Sequence Alignment



Sequence alignment

Introduction sequence alignment

Calculation of an alignment

- exercises

Different types of the algorithm

 exercises



Sequence alignment

Introduction sequence alignment

Calculation of an alignment

- exercises

Different types of the algorithm

 exercises



Sequence comparison

Pairwise sequence alignment

Multiple sequence alignment Database search

What is a pairwise alignment?

Two sequences are placed opposite each other to show similarity between the two

Insertions, deletions and mutations are introduced in both sequences to achieve this

CTCCTGAGGCAAATCTGTCAGTCCATCCTGGCTGAGTCCTCGCAGTCCCCGGCAGATCTTGAAGAAAAGA

Biological relevance

Insights into gene and protein function

High sequence similarity <u>usually</u> implies significant functional or structural similarity and therefore conserved functions

Examples

- Searching databases for related sequences and sub sequences
- Comparing two or more sequences for similarities
- Reconstructing long sequences of DNA from overlapping sequence fragments (sequence assembly)
- Exploring frequently occurring patterns of nucleotides
- Finding informative elements in protein and DNA sequences
- Find differences between sequences (sequence variation)

Database searching Heuristic - BLAST



Score Е Sequences producing significant alignments: (bits) Value gi|28804810|dbj|AB104818.1| Pipistrellus abramus HBA mRNA f... 115 1e-22 gi|49420|emb|X57029.1|CAGLOBINM M.auratus mRNA for alpha gl... 92 1e-15 3e-13 GUE gi|40225899|gb|BC032122.2| Homo sapiens hemoglobin, alpha 2... 84 gi|30047787|gb|BC050661.1| Homo sapiens hemoglobin, alpha 2... G 84 3e-13

*

Database searching Alignment

>gi|28804810|dbj|AB104818.1| Pipistrellus abramus HBA mRNA for alpha globin, complete cds Length = 554

Score = 115 bits (58), Expect = 1e-22
Identities = 142/170 (83%)
Strand = Plus / Plus

Query: 218 tggacgacctgccgtccgctctgtccgctctgtccgacctgcacgcttacaaactgcgtg 277

Sbjct: 248 tggacgacctgcccaccgccctgtccgccctgagcgacctgcacgcccacaagctgcgtg 307

Sbjct: 308 tggaccccgtcaacttcaagctcctgagccactgcctgctggtgaccctggcttgccacc 367

```
Score = 58.0 bits (29), Expect = 2e-05
Identities = 52/59 (88%), Gaps = 3/59 (5%)
Strand = Plus / Plus
```

>gi|49420|emb|X57029.1|CAGLOBINM Length = 418
M.auratus mRNA for alpha globin chain

Score = 91.7 bits (46), Expect = 1e-15 Identities = 151/186 (81%) Done ٠

Comparing sequences Dotplot







Reconstruct sequences Sequence assembly

Hierarchical shotgun sequencing



BAC library

Organized mapped large clone contigs

BAC to be sequenced

Shotgun clones



Reconstructing long sequences of DNA from overlapping sequence fragments

Finding informative elements Automated gene annotation



IGI/IPI, OTTO, humans

Finding informative elements Gene annotation



Finding informative elements Conservation



Finding informative elements Multiple sequence alignment



K-S

Finding informative elements Multiple sequence alignment



Thomas et al (2003), Nature

Sequence variation

Hemoglobin beta (HBB)

VS.

Hemoglobin beta in Sickle cell anemia (HBS)



Pairwise sequence comparison

Dotplot Edit transcript Sequence alignment

SequenSequence1Sequence1Sequence1Sequence1 Sequence2Seque2Sequence

Dot plots are two dimensional graphs, showing a comparison of two sequences

The two axis of the graph represent the two sequences being compared

Every region of the sequence is compared to every region of the other sequence

	Т	С	G	С	A	G	Т	С	С	С	С	G	G	С
С		Х		Х				Х	Х	Х	Х			Х
А					Х									
G			Х			Х						Х	Х	
С		Х		Х				Х	Х	Х	Х			Х
А					X									
G			Х			X						Х	X	
С		Х		Х				Х	Х	Х	Х			Х
А					X									
С		Х		Х				Х	Х	Х	Х			Х
Т	X						X							
С		Х		Х				Х	Х	Х	Х			Х
А					Х									
G			X			X						X	X	
С		Х		X				X	Х	Х	X			Х

Much noise without filter (four letter alphabet) Search for a given number of matches (<u>stringency</u>) with

Search for a given number of matches (*stringency*) within a given range (*window*)



Warning!

Window size smaller than region of similarity: segment pairs may have a score that is not sufficiently high enough

> ABCDEFGH XYZDXFGX

Window size to big: scores for unrelated amino acids may degrade the signal



Dotmatcher:

(windowsize = 5, threshold = 17.00 - 16/05/02)



Edit transcript

Transformation of one string into other string

- I insertion
- D deletion
- **R** replacement (substitution)
- M match

Edit transcript



example: S1 : vintner S2 : writers

 R	Ι	Μ	D	Μ	D	Μ	Μ	Ι	
V		i	n	t	n	е	r		
W	r	i		t		е	r	S	

Edit distance

Minimum number of operations to transform S1 into S2

Edit transcript with minimal operations: optimal transcript

Edit transcript vs Alignment

Edit transcript: shows mutation events Distance is minimized

Alignment: display relationship between strings Score is maximized

String alignment

S1 = qacdbdS2 = qawxb

q	а	С	-	d	b	d
q	а	W	X	I	b	-

String alignment - competition



Sequence alignment

Introduction sequence alignment

Calculation of an alignment

- exercises

Different types of the algorithm

 exercises



Dynamic programming

Dynamic programming

3 components:

- 1. Recurrence relation (*initialization*)
- 2. Tabular computation (*matrix filling*)
- 3. Traceback

Global alignment

Input:

Two sequences S1 and S2 of roughly the same length

S1 = acbcdb S2 = cadbd

Question:

What is the best similarity between them? Find best alignment

Some definitions

Two sequences: *S1[1..i]* and *S2[1..j]* with lengths: *n* and *m*

 S1:
 1 2 3 4 5 6 a c b c d b
 S2:
 1 2 3 4 5 c a d b d

Each character of *S1* (i = 0..n) is compared with *S2* (j = 0..m)

Recurrences global alignment

Base:

$$V(i,0) = \sum_{\substack{1 \le k \le j}} (S1(k), -)$$

Score from substitution matrix

For i>0 and j>0

$$V(i,j) = max [V(i-1, j-1) + s(S1(i), S2(j)),V(i-1, j) + s(S1(i), -),V(i, j-1) + s(-, S2(j))]$$

Substitution scores

Scoring matrix based on:

genetic code

- physical-chemical properties
- studies of molecular structure
- studies of evolution
- statistics




BLOSUM100 Scoring Matrix

	А	R	Ν	D	С	Q	Е	G	Н	I	L	К	Μ	F	Ρ	S	Т	W	Y	V	В	Ζ	Х	*
A_	8	-3	-4	-5	-2	-2	-3	-1	-4	-4	-4	-2	-3	-5	-2	1	-1	-6	-5	-2	-4	-2	-2	-10
R_		10	-2	-5	- 8	0	-2	-6	-1	-7	- 6	3	-4	-6	-5	-3	-3	-7	- 5	-6	-4	-1	-3	-10
N			11	1	- 5	-1	-2	-2	Θ	-7	-7	-1	-5	-7	-5	Θ	-1	- 8	- 5	-7	5	-2	-3	-10
D				10	- 8	-2	2	-4	-3	-8	- 8	-3	-8	-8	-5	-2	-4	-10	-7	-8	6	0	-4	-10
C					14	-7	-9	-7	-8	-3	-5	-8	-4	-4	-8	-3	-3	-7	-6	-3	-7	-8	-5	-10
Q-						11	2	-5	1	-6	-5	2	-2	-6	-4	-2	-3	-5	-4	-5	-2	5	-2	-10
E							10	-6	-2	-7	-7	0	-5	-8	-4	-2	-3	-8	-7	-5	Θ	7	-3	-10
G								9	-6	-9	- 8	- 5	-7	-8	-6	-2	-5	-7	-8	-8	-3	-5	-4	-10
H									13	-7	-6	-3	-5	-4	-5	-3	-4	-5	1	-7	-2	-1	-4	-10
I-										8	2	-6	1	-2	-7	-5	-3	-6	-4	4	-8	-7	-3	-10
L											8	-6	3	0	- /	-6	-4	-5	-4	0	-8	-6	-3	-10
K												10	-4	-6	-3	-2	-3	-8	-5	-5	-2	0	-3	-10
M													12	-1	-5	-4	-2	-4	-5	0	- /	-4	-3	-10
F														11	- /	-5	-5	0	4	-3	- /	- /	-4	-10
۲ <u>–</u>															12	-3	-4	-8	- /	-0	-5 1	-4	-4	-10
<u>১</u> т																9	2	- /	-5	-4 1	- T 2	-2	-2	-10
																	9	-/	- 5	- T	-2	-3	- 2	- 10 10
																			12	-5	-9	-7	-0	-10
ı V																			12	_J 8	-7	-5		-10
v R																				0	- 1	- J 0	-Δ	-10
7																					U	6	-2	-10
<u>~</u>																							-3	-10
*																								_ 1

Substitution scores

Common in protein alignments Also suggested for nucleotide sequences:

It is more likely that an A is replaced by a G than a T

Scoring DNA sequences: Most of the time identities are scored instead of putting weight to substitutions



S1

S2 (j)





S2



S2



S2



S2



S2



Scores: match = 2 mismatch = -1space = -1



	-	С	a	d	b	d
-	0	-1	-2	-3	-4	-5
a	-1	-1	1	← 0	← -1	← -2
C	-2	<u> </u>	↓ 0	× 0	-1	- 2
b	-3	† 0	` 0	- 1	<u> </u>	← 1
C	-4	↑ -1	<u>1</u>	` -1	1	<u> </u>
d	-5	[↑] -2	<u></u> -2	1	↓ 0	3
b	-6	^ -3	^ † -3	[†] 0	3	← 2

S2



Three possible moves in the *traceback:*

- *Diagonal*: letters are aligned
- *Left*: a gap is introduced in the left sequence (S1)
- *Up*: a gap in the top sequence (S2)

Trace back

- a c b c d b c a d b - d -







Recurrences global alignment

Base:
$$V(i,0) = \sum s(S1(k), -)$$

 $V(0,j) = \sum s(-, S2(k))$
 $1 \le k \le j$

V(i,j) = max [

V(i-1, j-1) + s(S1(i), S2(j)),

V(i-1, j) + s(S1(i), -),

V(i, j-1) + s(-, S2(j))

For i>0 and j>0



Sequence alignment

Introduction sequence alignment

Calculation of an alignment

exercises

Different types of the algorithm

 exercises



Ends-free alignment

Motivation: Shotgun sequence assembly

Genome

Shotgun Sequencing

Genome Assembly

Ends-free alignment



Recurrences ends-free alignment

Base: V(i,0) = 0V(0,j) = 0

```
For i>0 and j>0

V(i,j) = max [

V(i-1, j-1) + s(S1(i), S2(j)),

V(i-1, j) + s(S1(i), -),

V(i, j-1) + s(-, S2(j))

]
```

S2



S2



Scores: match = 2mismatch = -1space = -1



	-	g	C	t	a	a
-	0	0	0	0	0	0
C	0	-1	2	← 1	← 0	-1
g	0	<u> </u>	↓ 1	× 1	⊷ 0	-1
a	0	1	` 1	0	× 3	2
g	0	<u> </u>	← 1	- 0	[†] 2	<u> </u>
C	0	1	4	← 3	← 2	↓ 1
t	0	1 0	[†] 3	6	← 5	← 4

Trace back



S2

.

Trace back

cga g c t ----- g c t aa



	-	g	С	t	a	a
-	0	0	0	0	0	0
С	0	-1	2	← 1	← 0	- † -1
g	0	× 2	↓ 1	× 1	⊷ 0	-1
a	0	1	1	↓ 0	` 3	- 2
g	0	× 2	← 1	← 0	¹ 2	<u> </u>
С	0	1 1	× (4)	← 3	← 2	↓ 1
t	0	1 0	[†] 3	× 6	← 5	← 4

Recurrences ends-free alignment

Base: V(i,0) = 0V(0,j) = 0

For i>0 and j>0

V(i,j) = max [V(i-1, j-1) + s(S1(i), S2(j)),V(i-1, j) + s(S1(i), -),V(i, j-1) + s(-, S2(j))7



Local alignment

Finds regions of sequence with a high degree of similarity

Better at finding motifs, especially for sequences that are different overall

Will return only the best matching segment for a given pair of sequences

Global vs Local

Global : *Needleman-Wunsch* Local : *Smith-Waterman*

Local alignment far more meaningful in many biological applications

In some situations global alignment seems to be more effective in exposing important biological commonalities

Global vs Local

Example



Local alignment

Definition: Two sequences S1 and S2

Question: Find subsequences α and β of S1 and S2 whose similarity is maximum over all such pairs of subsequence

example:

S1 = g g t c t g a g

S2 = a a a c g a

Subseq α = c t g a

Subseq $\beta = c - g a$

Recurrences local alignment

Base:
$$V(i,0) = 0$$

 $V(0,j) = 0$












Tabular computation





S2

S1





Recurrences local alignment



Gapped alignments

Theory: single mutation can delete or insert several nucleotides or amino acids

Insertion of a gap must improve the quality of the alignment (raise quality score)

Different types of models



Gapped alignments

Affine is most used model

Gap opening/creation penalty (GOP): high Gap extending penalty (GEP): low

Multiple alignment

Multiple alignment



http://bibiserv.techfak.uni-bielefeld.de/visualign/

Motivation

It may reveal evolutionary commonalities

- Conserved motifs
- Common 2 and 3-dimensional structure
- Clues about common biological function
- It is used for characterizing families or superfamilies of proteins

Algorithms

Dynamic programming in higher dimensions

+ true optimization

 execution time & memory requirements grow exponentially

Algorithms

Perform all pairwise alignments Create consensus of two most similar pair Add more sequences to this generalized – *consensus* sequence

- + fast
- errors in early alignments affect rest of alignments

Correct mistakes

Remove sequences and realign them

Find short highly conserved regions, align them and repeat procedure on unaligned regions

Sequence alignment

Introduction sequence alignment

Calculation of an alignment

- exercises

Different types of the algorithm

 exercises



Algorithms for computing sequence alignments

Global

find overall alignment, rarely used

Ends-free

find overlap in sequence ends, sequence assembly

Local

find similar sub-sequences, often used

