

CSE 421 Algorithms

Richard Anderson
Lecture 20
Memory Efficient Longest Common
Subsequence

Longest Common Subsequence

- $C=c_1\dots c_g$ is a subsequence of $A=a_1\dots a_m$ if C can be obtained by removing elements from A (but retaining order)
- $LCS(A, B)$: A maximum length sequence that is a subsequence of both A and B

ocurranec **attacggct**
occurrence **tacgacca**

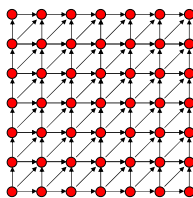
LCS Optimization

- $A = a_1a_2\dots a_m$
- $B = b_1b_2\dots b_n$
- $Opt[j, k]$ is the length of $LCS(a_1a_2\dots a_j, b_1b_2\dots b_k)$

Optimization recurrence

- If $a_j = b_k$, $Opt[j, k] = 1 + Opt[j-1, k-1]$
- If $a_j \neq b_k$, $Opt[j, k] = \max(Opt[j-1, k], Opt[j, k-1])$

Dynamic Programming Computation

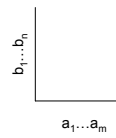


Storing the path information

```

A[1..m], B[1..n]
for i := 1 to m  Opt[i, 0] := 0;
for j := 1 to n  Opt[0, j] := 0;
Opt[0, 0] := 0;
for i := 1 to m
  for j := 1 to n
    if A[i] = B[j] { Opt[i, j] := 1 + Opt[i-1, j-1]; Best[i, j] := Diag; }
    else if Opt[i-1, j] >= Opt[i, j-1]
      { Opt[i, j] := Opt[i-1, j]; Best[i, j] := Left; }
    else { Opt[i, j] := Opt[i, j-1]; Best[i, j] := Down; }

```



Algorithm Performance

- $O(nm)$ time and $O(nm)$ space
- On current desktop machines
 - $n, m < 10,000$ is easy
 - $n, m > 1,000,000$ is prohibitive
- Space is more likely to be the bounding resource than time

Observations about the Algorithm

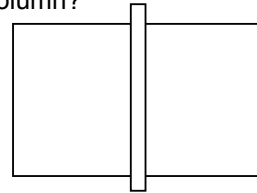
- The computation can be done in $O(m+n)$ space if we only need one column of the Opt values or Best Values
- The algorithm can be run from either end of the strings

Computing LCS in $O(nm)$ time and $O(n+m)$ space

- Divide and conquer algorithm
- Recomputing values used to save space

Divide and Conquer Algorithm

- Where does the best path cross the middle column?



- For a fixed i , and for each j , compute the LCS that has a_i matched with b_j

Constrained LCS

- $LCS_{i,j}(A,B)$: The LCS such that
 - a_1, \dots, a_i paired with elements of b_1, \dots, b_j
 - a_{i+1}, \dots, a_m paired with elements of b_{j+1}, \dots, b_n
- $LCS_{4,3}(abbacbb, cbbaa)$

A = **RR**SS**RT**TRTS
B = RTSRRSTST

Compute $LCS_{5,0}(A,B)$, $LCS_{5,1}(A,B)$, ..., $LCS_{5,9}(A,B)$

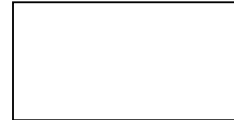
A = RRSSRTTRTS
 B = RTSRRSTST

Compute $LCS_{5,0}(A,B)$, $LCS_{5,1}(A,B)$, ..., $LCS_{5,9}(A,B)$

| j | left | right |
|---|------|-------|
| 0 | 0 | 4 |
| 1 | 1 | 4 |
| 2 | 1 | 3 |
| 3 | 2 | 3 |
| 4 | 3 | 3 |
| 5 | 3 | 2 |
| 6 | 3 | 2 |
| 7 | 3 | 1 |
| 8 | 4 | 1 |
| 9 | 4 | 0 |

Computing the middle column

- From the left, compute $LCS(a_1 \dots a_{m/2}, b_1 \dots b_j)$
- From the right, compute $LCS(a_{m/2+1} \dots a_m, b_{j+1} \dots b_n)$
- Add values for corresponding j's



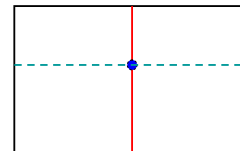
- Note – this is space efficient

Divide and Conquer

- $A = a_1, \dots, a_m$ $B = b_1, \dots, b_n$
- Find j such that
 - $LCS(a_1 \dots a_{m/2}, b_1 \dots b_j)$ and
 - $LCS(a_{m/2+1} \dots a_m, b_{j+1} \dots b_n)$ yield optimal solution
- Recurse

Algorithm Analysis

- $T(m,n) = T(m/2, j) + T(m/2, n-j) + cnm$



Prove by induction that
 $T(m,n) \leq 2cmn$

Memory Efficient LCS Summary

- We can afford $O(nm)$ time, but we can't afford $O(nm)$ space
- If we only want to compute the length of the LCS, we can easily reduce space to $O(n+m)$
- Avoid storing the value by recomputing values
 - Divide and conquer used to reduce problem sizes

Shortest Path Problem

- Dijkstra's Single Source Shortest Paths Algorithm
 - $O(m \log n)$ time, positive cost edges
- General case – handling negative edges
- If there exists a negative cost cycle, the shortest path is not defined
- Bellman-Ford Algorithm
 - $O(mn)$ time for graphs with negative cost edges

Lemma

- If a graph has no negative cost cycles, then the **shortest** paths are **simple** paths
- Shortest paths have at most $n-1$ edges

Shortest paths with a fixed number of edges

- Find the shortest path from v to w with exactly k edges

Express as a recurrence

- $\text{Opt}_k(w) = \min_x [\text{Opt}_{k-1}(x) + c_{xw}]$
- $\text{Opt}_0(w) = 0$ if $v=w$ and infinity otherwise

Algorithm, Version 1

```
foreach w
  M[0, w] = infinity;
M[0, v] = 0;
for i = 1 to n-1
  foreach w
    M[i, w] = min_x (M[i-1, x] + cost[x, w]);
```