Smith & Waterman Algorithm for Pairwise Local Alignment

Kawshik Kumar Paul(1705043) Monirul Haque Imon(1705054) Maisha Rahman Mim (1705060)

Department of Computer Science and Engineering Bangladesh University of Engieering and Technology Dhaka. Bangladesh

July 6, 2021



Lecture Outline

- Sequence Alignment
- Sequence Edits
- Openation of the second of
- Local Sequence Alignment
- 5 Smith & Waterman Algorithm
- 6 Example

- Sequence Alignment
- Sequence Edits
- Openie Programming
- Local Sequence Alignment
- 5 Smith & Waterman Algorithm
- 6 Example

Sequence Alignment

- Why do we need to align sequence?
- Evolutionary Relationships

Why do we need to align sequence?

- Comparing DNA/protein sequences for
 - Similarity
 - Homology
- Prediction of function
- Construction of phylogeny Shotgun assembly
 - End-space-free alignment / overlap alignment
- Finding motifs

Sequence Alignment

 Procedure of comparing to (Pairwise) or more (Multiple) sequences by searching for a series of individual characters that are in the same order in the sequence.



Sequence Alignment

Definition

Given two strings $x = x_1x_2...x_m$ and $y = y_1y_2...y_n$, an alignment is an assignment of gaps to positions 0 ... M in x and 0 ... N in y, so as to line up each letter in one sequence with either a letter or a gap in the other sequence.

A Simple Alignment

- Let us try to align two short nucleotide sequences:
 - -AATCTATA and AAGATA

- Without considering any gaps (insertions/deletions) there are 3 possible ways to align these sequences
 - AATCTATA AATCTATA AATCTATA AAGATA AAGATA AAGATA
- Which one is better?

What is a Good Alignment

```
A G G C T A G T T . A G C G A A G T T

A G G C G A A G T T

A G G C G A A G T T

Matches = 6

Mistatches = 3

Gap = 1

A G G C T A - G T T -

A G - C G A A G T T

Matches = 7

Mismatches = 1

Gaps = 3

A G G C - T - G T T -

A G - C G - A A G T T

Mismatches = 0

Gaps = 5
```

Scoring the Alignments

- We need to have a scoring mechanism to evaluate alignments
 - match score
 - mismatch score
- We can have the total score as:

$$\sum_{i=1}^{n} match or mismatch score at position i$$

 For the simple example, assume a match score of 1 and a mismatch score of 0:



Simple Alignment with Gaps

 Considering gapped alignments vastly increases the number of possible alignments:

• If gap penalty is -1, what will be the new scores?

- Sequence Alignment
- Sequence Edits
- Opening Programming
 Opening
- Local Sequence Alignment
- 5 Smith & Waterman Algorithm
- Example

Sequence Edits

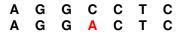
Lets do some sequence edits and view scores

Sequence Edits

Three types of sequence edits

- Mutations
- Insertions
- Operation
 Operation

Mutations



Insertions

Deletions

Scoring Function

Match: +m

Mismatch: -s

Gap: -d

Score $F = (\# \text{ matches}) \times m - (\# \text{ mismatches}) \times s - (\# \text{ gaps}) \times d$

Score Matrix

- Assign scores to each pair of symbol
 - Higher score means more similarity.

Score Matrix

- DNA
 - Match = +1
 - Mismatch = -3
 - Gap penalty = -5
 - Gap extension penalty = -2
- Protein sequences
 - Blossum62 matrix
 - Gap open penalty = -11
 - Gap extension penalty = -1

- Sequence Alignmen
- Sequence Edits
- Opening Programming
- 4 Local Sequence Alignment
- Smith & Waterman Algorithm
- Example

• How do we compute the best alignment?

Alignment is Additive

Observation:

The score of aligning $x_1...x_M$ and $y_1...y_N$ is additive

Say that
$$x_1...x_i$$
 $x_{i+1}...x_M$ aligns to $y_1...x_j$ $x_{j+1}...x_N$

The two scores add up:

$$F(x[1:M],y[1:N]) = F(x[1:i],y[1:j]) + F(x[i+1:M],y[j+1:N])$$

Types of Alignment

- Global
 - Strings of similar size
 - Genes with a similar structure
 - Larger regions with a preserved order (syntenic regions)
- Local
 - Finding similar regions among:
 - Dissimilar regions
 - Sequences of different lengths

Dynamic Programming

- Instead of evaluating every possible alignment, we can create a table of partial scores by breaking the alignment problem into subproblems.
- Consider two sequences CACGA and CGA
 - we have three possibilities for the first position of the alignment

First Position	Score	Remaining seqs				
С	+1	ACGA				
С	+1	GA				
_	-1	CACGA				
С	-1	GA				
С	-1	ACGA				
_	-1	CGA				

- Sequence Alignmen
- Sequence Edits
- Openie Programming
 Openie Programming
- Local Sequence Alignment
- Smith & Waterman Algorithm
- 6 Example

Local Sequence Alignment

- Suppose we have a long DNA sequence (eg 4000 bp) and we want to compare it with the complete yeast genome (12.5Mbp)
- What if only a portion of our query, say 200 bp length, has strong similarity to a gene in yeast.

Local Sequence Alignment Problem

Given two strings

$$x = x_1 x_M$$

$$y = y_1 x_N$$

Find substring $y^{,}$, $x^{,}$ whose similarity (optimal global alignment value) is maximum.

x = aaaacccccggggtta

y = ttcccgggaaccaacc

- Sequence Alignmen
- Sequence Edits
- Openie Programming
- Local Sequence Alignment
- 5 Smith & Waterman Algorithm
- 6 Example

Smith & Waterman Algorithm

- F(i,j)= optimal local similarity among suffixes A(1:i) and B(1:j)
- Recurrence relation
 - F(i,0) = 0
 - F(0, j) = 0
 - F(i,j) = max[0, F(i,j-1) + s(-,B(j)), F(i-1,j) + s(A(i),-), F(i-1,j-1) + s(A(i),B(j)]

- Sequence Alignmen
- Sequence Edits
- Opening Programming
- Local Sequence Alignment
- 5 Smith & Waterman Algorithm
- 6 Example

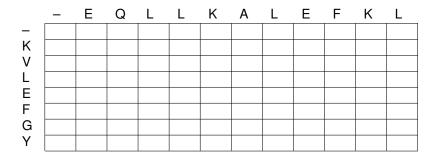
Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2



Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

$$F(i,0) = 0$$

$$F(0,j)=0$$

	_	E	Q	L	L	K	Α	L	Ε	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0											
V	0											
L	0											
Ε	0											
F	0											
G	0											
Υ	0											

Linear gap model

Gap = -1

Match= 4

Mismatch= -2

$$\begin{split} F(i,j) &= max \ [0, \ F(i,j\text{-}1) + s(\text{-},Q(j)), F(i\text{-}1,j) + s(P(i),\text{-}), F(i\text{-}1,j\text{-}1) + s(P(i),Q(j)] \end{split}$$

	_	Е	Q	L	L	K	Α	L	Ε	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0										
٧	0											
L	0											
Ε	0											
F	0											
G	0											
Υ	0											

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

$$\begin{split} F(i,j) &= max \ [0, \ F(i,j\text{-}1) + s(\text{-},Q(j)), F(i\text{-}1,j) + s(P(i),\text{-}), F(i\text{-}1,j\text{-}1) + s(P(i),Q(j)] \end{split}$$

	_	Е	Q	L	L	K	Α	L	Е	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0									
٧	0											
L	0											
Ε	0											
F	0											
G	0											
Υ	0											

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

$$\begin{split} F(i,j) &= max \ [0, \ F(i,j\text{-}1) + s(\text{-},Q(j)), F(i\text{-}1,j) + s(P(i),\text{-}), F(i\text{-}1,j\text{-}1) + s(P(i),Q(j)] \end{split}$$

	_	Е	Q	L	L	K	Α	L	Ε	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0								
٧	0											
L	0											
Ε	0											
F	0											
G	0											
Υ	0											

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

$$\begin{split} F(i,j) &= max \left[0, \, F(i,j\text{-}1) \, + \, s(\text{-},Q(j)), F(i\text{-}1,j) \, + \\ s(P(i),\text{-}), F(i\text{-}1,j\text{-}1) \, + \, s(P(i),Q(j)] \end{split}$$

	_	Ε	Q	L	L	K	Α	L	Ε	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0							
٧	0											
L	0											
Ε	0											
F	0											
G	0											
Υ	0											

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

$$\begin{split} F(i,j) &= max \ [0, \ F(i,j\text{-}1) + s(\text{-},Q(j)), F(i\text{-}1,j) + s(P(i),\text{-}), F(i\text{-}1,j\text{-}1) + s(P(i),Q(j)] \end{split}$$

	_	Ε	Q	L	L	K	Α	L	Е	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4						
٧	0											
L	0											
Ε	0											
F	0											
G	0											
Υ	0											

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

$$\begin{split} F(i,j) &= max \left[0, \, F(i,j\text{-}1) + s(\text{-},Q(j)), F(i\text{-}1,j) \right. + \\ s(P(i),\text{-}), F(i\text{-}1,j\text{-}1) + s(P(i),Q(j)] \end{split}$$

	_	Ε	Q	L	L	K	Α	L	Ε	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
٧	0											
L	0											
Ε	0											
F	0											
G	0											
Υ	0											

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

$$\begin{split} F(i,\,j) &= max\,[0,\,F(i,j\text{-}1)\,+\,s(\text{-},Q(j)),F(i\text{-}1,j)\,+\,\\ s(P(i),\text{-}),F(i\text{-}1,j\text{-}1)\,+\,s(P(i),Q(j)] \end{split}$$

	_	Е	Q	L	L	K	Α	L	E	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
Ε	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Υ	0	1	0	0	0	0	0	2	7	12	11	10

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

Alignment

Q:F P:F

	_	Е	Q	L	L	K	Α	L	Ε	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
Ε	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Υ	0	1	0	0	0	0	0	2	7	12	11	10

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

Alignment

Q:EF P:EF

Ε Q Κ Α Е F Κ Κ n n O E F G Υ

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

Alignment

Q: ...LEF

P: ...LEF

	_	Е	Q	L	L	K	Α	L	Е	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
Ε	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Υ	0	1	0	0	0	0	0	2	7	12	11	10

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

Alignment

Q: ..-LEF

P: ..VLEF

	_	Ε	Q	L	L	K	Α	L	Ε	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
Ε	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Υ	0	1	0	0	0	0	0	2	7	12	11	10

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

Alignment

Q: .A-LEF

P: .-VLEF

	_	Е	Q	L	L	K	Α	L	Ε	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
٧	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
Ε	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Υ	0	1	0	0	0	0	0	2	7	12	11	10

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

Alignment

Q: KA-LEF

P: K-VLEF

	_	Е	Q	L	L	K	Α	L	Ε	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
٧	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
Ε	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Υ	0	1	0	0	0	0	0	2	7	12	11	10

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

Alignment

Q: KA-LEF

P: K-VLEF

	_	Е	Q	L	L	K	Α	L	Ε	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
Ε	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Υ	0	1	0	0	0	0	0	2	7	12	11	10

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

Alignment

Q: K-ALEF

P: KV-LEF

	_	Е	Q	L	L	K	Α	L	Е	F	K	L
_	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
Ε	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Υ	0	1	0	0	0	0	0	2	7	12	11	10

Q = EQLLKALEFKL P = KVLEFGY

Linear gap model

Gap= -1

Match= 4

Mismatch= -2

Alignment

Q: KALEF P: KVLEF

Ε Q Κ Α Ε F Κ Κ n n E F G Υ

Thank you

Stay Home, Stay Safe