

Sequencing Alignment I

Lectures 16 – Nov 21, 2011
CSE 527 Computational Biology, Fall 2011
Instructor: Su-In Lee
TA: Christopher Miles
Monday & Wednesday 12:00-1:20
Johnson Hall (JHN) 022

1

Outline: Sequence Alignment

- What
- Why (applications)
 - Comparative genomics
 - DNA sequencing
- A simple algorithm
- Complexity analysis
- A better algorithm:
 - “Dynamic programming”

2

Sequence Alignment: What

- Definition
 - An arrangement of two or several biological sequences (e.g. protein or DNA sequences) highlighting their similarity
 - The sequences are padded with gaps (usually denoted by dashes) so that columns contain identical or similar characters from the sequences involved
- Example – pairwise alignment

T A C T A A G

T C C A A T

3

Sequence Alignment: What

- Definition
 - An arrangement of two or several biological sequences (e.g. protein or DNA sequences) highlighting their similarity
 - The sequences are padded with gaps (usually denoted by dashes) so that columns contain identical or similar characters from the sequences involved
- Example – pairwise alignment

T A C T A A G

| : | : | | :

T C C – A A T

4

Sequence Alignment: Why

- The most basic sequence analysis task
 - First aligning the sequences (or parts of them) and
 - Then deciding whether that alignment is more likely to have occurred because the sequences are related, or just by chance
- Similar sequences often have similar origin or function
- New sequence always compared to existing sequences (e.g. using BLAST)

5

Sequence Alignment

- Example: gene HBB
 - Product: hemoglobin
 - Sickle-cell anaemia causing gene
 - Protein sequence (146 aa)

```
MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAMGNPK  
VKAHGKKV LGGAFSDGLAHLDNLKGTFATLS ELHCDKLHVD PENFRLLGNV LVCVLAHHFG  
KEFTTPVQAA YQKVVAGVAN ALAHKYH
```

- BLAST (Basic Local Alignment Search Tool)
 - The most popular alignment tool
 - Try it! Pick any protein, e.g. hemoglobin, insulin, exportin,... BLAST to find distant relatives.
 - <http://www.ncbi.nlm.nih.gov/blast/>

6

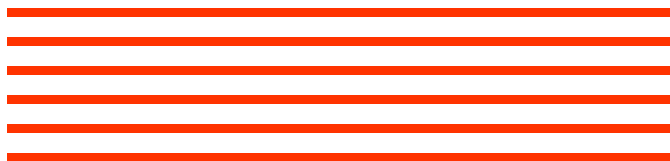
Sequence Alignment: Why

- The most basic sequence analysis task
 - First aligning the sequences (or parts of them) and
 - Then deciding whether that alignment is more likely to have occurred because the sequences are related, or just by chance
- Similar sequences often have similar origin or function
- New sequence always compared to existing sequences (e.g. using BLAST)
- Alignment algorithm is also used when sequencing DNA (reading DNA sequence)

7

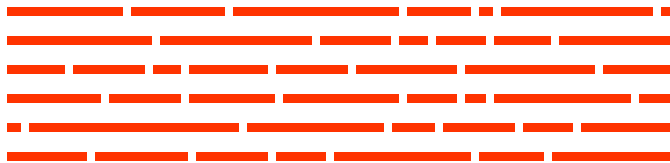
How Do We Read DNA?

- We replicate it

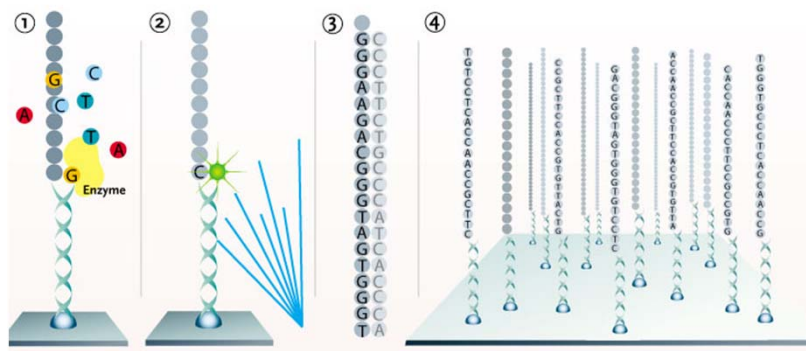


How Do We Read DNA?

- We replicate it
- We shred it



Reading Short DNA



- Use replication machinery with colored bases
- Take pictures of massively parallel reaction
- 10 million reads of 30 per day & \$1000

10

Example

The genome is:

TTATGGTCGGTGAGTGTGACTGGTGTTGTCTAA

The reads are:

GGTCGGTGAG
TGAGTGTGAC
TGGTGTTGTC
TGACTGGTTT
AATGGTCGGT
GAGTGTGACT
AAAAAAAAAA

11

Example

The genome is:

TTATGGTCGGTGAGTGTGACTGGTGTTGTCTAA

|||||
GGTCGGTGAG

The reads are:

GGTCGGTGAG
TGAGTGTGAC
TGGTGTTGTC
TGACTGGTTT
AATGGTCGGT
GAGTGTGACT
AAAAAAAAAA

12

Example

The genome is:

TTATGGTCGGT**TGAGTGTGAC**TGGTGTTGTCCTAA
 | | | | | | | |
 TGAGTGTGAC

The reads are:

GGTCGGTGAG
TGAGTGTGAC
TGGTGTTGTC
TGACTGGTTT
AATGGTCGGT
GAGTGTGACT
AAAAAAAAAA

13

Example

The genome is:

TTATGGTCGGTGAGTGTGAC**TGGTGTTGTC**CTAA
 | | | | | | | |
 TGGTGTTGTC

The reads are:

GGTCGGTGAG
TGAGTGTGAC
TGGTGTTGTC
TGACTGGTTT
AATGGTCGGT
GAGTGTGACT
AAAAAAAAAA

14

Example

The genome is:

TTATGGTCGGTGAGTGTGACTGGTGTGTGTCTAA
 | | | | | | | |
 TGACTGGTTT

The reads are:

GGTCGGTGAG
TGAGTGTGAC
TGGTGTTGTC
TGACTGGTTT
AATGGTCGGT
GAGTGTGACT
AAAAAAAAAA

15

Example

The genome is:

TATGGTCGGTGAGTGTGACTGGTGTGTGTCTAA
 | | | | | | | |
 AATGGTCGGT

The reads are:

GGTCGGTGAG
TGAGTGTGAC
TGGTGTTGTC
TGACTGGTTT
AATGGTCGGT
GAGTGTGACT
AAAAAAAAAA

16

Example

The genome is:

```
TTATGGTCGGTGAGTGTGACTGGTGTTGTCTAA
      GGTCGGTGAG
            TGAGTGTGAC
                  TGGTGTTGTC
                        TGACTGGTTT
AATGGTCGGT
            GAGTGTGACT
```

Assembled genome sequence is:

AATGGTCGGTGAGTGTGACTGGT**TT**TTGTC

17

Computational Problem

- **Goal:** Align reads to genome
- **Input:**
Many reads: l -long strings S_1, \dots, S_m
Approximate reference genome: string R
- **Output:**
 x_1, \dots, x_m along R where reads match, resp.
- **Complications:**
 - Errors
 - Differences

18

Sequence Alignment

- What
- Why (applications)
 - Comparative genomics
 - DNA sequencing
- A simple algorithm
- Complexity analysis
- A better algorithm:
 - “Dynamic programming”

19

Sequence Alignment: Key Issues

- What sorts of alignment should be considered
- The scoring system used to rank alignments
- The algorithm used to find optimal (or good) scoring alignments
- The statistical methods used to evaluate the significance of an alignment score

20

Terminology (CS, not necessarily Bio)

- *String*: ordered list of letters TATAAG
- *Prefix*: consecutive letters from front
empty, T, TA, TAT, ...
- *Suffix*: ... from end
empty, G, AG, AAG, ...
- *Substring*: ... from ends or middle
empty, TAT, AA, ...
- *Subsequence*: ordered, nonconsecutive
TT, AAA, TAG, ...

21

Sequence Alignment

S	a	c	b	c	d	b	
T		c	a	d	b	d	

S'	a	c	-	-	b	c	d	b
T'	-	c	a	d	b	-	d	-

Definition: An *alignment* of strings S , T is a pair of strings S' , T' (with spaces) s.t.

- (1) $|S'| = |T'|$, and $(|S| = \text{"length of } S\text{"})$
- (2) removing all spaces leaves S , T

22

Alignment Scoring

Mismatch = -1
Match = 2

S a c b c d b S' a c - - b c d b
T c a d b d T' - c a d b - d -

-1 2 -1 -1 2 -1 2 -1

Value = 3*2 + 5*(-1) = +1

- The *score* of aligning (characters or spaces) x & y is $\sigma(x,y)$.
- *Value* of an alignment $\sum_{i=1}^{|S'|} \sigma(S'[i], T'[i])$
- An *optimal alignment*: one of max value

23

Optimal Alignment: A Simple Algorithm

for all subseqs A of S, B of T s.t. $|A| = |B|$ **do**
 align A[i] with B[i], $1 \leq i \leq |A|$
 align all other chars to spaces
 compute its value
 retain the max
end

output the retained alignment

Example

S = abcd → A = cd
 T = wxyz → B = xz
 -ab**c**-**d** a-b**c**-**d**
 w--**x****y****z** -w-**x****y****z**

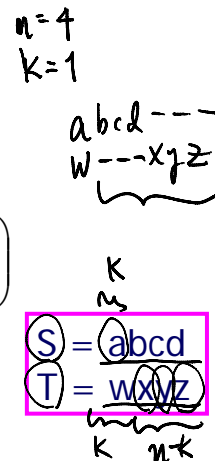
Outline: Sequence Alignment

- What
- Why (applications)
 - Comparative genomics
 - DNA sequencing
- A simple algorithm
- Complexity analysis
- A better algorithm:
 - "Dynamic programming"

25

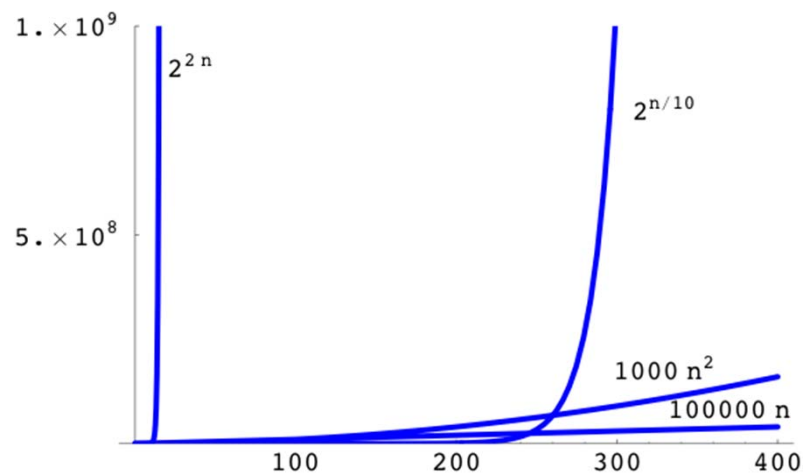
Complexity Analysis

- Assume $|S| = |T| = n$
- Cost of evaluating one alignment: $\geq n$
- How many alignments are there: $\sum_{k=0}^n \binom{2n}{n}$
 - pick n chars of S, T together ✓
 - say k of them are in S
 - match these k to the k unpicked chars of T
- Total time: $\geq n \binom{2n}{n} > 2^{2n}$, for $n > 3$
- E.g., for $n = 20$, time is $> 2^{40}$ operations



26

Polynomial vs exponential growth



27

Outline: Sequence Alignment

- What
- Why (applications)
 - Comparative genomics
 - DNA sequencing
- A simple algorithm
- Complexity analysis
- A better algorithm:
 - "Dynamic programming"

28

Alignment Scoring

Mismatch = -1
Match = 2

S a c b c d b S' a c - - b c d b
T c a d b d T' - c a d b - d -
 -1 2 -1 -1 2 -1 2 -1

$$\text{Value} = 3 \times 2 + 5 \times (-1) = +1$$

- The *score* of aligning (characters or spaces) x & y is $\sigma(x,y)$: e.g. $\sigma(a,-) = -1$, $\sigma(c,c) = 2$.
- *Value* of an alignment $\sum_{i=1}^{|S'|} \sigma(S'[i], T'[i])$
- An *optimal alignment*: one of max value
- A simple algorithm: complexity $> 2^{2n}$

29

Needleman-Wunsch Algorithm

- Align by "Dynamic programming"
- Key idea: Build up an optimal alignment using previous solutions for *optimal alignments of smaller subsequences*.
- Optimal alignment between S & T *ends* in 1 of 3 ways:
 - last chars of S & T aligned with each other
 - last char of S aligned with space in T
 - last char of T aligned with space in S
 - (never align space with space; $\sigma(-, -) < 0$)
 - In each case, the *rest* of S & T should be *optimally* aligned to each other

$\begin{bmatrix} \sim \sim \sim S[n] \\ \sim \sim \sim T[m] \end{bmatrix}$,
 $\begin{bmatrix} \sim \sim \sim S[n] \\ \sim \sim \sim - \end{bmatrix}$, or
 $\begin{bmatrix} \sim \sim \sim - \\ \sim \sim \sim T[m] \end{bmatrix}$
 Opt align of $S_1 \dots S_{n-1}$ & $T_1 \dots T_{m-1}$ Opt align of $S_1 \dots S_{n-1}$ & $T_1 \dots T_m$ Opt align of $S_1 \dots S_n$ & $T_1 \dots T_{m-1}$

Optimal Alignment in $O(n^2)$ via “Dynamic Programming”

- Input: $S, T, |S| = n, |T| = m$
- Output: **value** of optimal alignment
- Easier to solve a “harder” problem:
 $V(i,j)$ = value of optimal alignment of
 $S[1], \dots, S[i]$ with $T[1], \dots, T[j]$
for **all** $0 \leq i \leq n, 0 \leq j \leq m$.

31

General Case

Optimal align of $S[1], \dots, S[i]$ vs $T[1], \dots, T[j]$:

$$\begin{bmatrix} \sim\sim\sim\sim & S[i] \\ \sim\sim\sim\sim & T[j] \end{bmatrix}, \begin{bmatrix} \sim\sim\sim\sim & S[i] \\ \sim\sim\sim\sim & - \end{bmatrix}, \text{ or } \begin{bmatrix} \sim\sim\sim\sim & - \\ \sim\sim\sim\sim & T[j] \end{bmatrix}$$

Opt align of
 $S_1 \dots S_{i-1}$ & $T_1 \dots T_{j-1}$
Value = $V(i-1, j-1)$

Opt align of
 $S_1 \dots S_{i-1}$ & $T_1 \dots T_j$
Value = $V(i-1, j)$

Opt align of
 $S_1 \dots S_i$ & $T_1 \dots T_{j-1}$
Value = $V(i, j-1)$

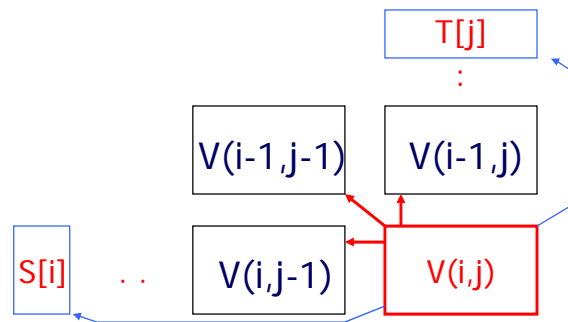
$$V(i,j) = \max \begin{cases} V(i-1,j-1) + \sigma(S[i], T[j]) \\ V(i-1,j) + \sigma(S[i], -) \\ V(i,j-1) + \sigma(-, T[j]) \end{cases},$$

for all $1 \leq i \leq n, 1 \leq j \leq m$.

32

Calculating One Entry

$$V(i,j) = \max \begin{cases} V(i-1,j-1) + \sigma(S[i], T[j]) \\ V(i-1,j) + \sigma(S[i], -) \\ V(i,j-1) + \sigma(-, T[j]) \end{cases}$$



33

Base Cases

- $V(i,0)$: first i chars of S all match spaces

$$V(i,0) = \sum_{k=1}^i \sigma(S[k], -)$$

- $V(0,j)$: first j chars of T all match spaces

$$V(0,j) = \sum_{k=1}^j \sigma(-, T[k])$$

34

Example

Mismatch = -1
Match = 2

- $V(i,j)$ = value of optimal alignment of $S[1], \dots, S[i]$ with $T[1], \dots, T[j]$

j	0	1	2	3	4	5	
i		c	a	d	b	d	←T
0	0	-1	-2	-3	-4	-5	
1	a	-1					
2	c	-2					
3	b	-3					
4	c	-4					
5	d	-5					
6	b	-6					
	↑S						

c
- Score(c,-) = -1

35

Example

Mismatch = -1
Match = 2

j	0	1	2	3	4	5	
i		c	a	d	b	d	←T
0	0	-1	-2	-3	-4	-5	
1	a	-1					
2	c	-2					
3	b	-3					
4	c	-4					
5	d	-5					
6	b	-6					
	↑S						

-
a Score(-,a) = -1

36

Example

Mismatch = -1
Match = 2

j	0	1	2	3	4	5
i		c	a	d	b	d
0	0	-1	-2	-3	-4	-5
1	a	-1				
2	c	-2				
3	b	-3				
4	c	-4				
5	d	-5				
6	b	-6				

←T

↑S

Score(-,c) = -1

37

Example

Mismatch = -1
Match = 2

T: ca
S: a

j	0	1	2	3	4	5
i		c	a	d	b	d
0	0	-1	-2	-3	-4	-5
1	a	-1	1			
2	c	-2				
3	b	-3				
4	c	-4				
5	d	-5				
6	b	-6				

←T

↑S

Score(-,c) = -1

Score(a,a) = +2

Score(a,-) = -1

Score(-,a) = -1

Score(c,a) = -1

Score(c,-) = -1

Score(b,-) = -1

Score(d,-) = -1

Score(-,d) = -1

Score(-,b) = -1

Score(-,c) = -1

Score(-,a) = -1

Score(-,d) = -1

Score(a,c) = -1

Score(a,b) = -1

Score(a,d) = -1

Score(c,b) = -1

Score(c,d) = -1

Score(b,d) = -1

Score(d,d) = -1

Score(d,b) = -1

Score(d,c) = -1

Score(d,a) = -1

Score(d,-) = -1

Score(b,c) = -1

Score(b,a) = -1

Score(b,d) = -1

Score(c,d) = -1

Score(a,d) = -1

Score(a,b) = -1

Score(a,c) = -1

Score(a,-) = -1

Score(-,d) = -1

Score(-,b) = -1

Score(-,c) = -1

Score(-,a) = -1

Score(-,-) = -1

38

Example

Mismatch = -1
Match = 2

	j	0	1	2	3	4	5	
i			c	a	d	b	d	←T
0		0	-1	-2	-3	-4	-5	
1	a	-1	-1	1				
2	c	-2	1					
3	b	-3						
4	c	-4						
5	d	-5						
6	b	-6						

↑
S

39

Example

Mismatch = -1
Match = 2

	j	0	1	2	3	4	5	
i			c	a	d	b	d	←T
0		0	-1	-2	-3	-4	-5	
1	a	-1	-1	1	0	-1	-2	
2	c	-2	1	0	0	-1	-2	
3	b	-3	0	0	-1	2	1	
4	c	-4	-1	-1	-1	1	1	
5	d	-5	-2	-2	1	0	3	
6	b	-6	-3	-3	0	3	2	

↑
S

Time = $O(mn)$

40

Finding Alignments: Trace Back

Arrows = (ties for) max in $V(i,j)$; 3 LR-to-UL paths = 3 optimal alignments

j	0	1	2	3	4	5	
i		c	a	d	b	d	←T
0	0	-1	-2	-3	-4	-5	
1	a	-1	-1	1	0	-1	-2
2	c	-2	1	0	0	-1	-2
3	b	-3	0	0	-1	2	1
4	c	-4	-1	-1	-1	1	1
5	d	-5	-2	-2	1	0	3
6	b	-6	-3	-3	0	3	2

↑
S

41

Complexity Notes

- Time = $O(mn)$, (value and alignment)
- Space = $O(mn)$
- Easy to get **value** in Time = $O(mn)$ and Space = $O(\min(m,n))$
- Possible to get value **and alignment** in Time = $O(mn)$ and Space = $O(\min(m,n))$ but tricky.

42

Significance of Alignments

- Is "42" a good score?
 - *Compared to what?*
- Usual approach: compared to a specific "null model", such as "random sequences"

43

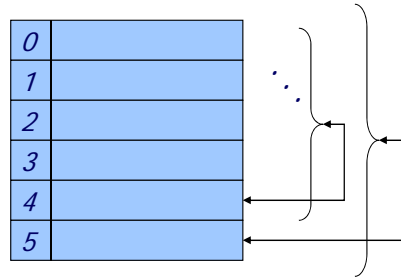
Overall Alignment Significance, II Empirical (via randomization)

- Generate N random sequences (say $N = 10^3 - 10^6$)
- Align x to each & score
- If k of them have better score than alignment of x to y, then the (empirical) probability of a chance alignment as good as observed x:y alignment is $(k+1)/(N+1)$
 - e.g., if 0 of 99 are better, you can say "estimated $p < .01$ "
- How to generate "random" sequences?
 - Scores are often sensitive to sequence composition
 - So uniform 1/20 or 1/4 is a bad idea
 - Even background p_i can be dangerous
 - Better idea: *permute* y N times

44

Generating Random Permutations

```
for (i = n-1; i > 0; i--){  
    j = random(0..i);  
    swap X[i] <-> X[j];  
}
```



All $n!$ permutations of the original data equally likely: Why? A specific element will be last with prob $1/n$; given that, a specific other element will be next-to-last with prob $1/(n-1)$, ...; overall: $1/(n!)$