

Is the similarity significant, or could it be due to chance?

Even if two proteins are unrelated, we would expect some similarity simply by chance.

Is the alignment score significantly higher than random?

Align random permutations of the sequences, and find the mean and standard deviation of the resulting distribution.

The z-score reflects the significance of a global similarity score.

$$z\text{-score} = \frac{\text{score} - \text{mean}}{\text{standard deviation}}$$

Larger values imply greater significance.

e-values and p-values

The expected number of HSPs with a score of at least S is given by the formula:

$$E = Kmn e^{-\lambda S}$$

Doubling the length of the query sequence (m) or the size of the database (n) should double the number of HSPs.

To obtain score $2x$, score x must be obtained twice in a row. So one expects E to decrease exponentially with score.

The probability of observing a score $\geq S$ is:

$$1 - \exp(-Kmn e^{-\lambda S})$$

This is the p-value.

BLAST

Basic Local Alignment Search Tool

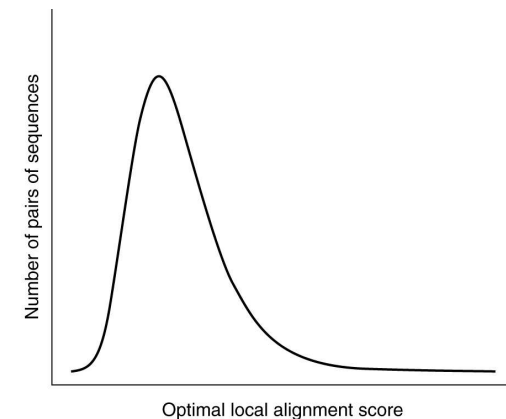
Less accurate than Smith-Waterman, but over 50 times faster.

1. Find ungapped matches of a small fixed length, w , that score at least T .
2. Extend matches in both directions in an attempt to find an alignment with a score exceeding S .

Segment pairs whose scores cannot be improved by extending or trimming are called high scoring pairs (HSPs).

Typical values for w are 3 when aligning proteins and 11 when aligning nucleic acids.

Extreme value distribution



FASTA

k-tuples, strings of length k.

k = 1 - 2 for proteins and 4-6 for nucleic acids.

Construct a look-up table with all k-tuples in the database.

Look up all k-tuples from the query string and mark matching database k-tuples. Sort matches by the difference in their indices (i-j).

Nearby matches on the same diagonal are joined to form an ungapped local alignment region.

Join nearby high scoring regions on different diagonals.

For the best regions, perform dynamic programming in a window around the region.

Possible substitution matrices for DNA

	A	C	G	T
A	2	-1	-1	-1
C	-1	2	-1	-1
G	-1	-1	2	-1
T	-1	-1	-1	2

	A	C	G	T
A	2	-2	-1	-2
C	-2	2	-2	-1
G	-1	-2	2	-2
T	-2	-1	-2	2

More realistic similarity measures

Not all substitutions are equally likely.

- A transition between two purines (A, G) or between two pyrimidines (C, T/U) is more common than a purine-pyrimidine transversion.
- Replacement of one amino acid residue by another with similar size or physiochemical properties is more common than replacement by a dissimilar amino acid residue.

Insertion/deletion of N contiguous amino acid residues or nucleotides is more likely than N independent insertion/deletion events.

Thus, we should have different penalties for opening gap and for extending a gap.

Relative likelihood and alignment score

Match model (M):

Sequences assumed to be dependent. Residues x_i and y_i at position i in the alignment occur together with probability $p_{x_i y_i}$.

Random model (R):

Sequences assumed to be independent. Residues x_i and y_i at position i in the alignment occur together with probability $q_{x_i} q_{y_i}$.

We can score an alignment using the log of the relative likelihood:

$$S = \log \left(\frac{Pr(x, y|M)}{Pr(x, y|R)} \right) = \log \frac{p_{x_1 y_1} p_{x_2 y_2} \cdots p_{x_n y_n}}{q_{x_1} q_{y_1} q_{x_2} q_{y_2} \cdots q_{x_n} q_{y_n}}$$
$$= \sum_{i=1}^n \log \left(\frac{p_{x_i y_i}}{q_{x_i} q_{y_i}} \right) = \sum_{i=1}^n s(x_i, y_i)$$

Percent accepted mutations

Expresses scores as log-odds values.

Score of mutation a-b is

$$\log \frac{\text{observed } a\text{-}b \text{ mutation rate}}{\text{mutation rate expected from amino acid frequencies}}$$

Frequencies of substitutions of each pair of amino acid residues, extracted from alignments of closely related proteins.

PAM1 reflects the amount of evolutionary change that yields an average of one mutation per 100 amino acids.

Can assume that no position has changed more than once.

Correct for different amino acid abundances.

PAM250

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	2																			
R	-2	6																		
N	0	0	2																	
D	0	-1	2	4																
C	-2	-4	-4	-5	4															
Q	0	1	1	2	-5	4														
E	0	-1	1	3	-5	2	4													
G	1	-3	0	1	-3	-1	0	5												
H	-1	2	2	1	-3	3	1	-2	6											
I	-1	-2	-2	-2	-2	-2	-3	-2	5											
L	-2	-3	-3	-4	-6	-2	-3	-4	-2	2	6									
K	-1	3	1	0	-5	1	0	-2	0	-2	-3	5								
M	-1	0	-2	-3	-5	-1	-2	-3	-2	2	4	0	6							
F	-4	-4	-4	-6	-4	-5	-5	-5	-2	1	2	-5	0	9						
P	1	0	-1	-1	-3	0	-1	-1	0	-2	-3	-1	-2	-5	6					
S	1	0	1	0	0	-1	0	1	-1	-1	-3	0	-2	-3	1	3				
T	1	-1	0	0	-2	-1	0	0	-1	0	-2	0	-1	-2	0	1	3			
W	-6	2	-4	-7	-8	-5	-7	-7	-3	-5	-2	-3	-4	0	-6	-2	-5	17		
Y	-3	-4	-2	-4	0	-4	-4	-5	0	-1	-1	-4	-2	7	-5	-3	-3	0	10	
V	0	-2	-2	-2	-2	-2	-2	-1	-2	4	2	-2	2	-1	-1	-1	0	-6	-2	4

PAM substitution matrices

Extrapolate to a family of PAM k matrices by multiplying the PAM1 matrix by itself k times.

Different PAM matrices are more suitable when comparing sequences that have diverged by different amounts.

The PAM250 matrix is commonly used.

250 mutations per 100 amino acids.

Sequences still 20% identical:

- some positions change many times, while others don't change at all.
- some positions change one or more times, then revert back to the original amino acid residue.

BLOSUM substitution matrices

Based on large collection of multiple alignments of similar ungapped segments.

$$\text{score}_{ab} = \log \frac{\text{observed relative frequency of aligned pairs } ab}{\text{expected probability of pair } ab}$$

Pairs are only counted between segments that are more than $x\%$ identical.

Different values of x give different BLOSUM matrices.

The BLOSUM62 matrix is commonly used.

BLOSUM62

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
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													
G	0	-2	0	-1	-3	-2	-2	6												
H	-2	0	1	-1	-3	0	0	-2	8											
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
S	1	-1	1	0	-1	0	0	-1	-2	-2	0	-1	-2	-1	4					
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-2	-1	1	5				
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

Graham Kemp, Chalmers University of Technology

Which substitution matrix should I use?

Use a matrix that corresponds to the evolutionary distance between the proteins being compared (usually not known!).

Low PAM matrices are good for finding short, strong similarities.

High PAM matrices are good for finding long, weak similarities.

BLOSUM matrices have been found to perform better for detecting weak homologies than the extrapolated PAM matrices.

Graham Kemp, Chalmers University of Technology