

CSE 421 Intro to Algorithms Winter 2000

Sequence Alignment

CSE 421, W '00, Ruzzo

1

Sequence Alignment

- What
- Why
- A Simple Algorithm
- Complexity Analysis
- A better Algorithm:
"Dynamic Programming"

CSE 421, W '00, Ruzzo

2

Sequence Similarity: What

GGACCA

TACTAAG

|:|:|:|:
TCC-AAT

CSE 421, W '00, Ruzzo

3

Sequence Similarity: Why

- Diff
- RCS
- Molecular Bio
 - Similar sequences often have similar origin or function
 - Similarity often recognizable after $10^8 - 10^9$ years

CSE 421, W '00, Ruzzo

4

Terminology

- **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, ...

CSE 421, W '00, Ruzzo

5

Sequence Alignment

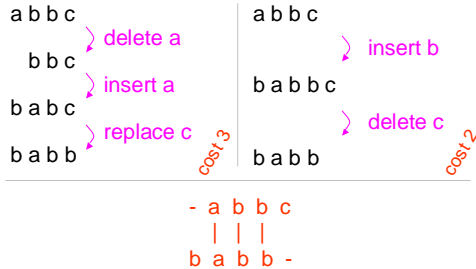
a c b c d b a c - - b c d b
 / \ / \ | | | |
c a d b d - c a d b - d -

- Defn:** An *alignment* of strings S, T is a pair of strings S', T' (with spaces) s.t.
- (1) $|S'| = |T'|$, and $(|S| = \text{"length of S"})$
 - (2) removing all spaces leaves S, T

CSE 421, W '00, Ruzzo

6

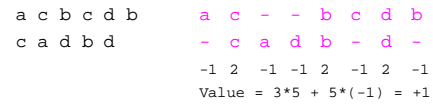
6.8: "Min_Edit_Distance"



CSE 421, W '00, Ruzzo

7

Alignment Scoring



- 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

CSE 421, W '00, Ruzzo

8

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

$S = abcd$ $A = cd$
 $T = wxyz$ $B = xz$
 $-abc-d$ $a-bc-d$
 $w--xyz$ $-w-xyz$

CSE 421, W '00, Ruzzo

9

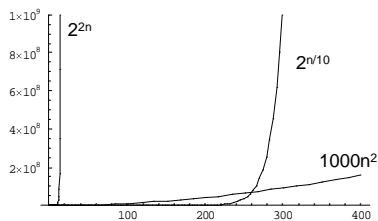
Analysis

- Assume $|S| = |T| = n$
- Cost of evaluating one alignment: $\geq n$
- How many alignments are there: $\geq \binom{2n}{n}$
 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

CSE 421, W '00, Ruzzo

10

Polynomial vs Exponential Growth



CSE 421, W '00, Ruzzo

11

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** $0 \leq i \leq n, 0 \leq j \leq m$.

CSE 421, W '00, Ruzzo

12

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

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

CSE 421, W '00, Ruzzo

13

General Case

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

$$\begin{array}{c}
 \left[\begin{array}{c} \sim \sim \sim S[i] \\ \sim \sim \sim T[j] \end{array} \right] \quad \left[\begin{array}{c} \sim \sim \sim S[i] \\ \sim \sim \sim - \end{array} \right] \quad \text{or} \quad \left[\begin{array}{c} \sim \sim \sim - \\ \sim \sim \sim T[j] \end{array} \right] \\
 \text{Opt align of } S_{i-1}, S_{i-1} \text{ \& } T_{j-1}, T_{j-1} \\
 V(i,j) = \max \left\{ \begin{array}{l} V(i-1, j-1) + \sigma(S[i], T[j]) \\ V(i-1, j) + \sigma(S[i], -) \\ V(i, j-1) + \sigma(-, T[j]) \end{array} \right\}
 \end{array}$$

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

CSE 421, W '00, Ruzzo

14

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	-1	1		
2	c	-2	1			
3	b	-3				
4	c	-4				
5	d	-5				
6	b	-6				

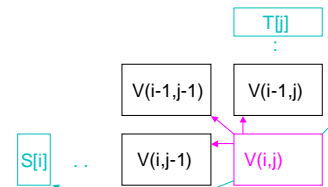
Time = $O(mn)$

CSE 421, W '00, Ruzzo

15

Calculating One Entry

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



CSE 421, W '00, Ruzzo

16

Finding Alignments: Trace Back

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

CSE 421, W '00, Ruzzo

17

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.

CSE 421, W '00, Ruzzo

18