

# Alignment IV

## BLOSUM Matrices

# BLOSUM matrices

- Blocks Substitution Matrix. Scores for each position are obtained from frequencies of substitutions in blocks of local alignments of protein sequences [Henikoff & Henikoff92].
- For example BLOSUM62 is derived from sequence alignments with no more than 62% identity.

# BLOSUM Scoring Matrices

- BLOck SUBstitution Matrix
- Based on comparisons of blocks of sequences derived from the Blocks database
- The Blocks database contains multiply aligned ungapped segments corresponding to the most highly conserved regions of proteins (local alignment versus global alignment)
- BLOSUM matrices are derived from blocks whose alignment corresponds to the BLOSUM-,matrix number

# Conserved blocks in alignments

AABCD A . . . BBCDA  
DABCD A . A . BBCBB  
BBBCDABA . BCCAA  
AAACD A C . DCBCDB  
CCBADAB . DBBDCC  
AAACAA . . . BBCCC

# Constructing BLOSUM $r$

- To avoid bias in favor of a certain protein, first eliminate sequences that are more than  $r\%$  identical
- The elimination is done by either
  - removing sequences from the block, or
  - finding a cluster of similar sequences and replacing it by a new sequence that represents the cluster.
- BLOSUM  $r$  is the matrix built from blocks with no more the  $r\%$  of similarity
  - E.g., BLOSUM62 is the matrix built using sequences with no more than 62% similarity.
  - Note: BLOSUM 62 is the default matrix for protein BLAST

# Collecting substitution statistics

1. Count amino acids pairs in each column; e.g.,
  - 6 AA pairs, 4 AB pairs, 4 AC, 1 BC, 0 BB, 0 CC.
  - Total = 6+4+4+1=15
2. Normalize results to obtain probabilities ( $p_x$ 's and  $q_{xy}$ 's)
3. Compute log-odds score matrix from probabilities:  
$$s(X,Y) = \log (q_{xy} / (p_x p_y))$$

A  
A  
B  
A  
C  
A

# Computing probabilities

Sum the scores for each columns across columns:

$$c_{ij} = \sum_k c_{ij}^{(k)}$$

Normalize the pair frequencies so they will sum to 1:

$$T = \sum_{i \geq j} c_{ij} = w \frac{n(n-1)}{2}$$

where  $w$  = number of columns  
 $n$  = number of sequences

$$q_{ij} = \frac{c_{ij}}{T}$$

# Computing probabilities

Calculate the expected probability of occurrence of the  $i$ th residue in an  $(i,j)$  pair:

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

The desired denominator is the expected frequency for each pair (assuming independence):

$$e_{ii} = p_i^2$$

$$e_{ij} = 2p_i p_j \quad (i \neq j)$$



# Computing probabilities

Each entry for  $(i,j)$  in the log odds matrix is then equal to  $q_{ij}/e_{ij}$

Log odds ratio:  $s_{ij} = \log_2 \frac{q_{ij}}{e_{ij}}$

Value stored for BLOSUM =  $2 s_{ij}$ , rounded to nearest integer (“half bit” units)

# Example

Matrix of  $c_{ij}$  values:

		A	I	L	S	T	V
sequence 1	A A I	A	I	L	S	T	V
sequence 2	S A L						
sequence 3	T A L						
sequence 4	T A V						
sequence 5	A A L						

  

	A	I	L	S	T	V
A	1+10					
I		0				
L		3	3			
S	2		0			
T	4			2	1	
V		1	3			0

$$T = \sum_{i \geq j} c_{ij} = 3 \left[ \frac{(5)(4)}{2} \right] = 30$$

# Example

Matrix of  $q_{ij}$  values:

	A	I	L	S	T	V
A	$11/30$					
I		0				
L		$3/30$	$3/30$			
S	$2/30$		0	0		
T	$4/30$			$2/30$	$1/30$	
V		$1/30$	$3/30$			0

=

	A	I	L	S	T	V
A	0.366					
I	0	0				
L	0	0.1	0.1			
S	0.066	0	0			
T	0.133	0	0	0.066	0.033	
V	0	0.033	0.1	0	0	0

Vector of  $p_i$  values:

$$p_A = \left(11 + \frac{6}{2}\right) / 30 = 14/30 = 0.46\bar{6}$$

$$p_I = \left(0 + \frac{4}{2}\right) / 30 = 2/30 = 0.06\bar{6}$$

$$p_L = \left(3 + \frac{6}{2}\right) / 30 = 6/30 = 0.2$$

$$p_S = \left(0 + \frac{4}{2}\right) / 30 = 2/30 = 0.06\bar{6}$$

$$p_T = \left(1 + \frac{6}{2}\right) / 30 = 4/30 = 0.13\bar{3}$$

$$p_V = \left(0 + \frac{4}{2}\right) / 30 = 2/30 = 0.06\bar{6}$$

# Example

Matrix of  $e_{ij}$  values:

	A	I	L	S	T	V
A	$(\frac{14}{30})^2$					
I	$2(\frac{14}{30})(\frac{2}{30})$	$(\frac{2}{30})^2$				
L	$2(\frac{14}{30})(\frac{6}{30})$	$2(\frac{2}{30})(\frac{6}{30})$	$(\frac{6}{30})^2$			
S	$2(\frac{14}{30})(\frac{2}{30})$	$2(\frac{2}{30})(\frac{2}{30})$	$2(\frac{6}{30})(\frac{2}{30})$	$(\frac{2}{30})^2$		
T	$2(\frac{14}{30})(\frac{4}{30})$	$2(\frac{2}{30})(\frac{4}{30})$	$2(\frac{6}{30})(\frac{4}{30})$	$2(\frac{2}{30})(\frac{4}{30})$	$(\frac{4}{30})^2$	
V	$2(\frac{14}{30})(\frac{2}{30})$	$2(\frac{2}{30})(\frac{2}{30})$	$2(\frac{6}{30})(\frac{2}{30})$	$2(\frac{2}{30})(\frac{2}{30})$	$2(\frac{4}{30})(\frac{2}{30})$	$(\frac{2}{30})^2$

Log odds ratio:

$$\text{e.g., } s_{AA} = \log_2 \frac{0.36\bar{6}}{\left(\frac{14}{30}\right)^2} = \log_2 1.6837 = 0.7516$$

BLOSUM value for AA =  $\text{round}(2 \cdot 0.7516) = 2$

Full matrix:

	A	I	L	S	T	V
A	2					
I	?	?				
L	?	4	3			
S	0	?	?	?		
T	0	?	?	4	2	
V	?	4	4	?	?	?

Note: undefined values result from unobserved pairs (would ordinarily not happen with real data)

# Comparison

- PAM is based on an evolutionary model using phylogenetic trees
- BLOSUM assumes no evolutionary model, but rather conserved "blocks" of proteins

