

Objectives

- Dynamic Programming
 - Sequence Alignment
 - Shortest Path

Mar 25, 2013

CSCI211 - Sprenkle

1

Review

- What is the sequence alignment problem?

Mar 25, 2013

CSCI211 - Sprenkle

2

String Similarity

- How similar are two strings?

- occurrence
- occurrence

- Measurements

- Gap (-): add a letter
- Mismatch

Which is the best alignment?

o c u r r a n c e -
o c c u r r e n c e

6 mismatches, 1 gap

o c - u r r a n c e
o c c u r r e n c e

1 mismatch, 1 gap

o c - u r r - a n c e
o c c u r r e - n c e

0 mismatches, 3 gaps

Mar 25, 2013

CSCI211 - Sprenkle

3

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 25, 2013

CSCI211 - Sprenkle

4

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

crossing

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

2 mismatches

Mar 25, 2013

CSCI211 - Sprenkle

5

Sequence Alignment Example

- $X = \text{CTACCG}$
- $Y = \text{TACATG}$
- 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
 $y_1 \ y_2 \ y_3 \ y_4 \ y_5 \ y_6$
- T A C A T G

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

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

Mar 25, 2013

CSCI211 - Sprenkle

6

Sequence Alignment Case Analysis

- Consider the last character of the strings X and Y: x_M and y_N
 - M and N are not necessarily equal
 - i.e., strings are not necessarily the same length
- What are the possibilities for x_M and y_N in terms of the alignment?



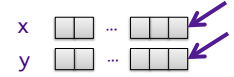
Mar 25, 2013

CSCI211 - Sprenkle

7

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



Formulate the optimal solution's value

Mar 25, 2013

CSCI211 - Sprenkle

8

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

What are the costs for these cases?



- $OPT(i, j) = \text{min cost of aligning strings } x_1 x_2 \dots x_i \text{ and } y_1 y_2 \dots y_j$

Mar 25, 2013

CSCI211 - Sprenkle

9

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

Mar 25, 2013

CSCI211 - Sprenkle

10

Sequence Alignment Cost Analysis

- Base costs? $\rightarrow i$ or j is 0
 - What happens when we run out of letters in one string before the other?

X = CTACCG
Y = TACTG

Mar 25, 2013

CSCI211 - Sprenkle

11

Sequence Alignment: Problem Structure

Gaps for remainder of Y

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

Ran out of 1st string

Ran out of 2nd string

Gaps for remainder of X

Mar 25, 2013

CSCI211 - Sprenkle

12

Sequence Alignment: Algorithm

Cost parameters

```

Sequence-Alignment( $m, n, x_1x_2\dots x_m, y_1y_2\dots y_n, \delta, \alpha$ )
  for  $i = 0$  to  $m$ 
     $M[i, 0] = i\delta$ 
  for  $j = 0$  to  $n$ 
     $M[0, j] = j\delta$ 

  for  $i = 1$  to  $m$ 
    for  $j = 1$  to  $n$ 
       $M[i, j] = \min(\alpha[x_i, y_j] + M[i-1, j-1],$ 
                     $\delta + M[i-1, j],$ 
                     $\delta + M[i, j-1])$ 

  return  $M[m, n]$ 

```

Mar 25, 2013

CSCI211 - Sprenkle

13

Example

X = bait**Y = boot**

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

		b	a	i	t
b					
o					
o					
t					

Mar 25, 2013

CSCI211 - Sprenkle

14

Example

X = bait**Y = boot**

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

		b	a	i	t
b	0	2	4	6	8
o	2				
o	4				
t	6				
t	8				

Mar 25, 2013

CSCI211 - Sprenkle

15

Example

X = bait**Y = boot**

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

		b	a	i	t
b	0	2	4	6	8
o	2	0	2	4	6
o	4				
o	6				
t	8				

Mar 25, 2013

CSCI211 - Sprenkle

16

Example

X = bait**Y = boot**

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

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

Mar 25, 2013

CSCI211 - Sprenkle

17

Example

X = bait**Y = boot**

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

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

Mar 25, 2013

CSCI211 - Sprenkle

18

Example

What is the value for the problem?
What is the solution?

X = bait

Y = boot

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

			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, 2013

CSCI211 - Sprenkle

19

Example

X = bait

Y = boot

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

			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, 2013

CSCI211 - Sprenkle

20

Sequence Alignment: Analysis

```

Sequence-Alignment(m, n, x1x2...xm, y1y2...yn,  $\delta$ ,  $\alpha$ )
  for i = 0 to m
    M[0, i] = i $\delta$ 
  for j = 0 to n
    M[j, 0] = j $\delta$ 

  for i = 1 to m
    for j = 1 to n
      M[i, j] = min( $\alpha[x_i, y_j] + M[i-1, j-1]$ ,
                     $\delta + M[i-1, j]$ ,
                     $\delta + M[i, j-1]$ )
  return M[m, n]

```

O(mn)

Costs?

Mar 25, 2013

CSCI211 - Sprenkle

21

Sequence Alignment: Algorithm

```

Sequence-Alignment(m, n, x1x2...xm, y1y2...yn,  $\delta$ ,  $\alpha$ )
  for i = 0 to m
    M[0, i] = i $\delta$ 
  for j = 0 to n
    M[j, 0] = j $\delta$ 

  for i = 1 to m
    for j = 1 to n
      M[i, j] = min( $\alpha[x_i, y_j] + M[i-1, j-1]$ ,
                     $\delta + M[i-1, j]$ ,
                     $\delta + M[i, j-1]$ )
  return M[m, n]

```

What are the space costs?

When computing M[i,j], which entries in M are used?

Mar 25, 2013

CSCI211 - Sprenkle

22

Sequence Alignment: Analysis

```

Sequence-Alignment(m, n, x1x2...xm, y1y2...yn,  $\delta$ ,  $\alpha$ )
  for i = 0 to m
    M[0, i] = i $\delta$ 
  for j = 0 to n
    M[j, 0] = j $\delta$ 

  for i = 1 to m
    for j = 1 to n
      M[i, j] = min( $\alpha[x_i, y_j] + M[i-1, j-1]$ ,
                     $\delta + M[i-1, j]$ ,
                     $\delta + M[i, j-1]$ )
  return M[m, n]

```

Space Cost: O(mn)

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

Mar 25, 2013

CSCI211 - Sprenkle

23

SEQUENCE ALIGNMENT IN
LINEAR SPACE

Mar 25, 2013

CSCI211 - Sprenkle

24

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( $m, n, x_1x_2\dots x_m, y_1y_2\dots y_n, \delta, \alpha$ )
  for  $i = 0$  to  $m$            # initialize first row
     $M[i, 0] = i\delta$ 
  for  $j = 1$  to  $n$ 
     $M[0, j] = j\delta$          # first gap

    for  $i = 1$  to  $m$ 
       $M[i, 1] = \min(\alpha[x_i, y_j] + M[i-1, 0],$ 
                     $\delta + M[i, 0],$ 
                     $\delta + M[i-1, 1])$ 

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

```

Any drawbacks?

Mar 25, 2013

CSCI211 - Sprenkle

25

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( $m, n, x_1x_2\dots x_m, y_1y_2\dots y_n, \delta, \alpha$ )
  for  $i = 0$  to  $m$            # initialize first row
     $M[i, 0] = i\delta$ 
  for  $j = 1$  to  $n$ 
     $M[0, j] = j\delta$          # first gap

    for  $i = 1$  to  $m$ 
       $M[i, 1] = \min(\alpha[x_i, y_j] + M[i-1, 0],$ 
                     $\delta + M[i, 0],$ 
                     $\delta + M[i-1, 1])$ 

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

```

Finds optimal value but will not be able to find alignment

Mar 25, 2013

CSCI211 - Sprenkle

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 25, 2013

CSCI211 - Sprenkle

27

Sequence Alignment: Linear Space

- Can we avoid using quadratic space?
 - Optimal value in $O(m)$ space and $O(mn)$ time.
 - Compute $\text{OPT}(i, \cdot)$ from $\text{OPT}(i-1, \cdot)$
 - BUT, no 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*
 - Section 6.7

Mar 25, 2013

CSCI211 - Sprenkle

28

IMPROVING SHORTEST PATH

Mar 25, 2013

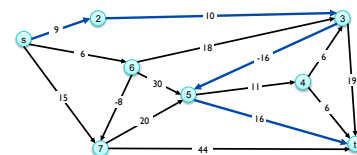
CSCI211 - Sprenkle

29

Shortest Paths

- Problem:** Given a directed graph $G = (V, E)$, with edge weights c_{vw} , find shortest path from node s to node t
 - allow negative weights

- Allows modeling other phenomena



Mar 25, 2013

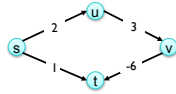
CSCI211 - Sprenkle

30

Shortest Paths: Failed Attempts

- Review: What was Dijkstra's algorithm?
 - Dijkstra can fail if negative edge costs

Shortest path from $s \rightarrow t$?



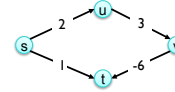
Mar 25, 2013

CSCI211 - Sprenkle

31

Shortest Paths: Failed Attempts

- Dijkstra. Can fail if negative edge costs



- Re-weighting. Adding a constant to every edge weight can fail

Why?

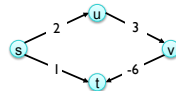
Mar 25, 2013

CSCI211 - Sprenkle

32

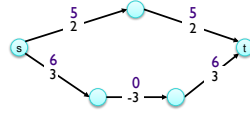
Shortest Paths: Failed Attempts

- Dijkstra. Can fail if negative edge costs



- Re-weighting. Adding a constant to every edge weight can fail

Why?



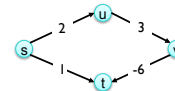
Mar 25, 2013

CSCI211 - Sprenkle

33

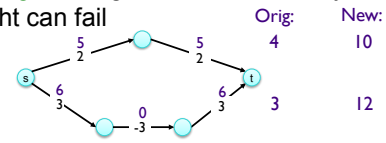
Shortest Paths: Failed Attempts

- Dijkstra. Can fail if negative edge costs



- Re-weighting. Adding a constant to every edge weight can fail

Why?

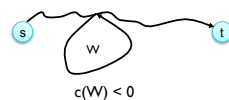
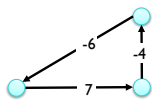


Mar 25, 2013

CSCI211 - Sprenkle

34

Shortest Paths: Negative Cost Cycles



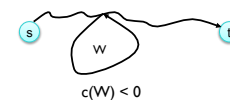
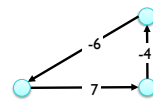
- If some path from s to t contains a negative cost cycle, there does **not** exist a shortest s - t path

Why?

- Otherwise, there exists one that is *simple* (i.e., does not repeat nodes)

What does this mean about number of edges in path?

Shortest Paths: Negative Cost Cycles



- If some path from s to t contains a negative cost cycle, there does **not** exist a shortest s - t path

- Otherwise, there exists one that is *simple* (i.e., does not repeat nodes)

➤ Path has at most $n-1$ edges, where n is # of nodes in graph

Mar 25, 2013

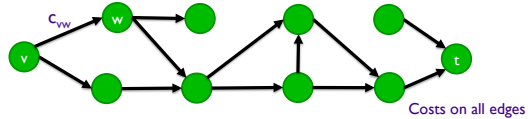
CSCI211 - Sprenkle

36

Towards a Recurrence

- $\text{OPT}(i, v)$: minimum cost of a v - t path P using **at most i** edges
 - This formulation eases later discussion
- Original problem is $\text{OPT}(n-1, s)$

Break down into subproblems based on i and v



Mar 25, 2013

CSCI211 - Sprenkle

37

Looking Ahead

- Wiki due Tuesday
 - Chap 6: 6.1-6.4
- Wednesday: Network Flows
- Thurs: CS Dept talk at 12:10
 - EC: 4 pts towards problem set grade for answering questions on Sakai
- PS8 due Friday

Mar 25, 2013

CSCI211 - Sprenkle

38