

Algorithms in Bioinformatics

Subtitle

- <u>Baum-Welch algorithm</u>
- Binarization of consensus partition matrices
- <u>BLAST</u>
- <u>Blast2GO</u>
- <u>Bowtie (sequence analysis)</u>
- <u>De novo sequence assemblers</u>
- <u>High-performance Integrated Virtual Environment</u>
- <u>Hirschberg's algorithm</u>

Island algorithm

- <u>Kabsch algorithm</u>
- <u>Microarray analysis techniques</u>
- <u>Needleman–Wunsch algorithm</u>
- <u>Neighbor joining</u>
- <u>Pairwise Algorithm</u>
- <u>Pseudo amino acid composition</u>
- <u>PSI Protein Classifier</u>

<u>Quartet distance</u>

- <u>Quasi-median networks</u>
- <u>Robinson–Foulds metric</u>
- <u>SAMtools</u>
- <u>SCHEMA (bioinformatics)</u>
- <u>Sequential pattern mining</u>
- <u>Smith–Waterman algorithm</u>
- <u>SPAdes (software)</u>

<u>TopHat (bioinformatics)</u>

- <u>UCLUST</u>
- <u>Ukkonen's algorithm</u>
- <u>UPGMA</u>
- <u>Velvet assembler</u>
- <u>ViennaRNA Package</u>
- <u>WPGMA</u>
- <u>Z curve</u>



Homology, orthology, paralogy

Homology: descent from a common ancestor



Orthology: descent from a common ancestor by genome division Paralogy: descent from a common ancestor by duplication within a genome





Global and Local Alginment

- Global alignments, which attempt to align every residue in every sequence, are most useful when the sequences in the query set are similar and of roughly equal size. A general global alignment technique is the <u>Needleman–Wunsch</u> <u>algorithm</u>, which is based on dynamic programming. (BLAST & FASTA)
- Local alignments are more useful for dissimilar sequences that are suspected to contain regions of similarity or similar sequence motifs within their larger sequence context. The <u>Smith–Waterman algorithm</u> is a general local alignment method also based on dynamic programming.

Local Alignment

Target Sequence

5' ACTACTAGATTACTTACGGATCAGGTACTTTAGAGGCTTGCAACCA 3'

Global Alignment

Global versus Local Alignment

A global alignment covers the entire lengths of the sequences involved

The Needleman-Wunsch algorithm finds the best global alignment between 2 sequences

A local alignment only covers parts of the sequences
The Smith-Waterman algorithm finds the best local alignment
between 2 sequences

Global alignment		Q	ĸ	E 	s I	G	P	S	S	S	Y	C I
	v	Q	Q	Е	S	G	L	v	R	т	т	С
Local alignment					E I	s I	G I					
					Е	S	G					

Pair wise and Multiple sequence alignment

- Pairwise sequence alignment methods are used to find the best-matching piecewise (local or global) alignments of two query sequences. Pairwise alignments can only be used between two sequences at a time.
- <u>Multiple sequence alignment</u> is an extension of pairwise alignment to incorporate more than two sequences at a time. Multiple alignment methods try to align all of the sequences in a given query set. Multiple alignments are often used in identifying <u>conserved</u> sequence regions across a group of sequences hypothesized to be evolutionarily related.

Pairwise versus Multiple Alignment

 So far we have considered the alignment of two sequences ('pairwise alignment')

> Q K E S G P S S S Y C | | | | | | | V Q Q E S G L V R T T C

 Alignment can be performed between three or more sequences ('multiple alignment')

> Q K E S G P S S S Y C | | | | | | | | V Q Q E S G L V R T T C | | | | | | | | | V Q K E S L L V R S T C

Dot-matrix

 To construct a dot-matrix plot, the two sequences are written along the top row and leftmost column of a twodimensional <u>matrix</u> and a dot is placed at any point where the characters in the appropriate columns match—this is a typical <u>recurrence plot</u>. Some implementations vary the size or intensity of the dot depending on the degree of similarity of the two characters, to accommodate conservative substitutions.



An Introduction to Bioinformatics Algorithms

www.bioalgorithms.info

Dot Matrices (cont'd)

- Identify diagonals above a threshold length
- Diagonals in the dot matrix indicate exact substring matching



Dynamic programming

The technique of <u>dynamic programming</u> can be applied to produce global alignments via the <u>Needleman-Wunsch</u> <u>algorithm</u>, and local alignments via the <u>Smith-Waterman</u> <u>algorithm</u>. In typical usage, protein alignments use a <u>substitution matrix</u> to assign scores to amino-acid matches or mismatches, and a <u>gap penalty</u> for matching an amino acid in one sequence to a gap in the other. Dynamic programming matrix:



Optimum alignment scores 11:

т	-	-	т	С	Α	т	Α
т	G	С	т	С	G	т	Α
+5	-6	-6	+5	+5	-2	+5	+5

Substitution Matrix – PAM & BLOSUM

PAM

- One of the first amino acid substitution matrices, the PAM (<u>Point Accepted Mutation</u>) matrix was developed by <u>Margaret Dayhoff</u> in the 1970s.
- This matrix is calculated by observing the differences in closely related proteins. The PAM1 matrix estimates what rate of substitution would be expected if 1% of the amino acids had changed. There are many PAM matrices like PAM 1 to PAM30)

	С	S	т	P	A	G	Ν	D	E	Q	н	R	К	м	Ι	L	V	F	Y	М	
C.	12																				Ų.
S	Ŭ.	- 4																			5
T	- 2	1	3																		Т
Ŧ	43	1	Q	¢																	Ð
A	+2	1	1	I	2																д
G	43	1	0	-1	1	5															G
11	-4		¢	-1	¢	0															
Ð	-5	Ū	0	-1	Q	1	24	4													D
E	+5	0	O.	+1	D.	Ō	1	3	4												E
Ō	-5	- 1	+1	Ó.	Ō	-1	1	2	2	d.											1
日		- 1	_ 1	n	-1	_ 2	2	1		9											н Н
R	- 4	n.	- 1	ñ	<u>_2</u>	<u>_ 9</u>	ā	<u>+ 1</u>	_ ī	ī		6									R
TC TC	_ F ,	ñ	ā.	_ 7	_ ,	_7	7	ň	ñ	T.	ñ		τ.								V
		25:	11	2.50	2.51			251	1.2	2.4	25	- 61 -	់ក	6	:::::	1,1,1			:::::	11111	
			i de la composition de la comp	25	말유						5										T S
崇			3		문구한	- - -		- 4 - 4		1.5				4							internet
14 17			781	56	[음음]		Es:	53		741		53	5 C							1993	·부···································
. <i>W</i>	74	, , , , , , , , , , , , , , , , , , ,		고 복 :	<u>. u</u>		- 4	:7:4:		<u>.</u> д.	- 4	二 24 (4-	- 4.	:. <u></u> ‡:	4.	· · · · · · · · · · · · · · · · · · ·			• : • : • :	
	- 4				- 4	- 2	- +			1	- 4			Ц. 	-						2
ľ	U	-3	-3	- 3	- 3		- 2	-4	-4	-4	U	-4	-4	- 4	- L	- 1	-2		10		ľ
	-8	-2	- 5	-0	-0-		-4		+/	+5	-3		-	-4	-5	- Z-	- D		U	17	
	C	S	Т	P	A	G	N	D	E	Q	H	R	К	M	Ι	L	V	F	Y	M	

BLOSUM (BLOck SUbstitution Matrix)

- Henikoff constructed these matrices using multiple alignments of evolutionarily divergent proteins.
- The probabilities used in the matrix calculation are computed by looking at "blocks" of conserved sequences found in multiple protein alignments.
- These conserved sequences are assumed to be of functional importance within related proteins and will therefore have lower substitution rates than less conserved regions. There are many BLOSUM (1–70)

BLOSUM Substitution Matrix





	С	s	т	P	A	G	N	D	E	Q	н	R	ĸ	м	I	L	v	F	Y	w	
С	9																				C
S	-1	4																			S
т	-1	1	5																		т
P	-3	-1	-1	7																	P
A	0	1	0	-1	4																A
G	-3	0	-2	-2	0	6															G
N	-3	1	0	-2	-2	0	6														N
D	-3	0	-1	-1	-2	-1	1	6													D
E	-4	0	-1	-1	-1	-2	0	2	5												E
Q	-3	0	-1	-1	-1	-2	0	0	2	5											Q
н	-3	-1	-2	-2	-2	-2	1	-1	0	0	8										н
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5									R
ĸ	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5								ĸ
м	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5							м
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4						I
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4					L
v	-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4	_			v
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6			F
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7		Y
w	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11	w
	C	S	т	P	A	G	N	D	E	Q	н	R	ĸ	M	I	L	v	F	Y	w	
										(a	1)										
Tar	aet se	auena	10											٨		r.	F		n		v
1 ar	get se	quen												~		-					912
Fun	ctiona	ally Id	dentic	al or	Relat	ed Er	nzym	es (pi	oper	ty A	protei	ins)		н		L	E		Р	1	L
Seq	uence	-relat	ed Pr	otein	s (~A	prot	eins)							L		L	Р		D	1	L
C																					
Sco	re													-1		0	0		-/	114	U
										11											

(b)

1998 P