise  
A small amount of information

BIOS477/877 L10 - 7

7

### How to quantify information

➤ **Information: a decrease in uncertainty**  
(unpredictability, a degree of surprise)

Device  $\Rightarrow$  A, B!, C!!, A, B!, C!!, A, B, ...

Prob(A) = 0.33  
Prob(B) = 0.33  
Prob(C) = 0.33

Three possible symbols  
More surprise, More information

BIOS477/877 L10 - 8

8

### How to quantify information

➤ **Information: a decrease in uncertainty**  
(unpredictability, a degree of surprise)

Device  $\Rightarrow$  A, A, A, A, A, B!!!, A, A, ...

Prob(A) = 7/8  
Prob(B) = 1/8

Two possible symbols  
Big surprise! A lot more information...?  
But not much surprise in getting the symbol A's

BIOS477/877 L10 - 9

9

### How to quantify information

➤ **Information: a decrease in uncertainty**

- Surprising answers convey more information!
- If each symbol is equally likely, the amount of information increases with the number of different symbols.
- The amount of information, or surprise of an answer, is inversely proportional to its probability.

$I(p) = \log_2(1/p)$  or  $I(p) = -\log_2 p$   
I: information, p: probability

BIOS477/877 L10 - 10

10

### Unit of information: bit

➤ **Bits:** the unit for values converted to base 2 logarithms (**nats:** the unit if base e is used)

$I(p) = -\log_2 p$

- If an answer is highly unexpected (e.g.,  $p = 0.1$ ),  
→  $I(0.1) = -\log_2 0.1 = 3.3$  bits (more information)
- For a very much expected answer (e.g.,  $p = 0.9$ ),  
→  $I(0.9) = -\log_2 0.9 = 0.15$  bits (less information)
- If there is only one possible answer (symbol):  
→  $p = 1, I(p) =$  ?

BIOS477/877 L10 - 11

11

### Shannon entropy

➤ Information can be represented by a series of symbols each with a certain probability

- **Shannon Entropy:** the average information per symbol

$H = -\sum p_i(\log_2 p_i)$

- If all  $n$  symbols are equally possible ( $p_i = p = 1/n$ ),  
→  $H = -\sum p(\log_2 p) = -(np \times \log_2 p)$   
=  $-\log_2 p$ , since  $np = 1$   
=  $-\log_2(1/n)$ , since  $p = 1/n$   
=  $\log_2(n)$

e.g.,  $H(1) = \log_2(1) = 0$  bit,  $H(2) = \log_2(2) = 1$ ,  $H(4) = \log_2(2^2) = 2$   
 $H(n)$ :  $n$  is the number of equally possible letters

BIOS477/877 L10 - 13

13

## Shannon entropy

- Information can be represented by a series of symbols each with a certain probability:
  - Shannon Entropy:** the average information per symbol

$$H = -\sum p_i (\log_2 p_i)$$

[Examples]

- For a random DNA sequence: ATGC ( $p = 0.25$  for all)  
 $H = -4 \times 0.25 \times \log_2(0.25)$  or  $\log_2(4) = 2$  bits
- For a AT-rich DNA sequence:  $p_A = p_T = 0.45$  and  $p_C = p_G = 0.05$   
 $H = \{-2 \times 0.45 \times (\log_2 0.45)\} + \{-2 \times 0.05 \times (\log_2 0.05)\}$   
 $= \{-2 \times 0.45 \times (-1.15)\} + \{-2 \times 0.05 \times (-4.32)\} = 1.47$  bits

BIOS477/877 L10 - 14

14

## Relative entropy (H)

- Expected score (E)**  
 $E = \sum_i \sum_j p_i p_j S(i,j)$  [ $p_i, p_j$ : expected freq. of AA<sub>i</sub>, AA<sub>j</sub>]
  - Relative entropy (H)**  
 $H = \sum_i \sum_j q_{ij} \lambda S(i,j)$  [ $q_{ij}$ : observed freq. of AA<sub>i</sub>, AA<sub>j</sub> pair]
- Since  $S(i,j) = 1/\lambda \log_2(q_{ij}/p_i p_j)$  or  $1/\lambda \log_e(q_{ij}/p_i p_j)$ ,
- $$H = \sum_i \sum_j q_{ij} \log_2(q_{ij}/p_i p_j) \text{ or } \sum_i \sum_j q_{ij} \log_e(q_{ij}/p_i p_j)$$
- $$= \sum_i \sum_j \{q_{ij} \log_2(q_{ij}) - q_{ij} \log_2\{p_i p_j\}\}$$

Connection to Shannon entropy:  
 $H = -\sum p_i (\log_2 p_i)$

Note: Both expected score and relative entropy have their units in bit (or nat).

BIOS477/877 L10 - 15

15

## Comparing scoring matrices

- Relative entropy (H) of a scoring matrix**  
 $H = \sum_i \sum_j q_{ij} \lambda S(i,j)$  [ $q_{ij}$ : observed freq. of AA<sub>i</sub>, AA<sub>j</sub> pair]
- The **average information** per residue pair for a scoring matrix
- Summarizes the behavior of the scoring matrix
- Measures the ability of the matrix to discriminate related from unrelated (nonrandom from random matching)  
 $\rightarrow H = 0$  when target distribution equals to background distribution  
 [ If  $q_{ij} = p_i p_j \rightarrow S(i,j) = 1/\lambda \log_e(q_{ij}/p_i p_j) = 1/\lambda \log_e(1) = 0$  ]  
 $\rightarrow H$  increases when the two distributions become more distinguishable
- Can be used to compare scoring matrices

BIOS477/877 L10 - 16

16

## Comparing scoring matrices

- Relative entropy (H)**  
 $H(\text{PAM1}) = 4.17$  bits  
 $H(\text{PAM50}) = 2.00$   
 $H(\text{PAM120}) = 0.98$   
 $H(\text{PAM160}) = 0.70$   
 $H(\text{PAM250}) = 0.36$   
 from Altschul (1991)
- $H$  decreases with increasing PAM#
- $H$  increases with increasing BLOSUM#
- Higher BLOSUM is generated with sequences that are more similar to one another  $\rightarrow$  More amino acid pairs are used

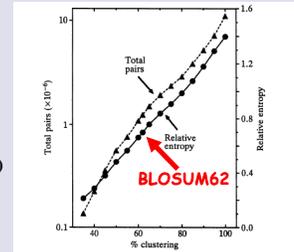


FIG. 1. Relationship between percentage clustering and total amino acid pair counts plotted on a logarithmic scale and relative entropy.

from Henikoff and Henikoff (1992)

BIOS477/877 L10 - 17

17

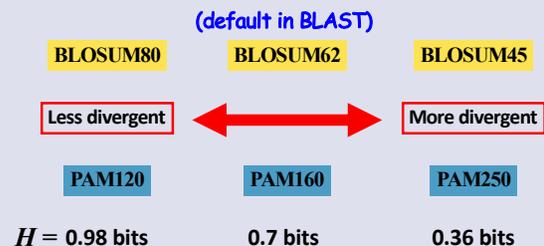
## Comparing scoring matrices

- Relative entropy (H) of a scoring matrix**  
 $H = \sum_i \sum_j q_{ij} \lambda S(i,j)$
- The average information per residue pair for a scoring matrix
- Decreases with increasing PAMn:  $H(\text{PAM1}) > H(\text{PAM120})$ 
  - All PAMn is extrapolated from PAM1
  - Higher PAMn is less specific, contains less information,  $\rightarrow$  has a lower  $H$
- Increases with increasing BLOSUMn:
  - $H(\text{BLOSUM45}) < H(\text{BLOSUM80})$
  - Higher BLOSUMn is generated using more data (fewer information is eliminated)  $\rightarrow$  has a higher  $H$   
 [e.g., BLOSUM100 is generated using the threshold 100%; only identical sequences are down-weighted for calculation]

BIOS477/877 L10 - 18

18

## BLOSUM and PAM matrices



BIOS477/877 L10 - 19

19

## BLOSUM62

```

# Matrix made by matblas from blosum62.iij
# * column uses minimum score
# BLOSUM Clustered Scoring Matrix in 1/2 Bit Units  $S(i,j) = 2\log_2(q_{ij}/e_{ij})$ 
# Blocks Database = /data/blocks_5.0/blocks.dat
# Cluster Percentage: >= 62
# Entropy = 0.6979, Expected = -0.5209
A R N D C Q E G H I L K M F P S T W Y V B E X *
A 4 -1 -2 -2 0 -1 -1 0 -2 -1 -1 -1 -1 -2 -1 1 0 -3 -2 0 -2 -1 0 -4
R -1 5 0 -2 -3 1 0 -2 0 -3 -2 2 -1 -3 -2 -1 -1 -3 -2 -3 -1 0 -1 -4
N -2 0 6 1 -3 0 0 0 1 -3 -3 0 -2 -3 -2 1 0 -4 -2 -3 3 0 -1 -4
D -2 2 1 6 -3 0 2 -1 -1 -3 -4 -1 -3 -3 -1 0 -1 -4 -3 -3 4 1 -1 -4
C 0 -3 -3 9 -3 -4 -3 -3 -1 -1 -3 -1 -2 -3 -1 -1 -2 -2 -1 -3 -3 -2 -4
Q -1 1 0 0 -3 5 2 -2 0 -3 -2 1 0 -3 -1 0 -1 -2 -1 -2 0 3 -1 -4
E -1 0 0 2 -4 2 5 -2 0 -3 -3 1 -2 -3 -1 0 -1 -3 -2 -2 1 4 -1 -4
G 0 -2 0 -1 -3 -2 -2 6 -2 -4 -4 -2 -3 -3 -2 0 -2 -2 -3 -3 -1 -2 -1 -4
H -2 0 1 -1 -3 0 0 -2 8 -3 -3 -1 -2 -1 -2 -1 -2 -2 2 -3 0 0 -1 -4
I -1 -3 -3 -3 -1 -3 -3 -4 -3 4 2 -3 1 0 -3 -2 -1 -3 -1 3 -3 -3 -1 -4
L -1 -2 -3 -4 -1 -2 -3 -4 -3 2 4 -2 2 0 -3 -2 -1 -2 -1 1 -4 -3 -1 -4
K -1 2 0 -1 -3 1 1 -2 -1 -3 -2 5 -1 -3 -1 0 -1 -3 -2 2 0 1 -1 -4
M -1 -1 -2 -3 -1 0 -2 -3 -2 1 2 1 5 0 -2 -1 -1 -1 -1 1 -3 -1 -1 -4
F -2 -3 -3 -3 -2 -3 -3 -3 -1 0 0 -3 0 6 -4 -2 -2 1 3 -1 -3 -3 -1 -4
P -1 -2 -2 -1 -3 -1 -1 -2 -2 -3 ...

```

Can be downloaded from [NCBI ftp site](https://www.ncbi.nlm.nih.gov/blast/matrix/)

BIOS477/877 L10 -20

20

## How to select a right scoring matrix

### Selecting the Right Similarity-Scoring Matrix

Pearson (2013)

**Relative entropy (H)**

Matrix	Gap penalty <sup>a</sup>	% Identity	Bits/position	Random alignment length	50-bit length
<i>SSEARCH version 36.3.6</i>					
BLOSUM50 <sup>b</sup>	10/2	25.3	0.21	160	238
BLOSUM62	11/1	28.9	0.40	86	125
VTML 140 <sup>c,d</sup>	13/2	23.9	0.25	139	200
VTML 140	10/1	28.4	0.44	82	114
VTML 120	11/1	32.1	0.54	62	93
VTML 80	10/1	40.5	0.74	47	68
VTML 40	13/1	64.7	1.92	18	26
VTML 20	15/2	86.1	3.30	11	15
VTML 10	16/2	90.9	3.87	9	13
<i>BLAST version 2.2.27+</i>					
BLOSUM50 <sup>b</sup>		29.4	0.39	85	128
BLOSUM62		29.6	0.41	82	122
BLOSUM80	10/1	32.0	0.48	69	104
PAM70	10/1	33.9	0.58	56	86
PAM30	9/1	45.9	0.90	34	56

- Default matrices (e.g., BLOSUM62) are good for identifying  $\geq 25\%$  identity.
- Deep scoring matrices (e.g., BLOSUM62, PAM250) require long sequence alignment to achieve significant scores (e.g., >50 bits).
- They are more likely to extend alignments outside of homologous region.

BIOS477/877 L10 -21

21

## How to select a right scoring matrix

### Selecting the Right Similarity-Scoring Matrix

Pearson (2013)

A	R	N	D	C	Q	E
A	7					
R	-7	8				
N	-6	-5	8			
D	-6	-12	-1	8		
C	-3	-7	-8	-14	12	
Q	-5	-2	-4	-4	-13	9
E	-5	-10	-5	-1	-14	-7

A	R	N	D	C	Q	E	
A	4						
R	-1	5					
N	-2	0	6				
D	-2	-2	1	6			
C	0	-3	-3	-3	9		
Q	-1	1	0	0	-3	5	
E	-1	0	0	2	-4	2	5

**Figure 3.5.2** Comparison of a "shallow" (VTML 20) and "deep" (BLOSUM62) scoring matrix. Both matrices are scaled in 1/2-bits. For the small part of the matrices shown here, the VTML20 matrix produces an average 2.80 half-bit identity score, and an average -0.59 nonidentity score (weighted by amino-acid abundance). In contrast, BLOSUM62 produces 1.86 for identities but only -0.06 for nonidentities. Thus, VTML20 targets shorter, higher-identity alignments, because it penalizes nonidentities much more strongly.

- Short alignments require shallow scoring matrices to be more significant.
- Shallower scoring matrices (e.g., PAM20) are more effective when searching over shorter evolutionary distances.

22

22

## BLOSUM62, RBLOSUM, and CorBLOSUM

CORRESPONDENCE

### BLOSUM62 miscalculations improve search performance

Styczynski et al.

Hess et al. BMC Bioinformatics (2010) 11:189  
DOI 10.1186/1471-2109-11-189-3

BMC Bioinformatics

RESEARCH ARTICLE

### Addressing inaccuracies in BLOSUM computation improves homology search performance

CorBLOSUM

Martin Hess<sup>1,2\*</sup>, Frank Keul<sup>1,3\*</sup>, Michael Goesele<sup>1</sup> and Kay Hamacher<sup>2</sup>

RESEARCH NOTE

### RBLOSUM performs better than CorBLOSUM with lesser error per query

Still the original BLOSUM62 is used...

Rengananayaki Govindarajan<sup>1</sup>, Bijl Christopher Leela and Achuthsankar S. Nair

Open Access

BIOS477/877 L10 -23

23

## More specific substitution matrices

RESEARCH ARTICLE

### PFASUM: a substitution matrix from Pfam structural alignments

Frank Keul<sup>1\*</sup>, Martin Hess<sup>1,2\*</sup>, Michael Goesele<sup>1</sup> and Kay Hamacher<sup>2</sup>

Open Access

BIOINFORMATICS

### Environment specific substitution tables improve membrane protein alignment

Jamie R. Hill<sup>1</sup>, Sebastian Kelm<sup>1</sup>, Jiye Shi<sup>2,3</sup> and Charlotte M. Deane<sup>1,\*</sup>

Published online 16 September 2021

Open Access

RESEARCH ARTICLE

### Evolutionary and functional lessons from human-specific amino acid substitution matrices

Tair Shauli<sup>1</sup>, Nadav Brandes<sup>2,3</sup> and Michal Linial<sup>2,3\*</sup>

Open Access

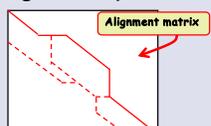
- See more on CANVAS (References)
- More substitution matrices reviewed in Trivedi & Nagarajaram (2020)
- EModelDB: a database of substitution matrix

BIOS477/877 L10 -24

24

## Pairwise alignment summary

- Alignment score depends on:
  - Scoring matrix (match, mismatch, Ts/Tv, BLOSUM, PAM, etc.)
  - Gap penalty
  - Alignment method (e.g., global or local)
- Alignment scores cannot be compared directly
  - if the scoring systems used are different
  - if sequences compared are different (e.g., longer alignments tend to have higher scores)
- Alignment scores are used:
  - for searching optimal alignments
  - from the alignment matrix
  - for a given pair of sequences
  - based on a given scoring system



Alignment matrix

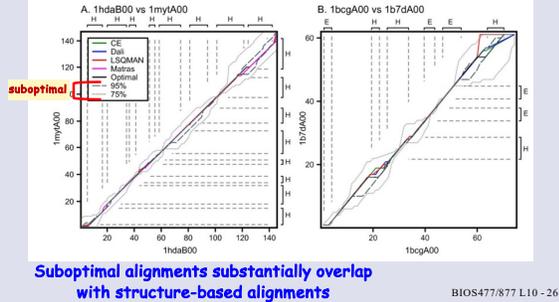
BIOS477/877 L10 -25

25

## Optimal vs. suboptimal alignment

Improving pairwise sequence alignment accuracy using near-optimal protein sequence alignments

Michael L. Sierk<sup>1</sup>, Michael E. Smoot<sup>2</sup>, Ellen J. Bass<sup>3</sup>, William R. Pearson<sup>4</sup>



BIOS477/877 L10 - 26

26

## Pairwise alignment summary

(continued)

- **Optimal alignments** and **biologically meaningful alignments** may not be the same
- Depending on the scoring system, **unreasonable alignments can become optimal**
  - We need to choose a better (**biologically reasonable**) scoring system considering:
    - level of divergence (scoring matrices)
    - gap penalty (affine, etc.)
    - algorithm (local, global, or semi-global)
  - Manual adjustment may be necessary
  - Need to test **statistical significance of the alignment** (is the alignment possible just by chance?)

BIOS477/877 L10 - 27

27

## Hypothesis testing

- For general tests
  - Two hypotheses
    - **Null-hypothesis**  
 $H_0$ : The previous (original) belief is true
    - **Alternative hypothesis**  
 $H_1$ : The previous (original) belief is false; The new theory is true
  - **S**: Test statistic
  - **Significance level** is chosen *a priori* (e.g., 0.05)
  - **P-value**:  $P(S|H_0 \text{ is true})$   
→ Probability of getting S if  $H_0$  is true
  - If  $P < \text{Significance level}$ , reject  $H_0$

BIOS477/877 L10 - 28

28

## How to calculate P-values

- **P-value**:  $P(S|H_0 \text{ is true})$ 
    - Calculated from the test statistic **S**  
→ Need to know the **probability distribution of the test statistic S under the null hypothesis,  $H_0$**
- Central Limit Theorem:**  
If the sample size is large enough, the sampling distribution of the mean of any independent, random variables will be **normal or nearly normal**.
- [Example experiment]**

  - Toss a coin 1000 times
  - Count the number of heads
  - Repeat 1000 times

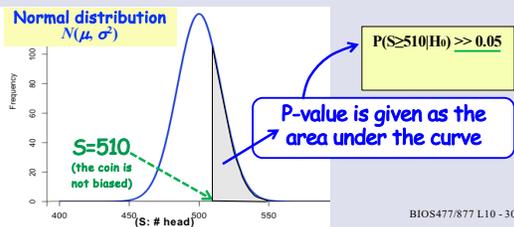
(Expect to see 500 heads/experiment)
- 

BIOS477/877 L10 - 29

29

## How to calculate P-values

- **P-value**:  $P(S|H_0 \text{ is true})$ 
  - Calculated from the test statistic **S**  
→ Need to know the **probability distribution of the test statistic S under the null hypothesis,  $H_0$**

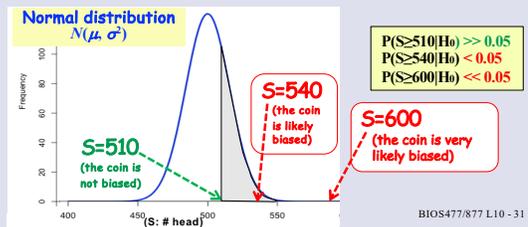


BIOS477/877 L10 - 30

30

## How to calculate P-values

- **P-value**:  $P(S|H_0 \text{ is true})$ 
  - Calculated from the test statistic **S**  
→ Need to know the **probability distribution of the test statistic S under the null hypothesis,  $H_0$**



BIOS477/877 L10 - 31

31

## How to test alignment scores

- Hypothesis testing **for sequence alignment**
  - Two hypotheses
    - Null-hypothesis  
 $H_0$ : Two sequences are not related (random)
    - Alternative hypothesis  
 $H_1$ : Two sequences are related
  - Test statistic: **alignment score ( $S$ )**
  - Significance level is chosen a priori (e.g., 0.05)
  - P-value:  $P(S|H_0 \text{ is true})$   
→ Probability of getting the alignment score  $S$  even if unrelated randomly selected sequences are aligned
  - If  $P < \text{Significance level}$ , reject  $H_0$   
(The score cannot be obtained just by aligning unrelated sequences)

BIOS477/877 L10 - 32

32

## Extreme value distribution

- P-value:  $P(S|H_0 \text{ is true})$ 
  - Calculated from the test statistic  $S$   
→ Need to know **the probability distribution of the test statistic  $S$  under the null hypothesis,  $H_0$**

Distribution of alignment scores follow

**Extreme Value Distribution (Gumbel distribution)**

The probability distribution of highest values in an experiment (e.g., optimal alignment scores)

BIOS477/877 L10 - 33

33

## Significance of alignment scores

EVD (Gumbel distribution):  
 $P(S < x) = \exp[-e^{-(x-\mu)/\beta}]$   
 $P(S \geq x) = 1 - \exp[-e^{-(x-\mu)/\beta}]$   
 $\beta$ : scale parameter  
 $\mu$ : location parameter

$P(S \geq x_1 | H_0) \gg 0.05$   
 $P(S \geq x_2 | H_0) < 0.05$   
 $P(S \geq x_3 | H_0) \ll 0.05$

Not significant      Significant!      Highly Significant!

BIOS477/877 L10 - 34

34

## Significance of alignment scores

- $P(S \geq x | H_0)$ : probability of getting the alignment score  $S \geq x$  by chance

**Karlin-Altschul equation** (Karlin and Altschul 1990)

**$P(S \geq x) = 1 - \exp[-Kmn e^{-\lambda x}] \approx Kmn e^{-\lambda x}$**

- $K$  and  $\lambda$ : parameters based on a given scoring matrix and the amino acid composition of the sequences
- $m$  and  $n$ : lengths of sequences aligned

→ Solved for ungapped local alignments  
 → Can be applied for gapped local alignments

➤ **E-value** =  $P(S \geq x | H_0) \times N$       **E-value  $\neq$  P-value**

$N$ : the number of sequences in the dataset  
**Expected number** of randomly selected sequences in the dataset that have alignment scores  $\geq x$

BIOS477/877 L10 - 35

35

## Estimation of $K$ , $\lambda$ , and P-value

- Estimation of  $K$  and  $\lambda$  from empirical distribution of random alignment scores (used in LALIGN and PRSS)
  - 1) The second sequence is shuffled many times  
→ Simulates random sequences
  - 2) Smith-Waterman local alignment score is calculated from each alignment:  $P(S \geq x | H_0)$
  - 3) The distribution is fitted to an extreme value distribution to obtain estimates of  $K$  and  $\lambda$
  - 4) P-value is estimated based on the  $K$  and  $\lambda$ , and the original alignment score  $x$ :  $P(S \geq x) \approx Kmn e^{-\lambda x}$

BIOS477/877 L10 - 36

36

## Simulation of random score distribution (example)

[Input sequences]  
 • RECA\_ECOLI (P0A7G6; 353 amino acids)  
 • RAD51\_YEAST (P25454; 400 amino acids)

Smith-Waterman local alignment score = 293  
 (BLOSUM50, gap opening: -10, gap extension: -1)

```

RECA_ECOLI      3 IDENK-QALAAALQITRQKQKQIMLGRSRMVEVITITGSLRDTA 51
RAD51_YEAST    124 ISEARADKLLNAAALVPMGPTAADPIM--RHSLELCTGSRKLDL 170
RECA_ECOLI     52 LQAGLQPMGRVETIYDPSGGKTYL-----PLQV-IAAQRGRTCAFT 94
RAD51_YEAST    171 LG-GDVYTGSIETLPGFPRFKSQLCHTLAVTCQPLDGGGRK-CLYI 218
RECA_ECOLI     95 DAEHALDP----YARLEQVID----NLLGSPQTPQDALEICDALAR- 135
RAD51_YEAST    219 DTEGTFRFRVRLVSIAGRFGLDFOALNNVAVANAYNADRLQLLDAAGM 268
RECA_ECOLI    136 --SGAVVIVVDSVAALTFPAEI--EETIGDSBMLAAMHGMAMRRLAG 181
RAD51_YEAST    269 MERRKFLIVDSVVALY-RITDFSGRGLSARQMLAKFM--RALSRLA- 314
RECA_ECOLI    182 NLRQNTLLIFINIRKMI--GVMPG-NDEPTTGGNALEFYAVRVDLR 228
RAD51_YEAST    315 --DQYFAYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVY 358
RECA_ECOLI    229 IGAVKREHVVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVY 269
RAD51_YEAST    359 --GFRGK--GQGLKLVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVY 392
  
```

Using [EMBOSS WATER @ EBL](http://EMBOSS.WATER@EBL)      BIOS477/877 L10 - 37

37

