

CISC 889 Bioinformatics (Spring 2004)

Sequence pairwise alignment (I)

CISC889, S04, Lec4, Liao

Contents

Alignment algorithms

- Needleman-Wunsch (global alignment)
- Smith-Waterman (local alignment)
- Heuristic algorithms
 - FASTA
 - BLAST

CISC889, S04, Lec4, Liao

Sequence Alignment

- Motivation
 - Sequence assembly: reconstructing long DNA sequences from overlapping sequence fragments
 - Annotation: assign functions to newly discovered genes
 - Raw genomic (DNA) sequences -> coding sequences (CDS), candidate for genes -> protein sequence -> function
 - Terminologies: cDNA, RNA, mRNA
 - Evolution: mutation -> sequence diversity (vs homology) -> (new) phenotype ?
 - Corner stone: sequence similarity -> sequence homology -> same function

CISC889, S04, Lec4, Liao

Alignment algorithms

- What is an alignment?

A one-to-one matching of two sequences so that each character in a pair of sequences is associated with a single character of the other sequence or with a null character (gap). Alignments are often displayed as two rows with an optional third row in between pointing out regions of similarity.
- Example:

```
>gi|7434520|pir|G64632 acetate kinase - Helicobacter pylori (strain 26695)
Length = 388

Score = 35.8 bits (81), Expect = 0.10
Identities = 21/51 (41%), Positives = 29/51 (56%), Gaps = 2/51 (3%)

Query:  1  VLVLCGSSSLKFAIIDAVNGEYLSGLAECF--HLPARIKWKMDGNKQE 49
          +LVLN GSSS+KE + D   + SGLAE      + + +IK + N QE
Sbjct:  3  ILVLNLGSSSIKFKLFDMKENKPLASGLAEKIGEEIGQLIKSHLHHNDQE 53
```
- Types of alignment:
 - pairwise vs multiple;
 - global vs local
- Algorithms
 - Rigorous
 - heuristic

CISC889, S04, Lec4, Liao

PAM matrices (Margaret Dayhoff, 1978)

- point accepted mutation or percent accepted mutation
 - unit of measurement of evolutionary divergence between two amino acid sequences
 - substitute matrices (scoring matrices)
- 1 PAM = one accepted point-mutation event per one-hundred amino acids

CISC889, S04, Lec4, Liao

PAM (cont'd)

caveat:

- Sequences s1 and s2 are x PAM divergent does not imply s1 and s2 have x percent sequence difference (should be equal or less).

facts:

1. even amino acid sequences that have diverged by 200 PAM units are expected to be identical in about 25% of their positions.
2. sequences that are 250 PAM units diverged can generally be distinguished from a pair of random sequences.

CISC889, S04, Lec4, Liao

PAM matrix is a 20 by 20 matrix, and each element p_{ij} represents the expected evolutionary exchange between the two corresponding amino acids for sequences that are a specific number of PAM units diverged. That is,

$$p_{ij} = \log[f(i,j)/f(i)f(j)]$$

where $f(i)$ and $f(j)$ are the frequencies that amino acids A_i and A_j appear in the sequences, and $f(i,j)$ the frequency that A_i and A_j are aligned.

By construction, PAM n matrices with $n > 1$ are extrapolated from those with lower n . Namely, it assumes that the frequencies of the amino acids remain constant over time and that the mutational processes causing replacements in an interval of one PAM unit operate the same for longer periods.

CISC889, S04, Lec4, Liao

Ideally, to align two sequences that are x PAM units diverged, the PAM x matrix should be used.

Practically, since you do not know *a priori* how much diverged given two sequences, you may need to try different PAM matrices.

Empirically, PAM 250 matrix is reported to be very effective for aligning sequences used in evolutionary studies.

CISC889, S04, Lec4, Liao

PROSITE is a "dictionary of sites and patterns in proteins" that is linked to the protein sequence database Swiss-Prot.

The signature patterns in PROSITE are written as regular expressions.

For example, G-[GN]-[SGA]-G-x-R-[SGA]-C-x(2)-[IV] is a PROSITE pattern.

BLOCKS is a database of protein motifs that is derived from the PROSITE library.

CISC889, S04, Lec4, Liao

Needleman-Wunsch algorithm (Global Pairwise optimal alignment, 1970)

To align two sequences $x[1...n]$ and $y[1...m]$,

i) if x at i aligns with y at j , a score $s(x_i, y_j)$ is added; if either x_i or y_j is a gap, a score of d is subtracted (penalty).

ii) The best score up to (i,j) will be

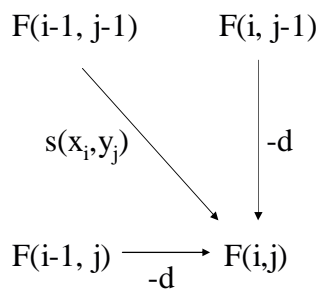
$$F(i,j) = \max \{ \begin{aligned} &F(i-1, j-1) + s(x_i, y_j), \\ &F(i-1, j) - d, \\ &F(i, j-1) - d \end{aligned} \}$$

CISC889, S04, Lec4, Liao

Needleman-Wunsch (cont'd)

iii) Tabular computing to get $F(i,j)$ for all $1 < i < n$ and $i < j < m$

Draw a diagram:



By definition, $F(n,m)$ gives the best score for an alignment of $x[1\dots n]$ and $y[1\dots m]$.

CISC889, S04, Lec4, Liao

iv) Trace-back

To find the alignment itself, we must find the path of choices (in applying the formulae of ii) when tabular computing that led to this final value.

- > Vertical move is gap in the column sequence.
- > Horizontal move is gap in the row sequence.
- > Diagonal move is a match.

CISC889, S04, Lec4, Liao

Example: Align HEAGAWGHEE and PAWHEAE.

Use BLOSUM 50 for substitution matrix and $d=-8$ for gap penalty.

(page 21 in the text)

		H	E	A	G	A	W	G	H	E	E
		0	-8	-16	-24	-32	-40	-56	-64	-72	-80
P		-8	-2	-9	-17	-25	-33	-42	-49	-57	-65
A		-16	-10	-3	-4	-12	-20	-28	-36	-44	-52
W		-24	-18	-11	-6	-7	-15	-5	-13	-21	-29
H		-32	-14	-18	-13	-8	-9	-13	-7	-3	-11
E		-40	-22	-8	-16	-16	-9	-12	-15	-7	3
A		-48	-30	-16	-3	-11	-11	-12	-12	-15	-5
E		-56	-38	-24	-11	-6	-12	-14	-15	-12	-9

CISC889, S04, Lec4, Liao

Time complexity: $O(nm)$

Space complexity: $O(nm)$

Big-O notation:

$f(x) = O(g(x)) \Rightarrow f$ is upper bound by g

$f(x) = \Omega(g(x)) \Rightarrow f$ is lower bound by g

$f(x) = \Theta(g(x)) \Rightarrow f$ is bound to g within constant factors

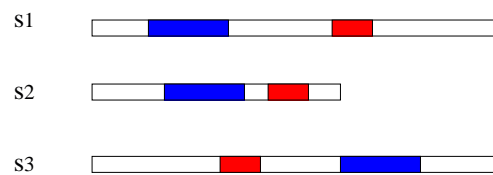
CISC889, S04, Lec4, Liao

Local pairwise optimal alignment

why need local alignment (vs global)?

- mosaic structure (functioning domains) of proteins

e.g., are these three sequences similar or not?



CISC889, S04, Lec4, Liao

Local alignment

- naive algorithm:
 - there are $\Theta(n^2 m^2)$ pairs of substrings; to align each pair as a global alignment problem will take $O(nm)$; the optimal local alignment will therefore take $O(n^3 m^3)$.
- Smith-Waterman algorithm (dynamic programming)
recurrence relationship
$$F(i,j) = \max \{ 0, \\ F(i-1, j-1) + s(x_i, y_j), \\ F(i-1, j) - d, \\ F(i, j-1) - d \}$$

Notes: 1) For this to work, the random match model must have a negative score.

2) The time complexity of Smith-Waterman is $\Theta(n^2 m^2)$.

CISC889, S04, Lec4, Liao

Example: Align HEAGAWGHEE and PAWHEAE.

Use BLOSUM 50 for substitution matrix and $d=-8$ for gap penalty.

(page 23 in the text)

		H	E	A	G	A	W	G	H	E	E
		0	0	0	0	0	0	0	0	0	0
P		0	0	0	0	0	0	0	0	0	0
A		0	0	0	5	0	5	0	0	0	0
W		0	0	0	0	2	0	20	12	0	0
H		0	10	2	0	0	0	12	18	22	14
E		0	2	16	8	0	0	4	10	18	28
A		0	0	8	21	13	5	0	4	10	20
E		0	0	0	13	18	12	4	0	4	16

CISC889, S04, Lec4, Liao

Heuristic alignment algorithms

- motivation: speed

sequence DB $\sim O(100,000,000)$ basepair

query sequence 1000 basepair

$O(nm)$ time complexity $\Rightarrow 10^{11}$ matrix cells in dynamic programming table

if 10,000,000 cells/second $\Rightarrow 10000$ seconds ~ 3 hours.

- heuristic versus rigorous

CISC889, S04, Lec4, Liao