

Administrative

CS 361, Lecture 26

Jared Saia
University of New Mexico

- This week and next, you can get one extra participation check by going to section and informing Nate that you're there and want a check.
- Sections are Thursday 5:30-6:20 DSH 134 and Friday 1:00-1:50 TAPY 218
- Good chance to get info on hw, projects and to review material for final

2

Outline

Administrative

- Administrative
- Dynamic Programming
- Course Evaluations

- The project is due this Thursday in class.
- This deadline is strict - late projects will receive no credit. (A partially completed project turned in on time will get some credit but a complete project turned in late will get no credit)

1

3

Project Deliverables

Each Group should turn in one group project consisting of:

- About 6-12 pages of text and figures
- 6-12 figures (please put multiple figures on one page where possible)
- Task list giving which tasks were performed by which group members. This should be signed by all members of the group.

4

Project Deliverables

Each *individual* should turn in the following:

- An *evaluation* of the contribution to the group project of every member of your group on a scale of 1(poor) to 10(excellent). To do this, write down the name of each member of your group (including yourself), and put a number between 1 and 10 next to each name. Please be honest and professional in your evaluation. Your evaluation will be one factor used to determine the project grade for each member of your group.

5

Project Comments

Before you turn in the project:

- Reread “The Top n Project Mistakes” section in the project description section on the course web page
- *You will loose points if you make these same mistakes*

6

Dynamic Programming

Dynamic Programming is different than Divide and Conquer in the following way:

- “Divide and Conquer” divides problem into independent subproblems, solves the subproblems recursively and then combines solutions to solve original problem
- Dynamic Programming is used when the subproblems are not independent, i.e. the subproblems share subsubproblems
- For these kinds of problems, divide and conquer does more work than necessary
- Dynamic Programming solves each subproblem once only and saves the answer in a table for future reference

7

The Pattern

- **Formulate the problem recursively.** Write down a formula for the whole problem as a simple combination of answers to smaller subproblems
- **Build solutions to your recurrence from the bottom up.** Write an algorithm that starts with the base cases of your recurrence and works its way up to the final solution by considering the intermediate subproblems in the correct order.

Note: Dynamic Programs store the results of intermediate subproblems. This is frequently *but not always* done with some type of table.

8

Edit Distance

- The *edit distance* between two words is the minimum number of letter insertions, letter deletions, and letter substitutions required to transform one word into another. For example, the edit distance between FOOD and MONEY is at most four:

FOOD → MOQD → MONYD → MONED → MONEY

9

String Alignment

Better way to display this process:

- Place the words one above the other in a table
- Put a gap in the first word for every insertion and a gap in the second word for every deletion
- Columns with two different characters correspond to substitutions
- Then the number of editing steps is just the number of columns that don't contain the same character twice

10

Example

- String Alignment for "FOOD" and "MONEY":

F	O	O		D
M	O	N	E	Y

- It's not too hard to see that we can't do better than four for the edit distance between "Food" and "Money"

11

Example II

- Unfortunately, it can be more difficult to compute the edit distance exactly. Example:

```
A L G O R   I   T H M
A L   T R U I S T I C
```

12

Key Observation

- If we remove the last column in an optimal alignment, the remaining alignment must also be optimal
- Easy to prove by contradiction: Assume there is some better subalignment of all but the last column. Then we can just paste the last column onto this better subalignment to get a better overall alignment.
- Note: The last column can be either: 1) a blank on top aligned with a character on bottom, 2) a character on top aligned with a blank on bottom or 3) a character on top aligned with a character on bottom

13

DP Solution

- To develop a DP algorithm for this problem, we first need to find a recursive definition
- Assume we have a m length string A and an n length string B
- Let $E(i, j)$ be the edit distance between the first i characters of A and the first j characters of B
- Then what we want to find is $E(n, m)$

14

Recursive Definition

- Say we want to compute $E(i, j)$ for some i and j
- Further say that the "Recursion Fairy" can tell us the solution to $E(i', j')$, for all $i' \leq i$, $j' \leq j$, except for $i' = i$ and $j' = j$
- Q: Can we compute $E(i, j)$ efficiently with help from our fairy friend?

15

Recursive Definition

There are three possible cases:

- **Insertion:** $E(i, j) = 1 + E(i - 1, j)$
- **Deletion:** $E(i, j) = 1 + E(i, j - 1)$
- **Substitution:** If $a_i = b_j$, $E(i, j) = E(i - 1, j - 1)$, else $E(i, j) = E(i - 1, j - 1) + 1$

16

Summary

Let $I(A[i] \neq B[j]) = 1$ if $A[i]$ and $B[j]$ are different, and 0 if they are the same. Then:

$$E(i, j) = \min \left\{ \begin{array}{l} E(i - 1, j) + 1, \\ E(i, j - 1) + 1, \\ E(i - 1, j - 1) + I(A[i] \neq B[j]) \end{array} \right\}$$

17

Base Case(s)

It's not too hard to see that:

- $E(0, j) = j$ for all j , since the j characters of B must be aligned with blanks
- Similarly, $E(i, 0) = i$ for all i

18

Recursive Alg

- We now have enough info to directly create a recursive algorithm
- The run time of this recursive algorithm would be given by the following recurrence:

$$T(m, 0) = T(0, n) = O(1)$$

$$T(m, n) = T(m, n - 1) + T(m - 1, n) + T(n - 1, m - 1) + O(1)$$

- Solution: $T(n, n) = \Theta(1 + \sqrt{2}^n)$, which is terribly, terribly slow.

19

Better Idea

- We can build up a $m \times n$ table which contains all values of $E(i, j)$
- We start by filling in the base cases for this table: the entries in the 0-th row and 0-th column
- To fill in any other entry, we need to know the values directly above, to the left and above and to the left.
- Thus we can fill in the table in the standard way: left to right and top down to ensure that the entries we need to fill in each cell are always available

		A	L	G	O	R	I	T	H	M
	0	→ 1	→ 2	→ 3	→ 4	→ 5	→ 6	→ 7	→ 8	→ 9
A	1	0	→ 1	→ 2	→ 3	→ 4	→ 5	→ 6	→ 7	→ 8
L	2	↓ 1	0	→ 1	→ 2	→ 3	→ 4	→ 5	→ 6	→ 7
T	3	↓ 2	↓ 1	↓ 1	→ 2	→ 3	→ 4	4	→ 5	→ 6
R	4	↓ 3	↓ 2	↓ 2	↓ 2	2	→ 3	→ 4	→ 5	→ 6
U	5	↓ 4	↓ 3	↓ 3	↓ 3	↓ 3	↓ 3	→ 4	→ 5	→ 6
I	6	↓ 5	↓ 4	↓ 4	↓ 4	↓ 4	3	→ 4	→ 5	→ 6
S	7	↓ 6	↓ 5	↓ 5	↓ 5	↓ 5	↓ 4	4	→ 5	→ 6
T	8	↓ 7	↓ 6	↓ 6	↓ 6	↓ 6	↓ 5	4	→ 5	→ 6
I	9	↓ 8	↓ 7	↓ 7	↓ 7	↓ 7	↓ 6	↓ 5	5	→ 6
C	10	↓ 9	↓ 8	↓ 8	↓ 8	↓ 8	↓ 7	↓ 6	↓ 6	6

20

Example Table

- Bold numbers indicate places where characters in the strings are equal
- Arrows represent predecessors that define each entry: horizontal arrow is deletion, vertical is insertion and diagonal is substitution.
- Bold diagonal arrows are “free” substitutions of a letter for itself
- Any path of arrows from the top left to the bottom right corner gives an optimal alignment (there are three paths in this example table, so there are three optimal edit sequences).

21

Analysis

- Let n be the length of the first string and m the length of the second string
- Then there are $\Theta(nm)$ entries in the table, and it takes $\Theta(1)$ time to fill each entry
- This implies that the run time of the algorithm is $\Theta(nm)$
- Q: Can you find a faster algorithm?

22

The code

```
EditDistance(A[1,..,m],B[1,..,n]){
  for (i=1;i<=m;i++){
    Edit[i,0] = i;}
  for (j=1;j<=n;j++){
    Edit[0,j] = j;}
  for (i=1;i<=m;i++){
    for (j=1;j<=n;j++){
      if (A[i]==B[j]){
        Edit[i,j] = min(Edit[i-1,j]+1,
                        Edit[i,j-1]+1,
                        Edit[i-1,j-1]);
      }else{
        Edit[i,j] = min(Edit[i-1,j]+1,
                        Edit[i,j-1]+1,
                        Edit[i-1,j-1]+1);
      }
    }
  }
  return Edit[m,n];}
```

23

Reconstructing an optimal alignment

- In this code, we do not keep info around to reconstruct the optimal alignment
- However, it is a simple matter to keep around another array which stores, for each cell, a pointer to the cell that was used to achieve the current cell's minimum edit distance
- To reconstruct a solution, we then need only follow these pointers from the bottom right corner up to the top left corner

24

In Class Exercise

- Create a string alignment table for the two strings "abba" and "bab". Put "abba" at the top of the table and "bab" on the left side
- Q1: ($i = 1, 2, \dots, 5$) What is the i -th row of your table
- Q6: What is the minimum edit distance and how many alignments achieve it?

25

Take Away

- To solve the string alignment problem, we did the following:
1) formulated the problem recursively 2) built a solution to the recurrence from the bottom up

26