

Objectives

Dynamic Programming: Computational Biology Applications

- Sequence Alignment in Linear Space

Mar 27, 2009

CS211

1

Edit Distance

[Levenshtein 1966, Needleman-Wunsch 1970]

- Gap penalty: δ
- Mismatch penalty: α_{pq}
 - If p and q are the same, then mismatch penalty is 0
- Cost = sum of gap and mismatch penalties

Parameters allow us to tweak cost

C	T	G	A	C	C	T	A	C	C	T	-	C	T	G	A	C	C	T	A	C	C	T
C	C	T	G	A	C	T	A	C	A	T	C	C	T	G	A	C	-	T	A	C	A	T
$\alpha_{TC} + \alpha_{GT} + \alpha_{AG} + 2\alpha_{CA}$											$2\delta + \alpha_{CA}$											

Mar 27, 2009

CS211

2

Sequence Alignment

Goal: Given two strings $X = x_1 x_2 \dots x_m$ and $Y = y_1 y_2 \dots y_n$ find alignment of minimum cost

An **alignment** M is a set of ordered pairs $x_i - y_j$ such that each item occurs in at most one pair and **no** crossings

The pair $x_i - y_j$ and $x_{i'} - y_{j'}$ **cross** if $i < i'$, but $j > j'$.

o	c	c	u	r	e	r	n	c	e	o	c	c	u	r	e	r	n	c	e
o	c	c	u	r	e	r	n	c	e	o	c	c	u	r	e	r	n	c	e

crossing 2 mismatches

Mar 27, 2009

CS211

3

Sequence Alignment Example

$X = \text{CTACCG}$

$Y = \text{TACTG}$

Solution: $M = x_2 - y_1, x_3 - y_2, x_4 - y_3, x_5 - y_4, x_6 - y_6$

x_1	x_2	x_3	x_4	x_5	x_6	
C	T	A	C	C	G	
	T	A	C	A	T	G
	y_1	y_2	y_3	y_4	y_5	y_6

What is the cost of M ?

$$\text{cost}(M) = \sum_{(x_i, y_j) \in M} \alpha_{x_i y_j} + \sum_{i: x_i \text{ unmatched}} \delta + \sum_{j: y_j \text{ unmatched}} \delta$$

mismatch gap

Recall: mismatch penalty is 0 if x_i and y_j are the same

Mar 27, 2009

CS211

4

Sequence Alignment Case Analysis

Consider last character of strings X and Y : x_M and y_N

- Case 1: x_M and y_N are aligned
- Case 2: x_M is not matched
- Case 3: y_N is not matched

Mar 27, 2009

CS211

5

Sequence Alignment Cost Analysis

Consider last character of strings X and Y : x_M and y_N

- Case 1: x_M and y_N are aligned
 - Pay mismatch for $x_M - y_N$ + min cost of aligning rest of strings
 - $\text{OPT}(M, N) = \alpha_{x_M y_N} + \text{OPT}(M-1, N-1)$
- Case 2: x_M is not matched
 - Pay gap for x_M + min cost of aligning rest of strings
 - $\text{OPT}(M, N) = \delta + \text{OPT}(M-1, N)$
- Case 3: y_N is not matched
 - Pay gap for y_N + min cost of aligning rest of strings
 - $\text{OPT}(M, N) = \delta + \text{OPT}(M, N-1)$

Mar 27, 2009

CS211

6

Sequence Alignment: Problem Structure

$$OPT(i, j) = \begin{cases} j\delta & \text{if } i = 0 \\ \min \begin{cases} \alpha_{x_i, y_j} + OPT(i-1, j-1) \\ \delta + OPT(i-1, j) \\ \delta + OPT(i, j-1) \end{cases} & \text{otherwise} \\ i\delta & \text{if } j = 0 \end{cases}$$

Gaps for remainder of Y

Gaps for remainder of X

Mar 27, 2009

CS211

7

Sequence Alignment: Analysis

```

Sequence-Alignment(m, n, x1x2...xm, y1y2...yn, δ, α)
  for i = 0 to m
    M[0, i] = iδ
  for j = 0 to n
    M[j, 0] = jδ

  for i = 1 to m
    for j = 1 to n
      M[i, j] = min(α[xi, yj] + M[i-1, j-1],
                  δ + M[i-1, j],
                  δ + M[i, j-1])
  return M[m, n]

```

$O(mn)$

Observation: to calculate the current value, we only need the row above us and the entry to the left

Mar 27, 2009

CS211

8

SEQUENCE ALIGNMENT IN LINEAR SPACE

9

Sequence Alignment: $O(m)$ Space

Collapse into an $m \times 2$ array

- $M[i, 0]$ represents previous row; $M[i, 1]$ -- current

```

Space-Efficient-Alignment(m, n, x1x2...xm, y1y2...yn, δ, α)
  for i = 0 to m
    # initialize first row
    M[i, 0] = iδ
  for j = 1 to n
    # first gap
    M[0, 1] = jδ

  for i = 1 to m
    M[i, 1] = min(α[xi, y1] + M[i-1, 0],
                δ + M[i, 0],
                δ + M[i-1, 1])

  for i = 1 to m
    # copy current row into previous
    M[i, 0] = M[i, 1]
  return M[m, 1]

```

Mar 27, 2009

CS211

Any drawbacks?

10

Sequence Alignment: $O(m)$ Space

Collapse into an $m \times 2$ array

- $M[i, 0]$ represents previous row; $M[i, 1]$ -- current

```

Space-Efficient-Alignment(m, n, x1x2...xm, y1y2...yn, δ, α)
  for i = 0 to m
    M[i, 0] = iδ
  for j = 1 to n
    M[0, 1] = jδ

  for i = 1 to m
    M[i, 1] = min(α[xi, y1] + M[i-1, 0],
                δ + M[i, 0],
                δ + M[i-1, 1])

  for i = 1 to m
    M[i, 0] = M[i, 1]
  return M[m, 1]

```

Mar 27, 2009

CS211

Finds optimal value but won't be able to find the alignment

11

Why Do We Care About Space?

For English words or sentences, probably doesn't matter

Matters for Biological sequence alignment

- Consider: 2 strings with 100,000 symbols each
 - Processor can do 10 billion primitive operations
 - BUT dealing with a 10 GB array

Mar 27, 2009

CS211

12

Sequence Alignment: Linear Space

Can we avoid using quadratic space?

- Optimal value in $O(m)$ space and $O(mn)$ time.
 - Compute $OPT(i, \cdot)$ from $OPT(i-1, \cdot)$
 - BUT, no longer a simple way to recover alignment itself

Theorem. [Hirschberg 1975] Optimal alignment in $O(m + n)$ space and $O(mn)$ time.

- Clever combination of divide-and-conquer and dynamic programming
- Inspired by idea of Savitch from complexity theory

Mar 27, 2009

CS211

13

Recall Our Example

$X = \text{bait}$

$Y = \text{boot}$

$\alpha = 1$, for vowel mismatch
 $\alpha = 2$, for other mismatches
 $\delta = 2$

$i=4$ j

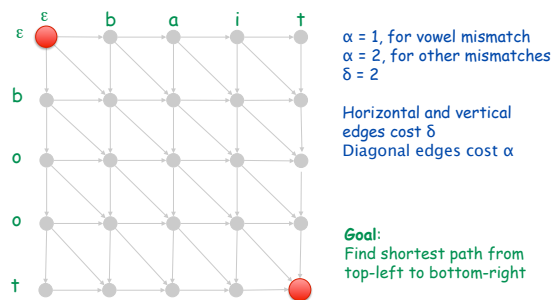
		b	a	i	t
	0	2	4	6	8
b	2	0	2	4	6
o	4	2	1	3	5
o	6	4	3	2	4
t	8	6	5	4	2

Mar 25, 2009

CS211

14

Mapping to a Graph Problem

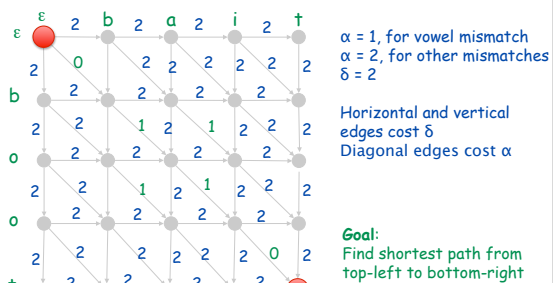


Mar 27, 2009

CS211

15

Mapping to a Graph Problem



Mar 27, 2009

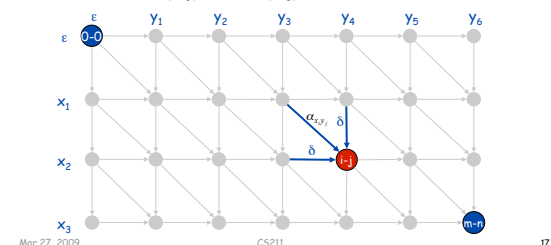
CS211

16

Sequence Alignment: Forward

Edit distance graph

- Let $f(i, j)$ be shortest path from $(0,0)$ to (i, j)
- Observation:** $f(i, j) = OPT(i, j)$



Mar 27, 2009

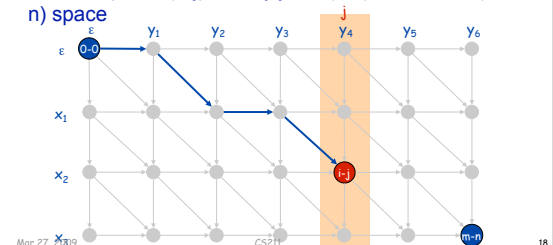
CS211

17

Sequence Alignment: Forward

Edit distance graph

- Let $f(i, j)$ be shortest path from $(0,0)$ to (i, j)
- Can compute $f(\cdot, j)$ for any j in $O(mn)$ time and $O(m + n)$ space



Mar 27, 2009

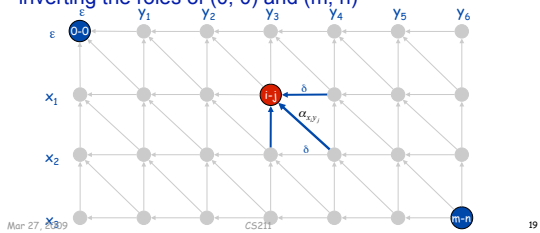
CS211

18

Sequence Alignment: Backward

Edit distance graph

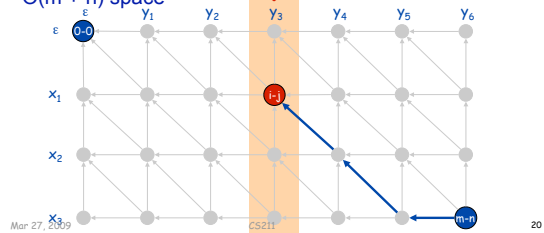
- Let $g(i, j)$ be shortest path from (m, n) to (i, j)
- Can compute by reversing the edge orientations and inverting the roles of $(0, 0)$ and (m, n)



Sequence Alignment: Backward

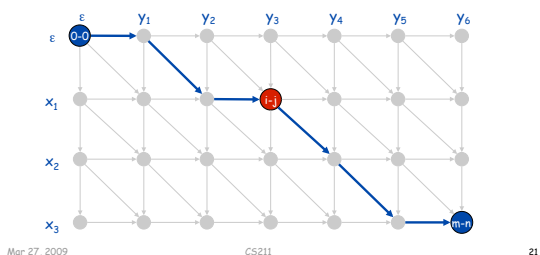
Edit distance graph

- Let $g(i, j)$ be shortest path from (m, n) to (i, j)
- Can compute $g(\cdot, j)$ for any j in $O(mn)$ time and $O(m + n)$ space



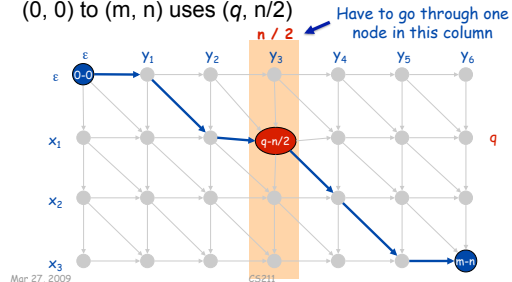
Sequence Alignment: Linear Space

Observation 1. The cost of the shortest path that uses (i, j) is $f(i, j) + g(i, j)$



Sequence Alignment: Linear Space

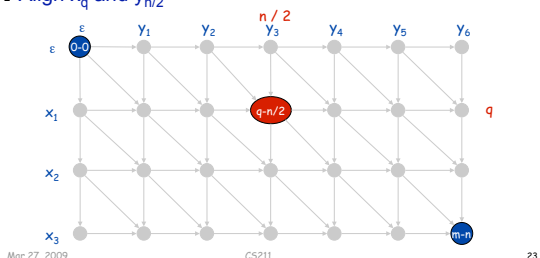
Observation 2. Let q be an index that minimizes $f(q, n/2) + g(q, n/2)$. Then, the shortest path from $(0, 0)$ to (m, n) uses $(q, n/2)$



Sequence Alignment: Linear Space

Divide: find index q that minimizes $f(q, n/2) + g(q, n/2)$ using DP

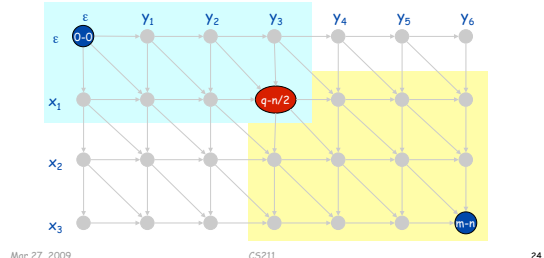
- Align x_q and $y_{n/2}$



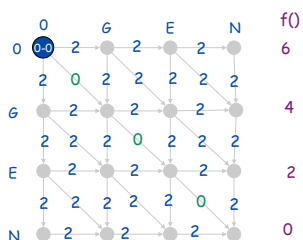
Sequence Alignment: Linear Space

Conquer: recursively compute optimal alignment in each piece

- Reuse working space from one recursive call to next



Space-efficient alignment: Left

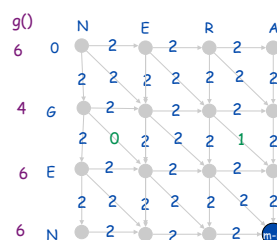


Mar 27, 2009

CS211

31

Backwards Space Efficient

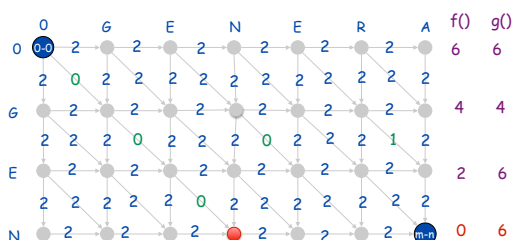


Mar 27, 2009

CS211

32

Example

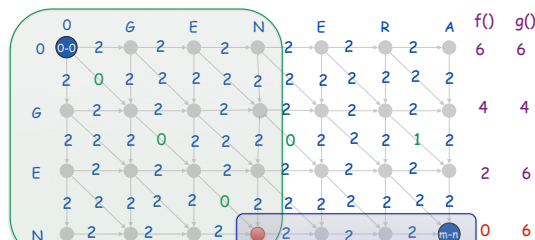


Mar 27, 2009

CS211

33

Example



Mar 27, 2009

CS211

34

Divide and Conquer Sequence Alignment: Analysis

P contains node on shortest corner-to-corner path
 Divide-and-Conquer-Alignment (X, Y)
 m = length of X
 n = length of Y
 if m ≤ 2 or n ≤ 2
 compute optimal alignment using Alignment(X, Y)
 return
 Space-Efficient-Alignment(X, Y[1:n/2])
 Backward-Space-Efficient-Alignment(X, Y[n/2+1:n])
 q = index that minimizes f(q, n/2) + g(q, n/2)
 add(q, n/2) to P
 Divide-and-Conquer-Alignment(X[1:q], Y[1:n/2])
 Divide-and-Conquer-Alignment(X[q:m], Y[n/2+1:n])
 return P

What is the recurrence relation?

Mar 27, 2009

CS211

35

Sequence Alignment: Running Time Analysis Warmup

Theorem. Let $T(m, n)$ = max running time of algorithm on strings of length at most m and n.
 $T(m, n) = O(mn \log n)$.

$$T(m, n) \leq 2T(m, n/2) + O(mn) \Rightarrow T(m, n) = O(mn \log n)$$

Remark. Analysis is not tight because two sub-problems are of size (q, n/2) and (m - q, n/2).

Mar 27, 2009

CS211

36

Sequence Alignment: Running Time Analysis

Theorem. Let $T(m, n)$ = max running time of algorithm on strings of length m and n . $T(m, n) = O(mn)$

Recurrence Relation:

$$\begin{aligned} T(m, 2) &\leq cm \\ T(2, n) &\leq cn \\ T(m, n) &\leq cmn + T(q, n/2) + T(m-q, n/2) \end{aligned}$$

Solve using substitution:

$$\begin{aligned} T(m, n) &\leq T(q, n/2) + T(m-q, n/2) + cmn \\ &\leq 2cqn/2 + 2c(m-q)n/2 + cmn \\ &= cqn + cmn - cqn + cmn \\ &= 2cmn \end{aligned}$$

Mar 27, 2009

CS211

37

Review: Problem Set 4

Recurrence Relations

Mar 27, 2009

CS211

38

1a. $T(n) = 5T(n/2) + O(n)$

Using formula we figured out $O(n^{\log_2 5})$ for $q > 2$:
 $O(n^{\log_2 5})$

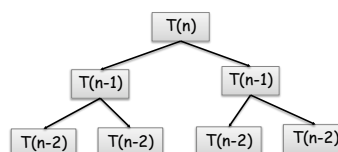
Mar 27, 2009

CS211

39

1b. $T(n) = 2T(n-1) + O(1)$

Unrolling recurrence



Most common mistake:
 $O(n)$

How many levels?
 How many problems at bottom level?
 What is the cost of each problem?

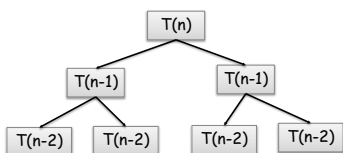
Mar 27, 2009

CS211

40

1b. $T(n) = 2T(n-1) + O(1)$

Unrolling recurrence



Most common mistake:
 $O(n)$

Recurrence: Like
 Fibonacci sequence

How many levels? $n-1$
 How many problems at bottom level? 2^n
 What is the cost of each problem? $c \cdot O(1)$

Therefore, bottom level is 2^n and $T(n) \in O(2^n)$

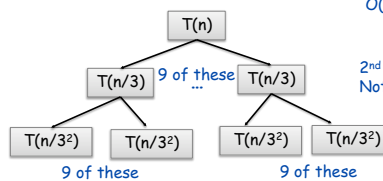
Mar 27, 2009

CS211

41

1c. $T(n) = 9T(n/3) + O(n^2)$

Unrolling the recurrence



Most common mistake:
 $O(n)$

2nd most common mistake:
 Not reducing problem

How many levels?
 How many problems at each level?

Mar 27, 2009

CS211

42

1c. $T(n) = 9T(n/3) + O(n^2)$

Unrolling the recurrence

Most common mistake: $O(n)$

How many levels? $\log_3 n$
 How many problems at each level? 9^l
 What is the cost of each level?

Mar 27, 2009 CS211 43

1c. $T(n) = 9T(n/3) + O(n^2)$

Unrolling the recurrence

Most common mistake: $O(n)$

How many levels? $\log_3 n$
 How many problems at each level? 9^l
 What is the cost of each level? cn^2

$O(n^2 \log_3 n)$

Mar 27, 2009 CS211 44

Problem Set 6

- 1 Standard DP Problem
- 1 DP Problem *and* program implementing algorithm
 - Goal: "pretty print" text so that the *slack* between text and right margin is minimized
 - Slack: number of spaces between text and right margin
 - Write a Python program that *pretty prints* text given a maximum length of a line
 - Template Python program, test files on Course web site

Mar 27, 2009

CS211

45

Unix Details

Python program (pp) must be *executable*

- `chmod u+x pp`
 - Make pp executable by user

Data file and example output are in tgz format

- Tar, gzip
- Need to unzip and extract file

Mar 27, 2009

CS211

46

Unix Details

Python program (pp) must be *executable*

- `chmod u+x pp`
 - Make pp executable by user

Data file and example output are in tgz format

- Tar, gzip
- Need to unzip and extract file
- `tar xzf neruda.tgz`

Mar 27, 2009

CS211

47