

Bioinformatics: Sequence Alignment

Swakkhar Shatabda

Department of Computer Science and Engineering
BRAC University



Book Reference

Bioinformatics Algorithms, An Active Learning Approach , Vol 1, Chapter 5



Inspiring Excellence

Sequence Alignment

Biological Question

- How can we find similarity between two sequences?
- Why is it important?
- Similar Sequence → Similar Structure → Similar Function
- The purpose of sequence alignment is to line up all residues in the inputted sequences for maximal level of similarity, in the sense of their functional or evolutionary relationship.

ATGCATGC

A**TGCATGC**–

A**TGC**–**TTA**–

TGCATGCA

–**TGCATGC**A

–**TGC****ATT**A

Pairwise Sequence Similarity:

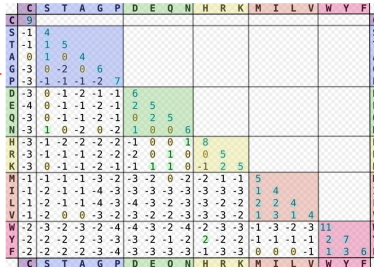
<https://www.ebi.ac.uk/jdispatcher/psa>



Inspiring Excellence

Pairwise Sequence Alignment

```
#=====
#
# Aligned_sequences: 2
# 1: test1
# 2: test2
# Matrix: EBLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
#
# Length: 105
# Identity:      96/105 (91.4%)
# Similarity:    99/105 (94.3%)
# Gaps:          0/105 ( 0.0%)
# Score: 522.0
#
#
#=====
```



test1	1	MGDVEKGKKIFIMKCSQCHTVEKGGKHKTGPNLHGLFGRKTGQAPGYSYT	50
		: . : : : .	
test2	1	MGDVEKGKKIFVQKCAQCCHTVEKGGKHKTGPNLHGLFGRKTGQAAGFSYT	50
test1	51	AANKNKGIIWGEDTLMEYLENPKKYIPGTMIFVGIIKKKEERADLIAYLK	100
		. . : : . .	
test2	51	DANKNKGITWGEDTLMEYLENPKKYIPGTMIFAGIIKKGERADLIAYLK	100
test1	101	KATNE	105
test2	101	KATNE	105

Sequence Alignment Problem

- Input Data:
 - Two sequences $S1$ and $S2$
- Parameter (s)
 - A scoring function f for
 - Substitutions
 - Gaps
- Output
 - The optimal alignment of $S1$ and $S2$ that has the maximal score.

$$\arg \max_{align} (f(align(S1, S2)))$$

- Enumerating all possible alignments is not feasible.
- A residue can either align to another residue or to a gap.
- We can use **dynamic programming**.



Inspiring Excellence

Formulation

- Align two sequences, x and y
- $F(i, j)$ is the score of the best alignment between $x_{1\dots i}$ and $y_{1\dots j}$
- $s(A, B)$ is the score for substituting A with B
- d is the (linear) gap penalty

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

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) & \mathbf{x_i \text{ aligned to } y_j} \\ F(i-1, j) + d & \mathbf{x_i \text{ aligned to a gap}} \\ F(i, j-1) + d & \mathbf{y_j \text{ aligned to a gap}} \end{cases}$$

Dynamic Programming

- Break the problem into smaller sub-problems.
- Solve these sub-problems optimally recursively.
- Use these optimal solutions to construct an optimal solution for the original problem.

Sequence Alignment: Example

Input Sequence 1: AAG

Input Sequence 2: AGC

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

For simplicity, let's set (i.e. **linear gap penalty**)
gap OPEN (d) = gap EXTEND (e) = -5

GAC-AT

C-ACAT

$$(-7) + (-5) + (-7) + (-5) + 2 + 2 = -20$$



Inspiring Excellence

Sequence Alignment: Example

		A	A	G
A				
G				
C				



Inspiring Excellence

Sequence Alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

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

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

Find the optimal alignment of AAG and AGC.

Use a linear gap penalty of $d=-5$.

		A	A	G
	0			
A				
G				
C				

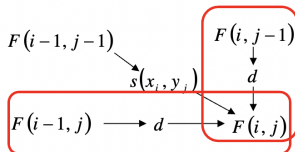


Inspiring Excellence

Sequence Alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

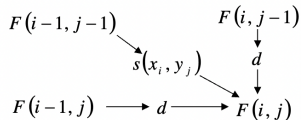
Find the optimal alignment of AAG and AGC.
Use a linear gap penalty of $d=-5$.



		A	A	G
	0	-5	-10	-15
A	-5			
G	-10			
C	-15			

Sequence Alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2



Find the optimal alignment of AAG and AGC.
Use a linear gap penalty of $d=-5$.

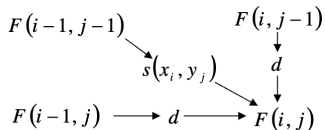
		A	A	G
	0	-5	-10	-15
A	-5	2	-3	-8
G	-10	-3	-3	-1
C	-15	-8	-8	-6

Sequence Alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

Find the optimal alignment of AAG and AGC.
Use a linear gap penalty of $d=-5$.

		A
	0	-5
A	-5	2

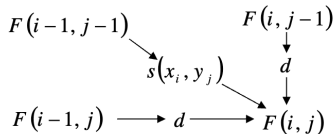


Sequence Alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

Find the optimal alignment of AAG and AGC.
Use a linear gap penalty of $d=-5$.

		A
	0	-5
A	-5	2



$$-5 + (-5) = -10$$

$$0 + 2 = 2$$

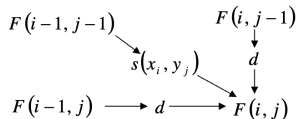
$$-5 + (-5) = -10$$



Inspiring Excellence

Sequence Alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2



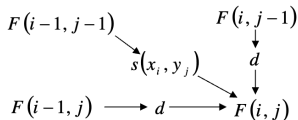
Find the optimal alignment of AAG and AGC.
Use a linear gap penalty of $d=-5$.

		A	A	G
	0	-5	-10	-15
A	-5	2	-3	-8
G	-10	-3	-3	-1
C	-15	-8	-8	-6

Sequence Alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

Find the optimal alignment of AAG and AGC.
Use a linear gap penalty of $d=-5$.



		A	A	G
	0	-5	-10	-15
A	-5	2	-3	-8
G	-10	-3	-3	-1
C	-15	-8	-8	-6

Sequence Alignment: Example

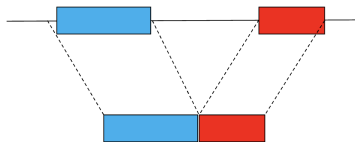
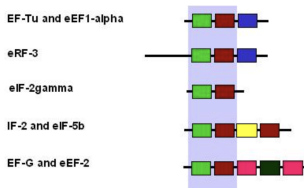
- Trace back to the **upper left**. Each arrow introduces **one** symbol at the end of each aligned sequence.

A A G -
- A G C
A A G -
A - G C

		A	A	G
	0	-5		
A		2	-3	
G				-1
C				-6

Limitations of Global Alignment

- Two functionally related proteins might be significantly/ largely different in their whole sequences but share similar important functional domains. The sequence fragment of the functional domain might be very conservative across different proteins in the same protein family, and determine the biological function.
- Secondly, the new discovery of introns required the DNA sequence alignment algorithms to be able to handle large deletions and interspersed conserved fragments (exons) and variable fragments (introns). Needleman-Wunsch → Smith-Waterman



Sequence Alignment: Local

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

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

Global alignment

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

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \\ 0 \end{cases}$$

Local alignment

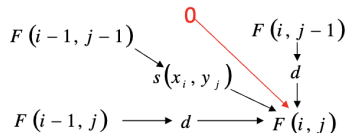


Inspiring Excellence

Sequence Alignment: Local

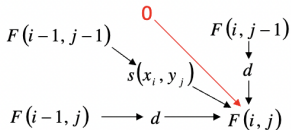
$$F(0,0) = 0$$

$$F(i,j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \\ 0 \end{cases}$$



Sequence Alignment: Local

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2



Find the optimal **local alignment** of AAG and AGC.
Use a linear gap penalty of $d = -5$.

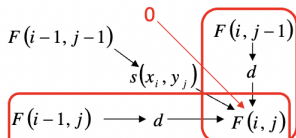
		A	A	G
A				
G				
C				

Sequence Alignment: Local

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

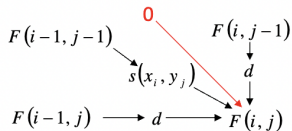
Find the optimal **local alignment** of AAG and AGC.
Use a linear gap penalty of $d = -5$.

		A	A	G
	0	0	0	0
A	0			
G	0			
C	0			



Sequence Alignment: Local

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

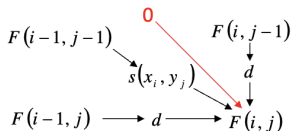


Find the optimal **local alignment** of AAG and AGC.
Use a linear gap penalty of $d = -5$.

		A	A	G
	0	0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0

Sequence Alignment: Local

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2



Find the optimal **local alignment** of AAG and AGC.
Use a linear gap penalty of $d = -5$.

		A	A	G
	0	0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0

Sequence Alignment: Local

- Trace back begins at **the highest score** in the matrix and continues **until you reach 0**.

A G
A G

		A	A	G
	0	0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0



Inspiring Excellence

Sequence Alignment: Local

- And also the **secondary best** alignment

A
A

		A	A	G
	0	0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0

Sequence Alignment: Local

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

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

A	A	G	-	A	A	G	-
-	A	G	C	A	-	G	C

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

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \\ 0 \end{cases}$$

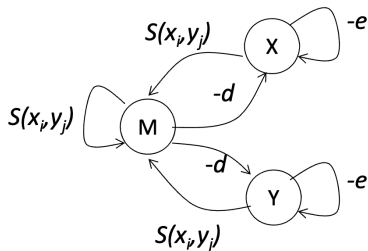
A	G
A	G

A
A



Inspiring Excellence

Affine Gap Penalty: Alignment as a series of state(s)



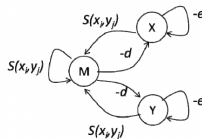
M	Match (<i>not necessarily identical</i>)
X	Insert at sequence X (delete at sequence Y)
Y	Insert at sequence Y (delete at sequence X)

d	Gap open
e	Gap Extension

Affine Gap Penalty: Alignment as a series of state(s)

- $M(i,j)$ is the score of the best alignment between $x_{1...i}$ and $y_{1...j}$, given x_i **aligned to** y_j
- $X(i,j)$ is the score of the best alignment between $x_{1...i}$ and $y_{1...j}$, given x_i **aligned to a gap**
- $Y(i,j)$ is the score of the best alignment between $x_{1...i}$ and $y_{1...j}$, given y_j **aligned to a gap**

$$M(i, j) = \max \begin{cases} M(i-1, j-1) + s(x_i, y_j) \\ X(i-1, j-1) + s(x_i, y_j) \\ Y(i-1, j-1) + s(x_i, y_j) \end{cases}$$



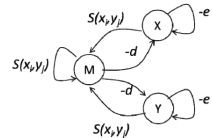
$$X(i, j) = \max \begin{cases} M(i-1, j) - d \\ X(i-1, j) - e \end{cases}$$

$$Y(i, j) = \max \begin{cases} M(i, j-1) - d \\ Y(i, j-1) - e \end{cases}$$

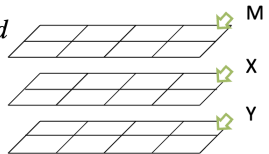
Affine Gap Penalty: Alignment as a series of state(s)

$$x_i \text{ aligned to } y_j \quad M(i, j) = \max \begin{cases} M(i-1, j-1) + s(x_i, y_j) & \text{after a match} \\ X(i-1, j-1) + s(x_i, y_j) \\ Y(i-1, j-1) + s(x_i, y_j) \end{cases} \quad \text{after a gap}$$

$$x_i \text{ aligned to a gap} \quad X(i, j) = \max \begin{cases} M(i-1, j) - d \\ X(i-1, j) - e \end{cases}$$

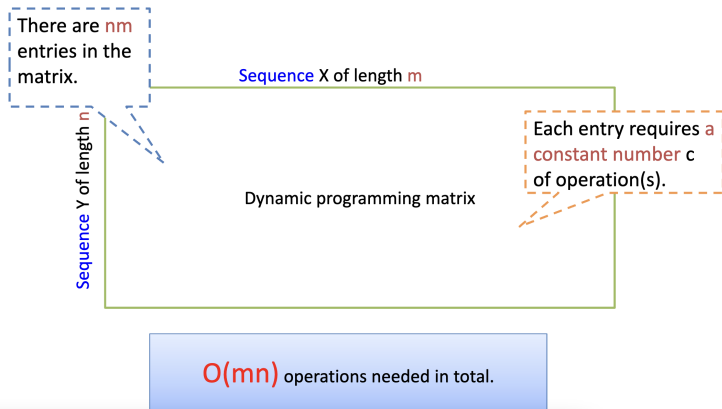


$$y_j \text{ aligned to a gap} \quad Y(i, j) = \max \begin{cases} M(i, j-1) - d \\ Y(i, j-1) - e \end{cases}$$



Runtime?

- How to do the search on 550K proteins or billions of reads?
- BLAST: <https://blast.ncbi.nlm.nih.gov/Blast.cgi>



Multiple Sequence Alignment

- If sequence similarity is weak, pairwise alignment may not identify biologically related sequences.
- simultaneous comparison of many sequences often allows us to find similarities that pairwise sequence comparison fails to reveal.
- Bioinformaticians sometimes say that while pairwise alignment whispers, multiple alignment shouts.

```
YAFDLGYTCMFPVLLGGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTA
-AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYI
IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPPALLKQCLVSA
```

```
SFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTE-FINHYGPTEATIGA
-YEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS
----PTMISSLEILFAAGDRLSSQDAILARRAVGSGV-Y-NAYGPTENTVLS
```

```
YAFDLGYTCMFPVLLGGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTA
-AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYI
IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPPALLKQCLVSA
```

```
SFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTE-FINHYGPTEATIGA
-YEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS
----PTMISSLEILFAAGDRLSSQDAILARRAVGSGV-Y-NAYGPTENTVLS
```



Inspiring Excellence

Multiple Sequence Alignment

A	T	-	G	T	T	a	T	A	
A	g	C	G	a	T	C	-	A	
A	T	C	G	T	-	C	T	c	
0	1	2	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	7	8
0	1	2	3	4	5	5	6	7	8

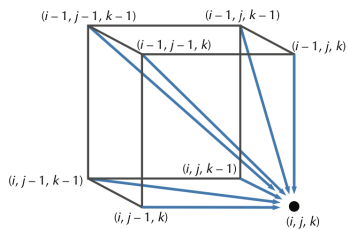
$(0, 0, 0) \rightarrow (1, 1, 1) \rightarrow (2, 2, 2) \rightarrow (2, 3, 3) \rightarrow (3, 4, 4) \rightarrow (4, 5, 5) \rightarrow (5, 6, 5) \rightarrow$
 $(6, 7, 6) \rightarrow (7, 7, 7) \rightarrow (8, 8, 8)$

- A multiple alignment of t strings v_1, \dots, v_t , also called a t -way alignment, is specified by a matrix having t rows, where the i -th row contains the symbols of v_i in order, interspersed with space symbols.



Inspiring Excellence

Multiple Sequence Alignment



$$s_{i,j,k} = \max \begin{cases} s_{i-1,j,k} & + \text{Score}(v_i, -, -) \\ s_{i,j-1,k} & + \text{Score}(-, w_j, -) \\ s_{i,j,k-1} & + \text{Score}(-, -, u_k) \\ s_{i-1,j-1,k} & + \text{Score}(v_i, w_j, -) \\ s_{i-1,j,k-1} & + \text{Score}(v_i, -, u_k) \\ s_{i,j-1,k-1} & + \text{Score}(-, w_j, u_k) \\ s_{i-1,j-1,k-1} & + \text{Score}(v_i, w_j, u_k) \end{cases}$$

- In the case of t sequences of length n , the alignment graph consists of approximately n^t nodes, and each node has up to $2^t - 1$ incoming edges, yielding a total runtime of $O(n^t 2^t)$
- As t grows, the dynamic programming algorithm becomes impractical.