

or why biologists need computers

http://www.bioinformatics.uni-muenster.de/teaching/courses-2015/bioinf1/index.hb

1

Prof. Dr. Wojciech Makałowski Institute of Bioinformatics

INTRODUCTION TO SEQUENCE ANALYSIS

dot plots, alignments, and similarity searches



2





EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THISISANANCESTRALSEQUENCE THISISANMNCESTRALSEQUENCE THISISANMNCESTRAWSEQUENCE THISISANMPCESTRAWSEQUENCE



EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THISISCNMPESTRAWSEQUENCE

Gene duplication or speciation!

THISISCNMPESTRAWSEQUENCE

7

EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THISISCNMPESTRAWSEQUENCE THISISCOMPEETRAWSEQUENCE

THISISCNMPESTRAWSEQUENCE THISISNMPERSXTRASEQUENCE

Please note deletion of "C" and "W"

compensated by insertion of "R" and "X"

9

EVOLUTIONARY BASIS OF

SEQUENCE ANALYSES

THISISCOMPLETLNAWSEQUENCE

EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

8

THISISCOMPEETLAWSEQUENCE

Please note insertion of "C"

THISISCNMPEEXTRASEQUENCE

10

EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THISISCOMPLETLNAWSEQUENCE

THISISCSUPEEXTRASEQUENCE

THISISCSMPEEXTRASEQUENCE







ALGORITHM

A step-by-step problemsolving procedure, especially an established, recursive computational procedure for solving a problem in a finite number of steps.





19

20

DOT MATRIX PLOTS

- Sensitive qualitative indicators of similarity
- ✤ Better than alignments in some ways
 - · rearrangements
 - repeated sequences
- Rely on visual perception (not quantitative)
- Useful for RNA structure

21

DOT MATRIX PLOTS

- Simplest method put a dot wherever sequences are identical
- A little better use a scoring table, put a dot wherever the residues have better than a certain score (especially useful for amino acid sequence comparison)
- Or, put a dot wherever you get at least n matches in a row (identity matching, compare/word)
- +> Even better filter the plot

22

WINDOWED SCORES ALGORITHM

- 1. calculate a score within a window of a given size, for example six
- 2. plot a point if score is over a threshold (stringency), for example 70%
- 3. move the window over a given step, for example one
- 4. repeat step one to three till the end of sequence

WINDOWED SCORES EXAMPLE

Let's compare two nucleotide sequences

ACCTTGTCCTCTTTGCCC ACGTTGACCTGTAACCTC

using following parameters: window size = 9, step = 3, threshold = 4











А

A C C T T G T C C T C T T G C C C

















DOT PLOT EXAMPLES -REARRANGEMENTS deletion duplication Low 40 inversion

51

THISISANANCEST-R--ALSEQUENCE
THISISCOMP-LETELYNEWSEQUENCE
THISISSU-PEREXTRA-SEQUENCETHISISCOMP-LETELYNEWSEQUENCE
THISISSU-PEREXTRA-SEQUENCE6162ALIGNMENT PROBLEMALIGNMENT PROBLEMTHISISSOARCES-T-R--ALSEQUENCE
THISISSU-PEREXT-R--A-SEQUENCETHISISCOMP-LE-TELYNEWSEQUENCE
THISISSU-PEREXT-R--A-SEQUENCEFHISISSOARCES-T-R--ALSEQUENCE
THISISSU-PEREXT-R--A-SEQUENCETHISISCOMP-LE-TELYNEWSEQUENCE
THISISSU-PEREXT-R--A-SEQUENCE

63

ALIGNMENT PROBLEM

64

ALIGNMENT PROBLEM

ALIGNMENT

- Any assignment of correspondences that preserves the order of residues within the sequence is an alignment
- \cdot It is <u>the</u> basic tool of bioinformatics
- Computational challenge introduction of insertions and deletions (gaps) that correspond to evolutionary events
- · ★ We must define criteria so that an algorithm can choose the <u>best</u> alignment

ALIGNMENT AN EXAMPLE

Let's compare two strings gctgaacg and ctataatc

an uninformative alignment

an alignment without gaps

an alignment with gaps gctga-a--cg --ct-ataatc

another alignment with gaps gctg-aa-cg -ctataa-tc

SCORING SCHEMES A scoring system must account for residue substitution, and insertions or deletions (indels) Indels (gaps) will have scores that depend on their length For nucleic acid sequences, it is common to use a simple scheme for substitutions, e.g. +1 for a match, -1 for a mismatch More realistic would be to take into account nucleotide frequencies (sequence composition) and fact that transitions are more frequent than transversions

67

68

GAP SCORING SYSTEMS

- •> non-affine model each gap position treated the same, e.g. match = 4, mismatch = -3, gap -4
- * affine model first gap position penalized more than others, e.g. match = 4, mismatch = -3, gap opening = -8, gap = -4

GAP SCORING AN EXAMPLE

non-affine gapping score - the second alignment is "better"

69

70

GAP SCORING AN EXAMPLE

affine gapping score - the first alignment is "better"

GGTGCCAC-TCCAC----CTG AGTGCCACCCCCAATGCCGCTG -3 4 4 4 4 4 4 4 4 - 7

GGTGCCAC-TCCA---C--CTG AGTGCCACCCCCAATGCCGCTG -3 4 4 4 4 4 4 4 4 4 2-3 4 4 4 -12 -4 4 4 -12 -4 4 4 4 4 = 2

GAP SCORING AN EXAMPLE

Equivalent alignments

GGTGCCAC-TCCA---C--CTG AGTGCCACCCCCAATGCCGCTG 3 4 4 4 4 4 4 4 4 4 123 4 4 4 12 4 4 4 12 4 4 4 4 4 = 2

GGTGCCACT-CCA---C--CTG AGTGCCACCCCCAATGCCGCTG -3 4 4 4 4 4 4 4 -3 -12 4 4 4 -12 -4 4 -12 -4 4 4 4 = 2

AMINO ACID SCORING SYSTEMS

* more complicated than nucleotide matrices

- ✤ first, we can align two homologous protein sequences and count the number of any particular substitution, for instance Serine to Threonine
- * a likely change should score higher than a rare one
- we have to take into account that several the same position mutated several times after sequence divergence - this could bias statistics

73

APPROXIMATE RELATION BETWEEN PAM AND SEQUENCE IDENTITY

PAM	0	30	80	110	200	250
AA sequence identity (%)	100	75	50	60	25	20

PAM matrix is expressed as log-odds values multiplied by 10 simply to avoid decimal points

75

PAM MATRIX CALCULATION

74

AMINO ACID SCORING

• to avoid this problem one can compare very similar sequences so one can assume that no

Margret Dayhoff introduced the PAM system

identical residues

identical residues

P 1 PAM - two sequence have 99%

IO PAM - two sequence have 90%

position has changed more than once

(Percent of Accepted Mutations)

score of substitution i $\langle - \rangle j = \log$

SYSTEMS

observed i <-> j mutation rate

mutation rate expected from amino acids frequencies

For instance, a value 2 implies that in related sequences the mutation would be expected to occur 1.6 times more frequently than random.

The calculation: The matrix entry 2 corresponds to the actual value 0.2 because of the scaling. The value 0.2 is \log_{10} of the relative expectation value of the mutation. Therefore, the expectation value is $10^{0.2} = 1.6$

76

AMINO ACID MATRICES

- Problem with PAM schema lies in that the high number matrices are extrapolated from closely related sequences
- Henikoffs developed the family of BLOSUM matrices based on the BLOCKS database of aligned protein sequences, hence the name BLOcks SUbstitution Matrix
- ★ observed substitution frequencies taken from conserved regions of proteins (blocks), not the whole proteins as in case of Dayhoff's work
- * two avoid overweighting closely related sequences, the Hennikoffs replaced groups of proteins that have sequence identities higher than a threshold by either a single representative or a weighted average, e.g. for the commonly used BLOSUM62 matrix the threshold is 62%

· ➢ NOTE reversed numbering of PAM and BLOSUM matrices

80

81

SCORING RECOMMENDATIONS

• ⊱ nucleotide sequence comparison

• ★ match +10, mismatch -3, gap opening -50, gap extension -5

• amino acid sequence comparison

- For general use (e.g. unknown sequence similarity) - BLOSUM62
- for diverged proteins PAM250 or BLOSUM30
- ·⊱ for similar sequences PAM15 or BLOSUM80

Incrementally extend

Remember the best

sub-path leading to each point on the

+1

-1

-1

the path

lattice

Match:

Gap:

Mismatch:

SEQUENCE SIMILARITY SEARCH

104

BASICS OF DATABASE SEARCH

- $\cdot \succ \quad \text{Database searching is fundamentally different from alignment}$
- ★ The goal is to find homologous sequences (often more than one), not to establish the correct one-to-one mapping of particular residues
- $\cdot \succ \;$ Usually, this is a necessary first step to making an information map between two sequences
- · → Database searching programs were originally thought of as approximations to dynamic programming alignments

105

BASICS OF DATABASE SEARCH

basic terminology:

query - sequence to be used for the database search

subject - sequence found in the database that meets some similarity criteria

hit - local alignment between query and subject

106

BASICS OF DATABASE SEARCH

Through the influence of BLAST and FASTA, database searching programs have converged to a basic format

- a. a graphical depiction of the results
- b. a list of top scoring sequences from the databases
- **c**. a series of alignments for some of the top scoring sequences

Related sequences have "diagonals" with high similarity

BLAST

Basic Local Alignment Search Tool

References: Altschul, S.F., Gish, W., Miller, W., Myers, E.W. & Lipman, D.J. (1990) "Basic local alignment search tool." J. Mol. Biol. 215;403-410. Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W. & Lipman, D.J. (1997) "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs." Nucleic Acids Res. 25:3389-3402

109

NUCLEOTIDE BLAST ALGORITHM

- 1. Break down query sequence into overlapping words.
- 2.Scan databases for exact matches of size W (BLASTn) or 110110 pattern (MegaBlast).
- **3.**Try to extend the word matches into the complete maximal scoring pair (MSP). Significance is easily calculated from Karlin-Altschul equation.
- 4. Perform local dynamic programming alignment around MSP regions

110

BLAST - extend word matches

Most expensive step in BLAST algorithm

Extend to end of high scoring segment pair, or HSP. HSPs approximate maximal segment pairs or MSPs. They are only approximate because extension does not continue until running score reaches zero - drop off value concept.

After initial hit was found BLAST tries so called extension - an alignment is extended until the maximum value of the score drops by x, hence name x dropoff value

112

PROTEIN BLAST ALGORITHM

* Break down query sequence into overlapping words and create a lookaup table.

111

- $\cdot \\ \bullet$ For each word, determine a neighborhood of words that, if found in another sequence, would likely to be part of a significant maximum scoring pair (MSP).
- \cdot Scan databases for neighborhood words.
- * If two words are found on the same diagonal within a specified distance, try to extend the word matches into the complete MSP. Significance is (relatively) easy calculated from Karlin-Altschul equation.
- Perform local dynamic programming alignment around MSP regions
- * first step of BLASTp is controlled by three parameters and a score matrix
- ↔ w word length (k-tuple in FASTA terminology); default value is 3 (lowest possible is 2); two words on the same diagonal are required
- f score threshold; unlike FASTA BLAST allows mismatches at this step but overall score of the "mini-alignment" has to be above the threshold - the concept of "neighborhood words"

BLASTp - neighborhood words

Example - ITV triplet

	BLOSUM62	PAM230	
ITV - ITV	4+5+4 = 13	5+3+5 = 13	
ITV - MTV	1+5+4 = 10	2+3+5 = 10	
ITV - ISV	4+1+4 = 9	2+3+5 = 10	
ITV - LTV	2+5+4 = 11	2+3+5 = 10	
ITV - LSV	2+1+4 = 7	2+3+5 = 10	
ITV - MSV	1+1+4 = 6	2+3+5 = 10	
ITV - IAV	4+0+4 = 8	5+1+5 = 11	
ITV - MAV	1+0+4 = 5	2+1+5 = 8	
ITV - ITL	4+5+1 = 10	5+3+2 = 10	
ITV - LAV	2+0+4 = 6	2+1+5 = 8	

BLASTp - neighborhood words Threshold f = 11 (default for BLASTp) f=10 BLOSUM62 PAM230 BLOSUM62 PAM230 ITV - ITV ITV - ITV 4+5+4 = 135+3+5 = 4+5+4 = 13 5+3+5 = 13 ITV - MTV 1+5+4 = 10 2+3+5 = 10 ITV - MTV 1+5+4 = 102+3+5 = 10ITV - ISV 4+1+4 = 9 2+3+5 = 10 ITV - ISV 4+1+4 = 9 2+3+5 = 10 2+5+4 = 11 2+3+5 = 10 2+5+4 = 11 2+3+5 = 10 2+1+4 = 7 2+3+5 = 10 ITV - LTV ITV - LTV 2+1+4 = 7 2+3+5 = 10 ITV - LSV ITV - LSV ITV - MSV 1+1+4 = 6 2+3+5 = 10 ITV - MSV 1+1+4 = 6 2+3+5 = 10 ITV - IAV 4+0+4 = 8 5+1+5 = 11 ITV - IAV 4+0+4 = 8 5+1+5 = 11+0+4 = 5 2+1+5 = 8 ITV - MAV ITV - MAV 1+0+4 = 5 2+1+5 = 8 4+5+1 = 10 5+3+2 = 10 ITV - ITL 4+5+1 = 10 5+3+2 = 10 ITV - ITL 2+0+4 = 6 2+1+5 = 8 2+0+4 = 6 2+1+5 = 8 ITV - LAV ITV - LAV Pairs marked in blue would initiate an alignment extension

115

BLAST - FINAL STEP

- Smith-Waterman algorithm (local dynamic programming), discussed before but limited to regions that include the HSPs
- Significance of alignment with gaps can be evaluated using K and λ estimated from alignments of random sequences with same gap penalty and scoring parameters
- In spite of claims of being "mathematically rigorous" these parameters can only be estimated empirically

116

KARLIN-ALTCHUL STATISTICS

High scores of local alignments between two random sequences follow Extreme Value Distribution

117

KARLIN-ALTCHUL STATISTICS

For ungapped alignments their expected number with score S or greater equals

E = Kmne^{-λS}

K i λ , are parameters related to a search space and scoring system, and m, n represent a query and database length, respectively.

Score can be transformed to a bit-score according to formula S'= bitscore = (λ S - InK)/In2, then

 $E = mn2^{-S^2}$

118

KARLIN-ALTCHUL STATISTICS

- for ungapped alignments parameters K and λ are calculated algebraically but for gapped alignment a solid theory doesn't exist and these parameters are calculated by simulation which has to be run for every combination of scoring system including gap penalties
- ★ more at <u>http://www.ncbi.nlm.nih.gov/BLAST/</u> <u>tutorial/Altschul-1.html</u>

BLAST - KNOWN PROBLEMS

- ☆ Significance is calculated versus theoretic distribution using Karlin-Altschul equation not real sequences.
- ·⊱ Assumes sequences are random
- \ast Assume database is one long sequence length effects are not corrected for
- \ast Statistics are very inaccurate for short queries (ca. 20 characters).
- ★ Be careful when you change BLAST parameters, some of them should be coordinated, e.g. match/mismatch penalty and Xdrop off value
- * nucleotide BLAST default parameters tuned up for speed not sensitivity [Gotea, Veeramachaneni, and Makalowski (2003) Mastering seeds for genomic size nucleotide BLAST searches. Nucleic Acids Res. 31(23):6935-41]

BLAST ALGORITHM IMPLEMENTATON

Program	Query	Database type	
blastn	nt	nt	
megablast	nt	nt	
blastp	aa	aa	
blastx	nt	aa	
tblastn	aa	nt	
tblastx	nt	aa	
blast2seq	nt, aa	nt, aa	

121

BIOINFORMATICS CREED

- Remember about biology
- Do not trust the data
- Use comparative approach
- Use statistics
- Know the limits
- Remember about

