

Edit Distance Problem

- Given two strings x and y what is the cheapest way to convert x into y
- By conversion we mean character edits:
 - Insert character c
 - Delete character c
 - Replace character c with c'
- Each one of these operations has a certain cost that is also a function of the character being inserted, deleted or replaced.
- The cost information is usually previously provided in a table
- Used on auto-correct string features and comparing DNA sequences
spelling corrections

- Another problem encompassed by Edit Distance is the longest common subsequence problem, e.g.:

X string: H I E R O G L Y P H O L O G Y
 Y string: M I C H A E L A N G E L O

Question:
 What is the longest common subsequence?

Not subtracting! ~~A~~

- We are allowed to drop any set of letters from X.
- " " " " " " " " " " " "
- We are not allowed to replace any characters
- I want X and Y in the end to be equal

- Longest common subsequence in the example:

H	I	E	R	O	G	L	Y	P	H	O	L	O	G	Y	}	H	E	L	L	O
M	I	C	H	A	E	L	A	N	G	E	L	O								
- We can model the longest common subsequence problem as an edit distance problem in the following way:
 - cost of insert/delete = 1
 - cost of replace = $\begin{cases} 0, & \text{if } c = c' \\ \infty, & \text{otherwise} \end{cases}$
- We want to minimize the number of insertions and deletions
- Lets focus on the more general edit distance problem
- Ideally we should be able to use our "subproblems for strings" strategies, i.e. { suffixes, prefixes, substrings }
- But now we have two strings, which is annoying
- Lets see how we can solve this:
 - ① subproblem = solve edit distance problem on $x[i:]$ and $y[j:] \quad \forall i, j$
 - # subproblems = n choices for x suffixes and n choices for y suffixes
 - = $n \times n = n^2 \Rightarrow \Theta(n^2)$

If x and y have different lengths then:
 $\Theta(|x| \cdot |y|)$

② Guess / Try :

string u: _____
 i

string y: _____
 j

Question:

I want to convert u into y. What should I look into here?

• Lets focus on the first character.

Question:

What are the possible ways to deal with the first character of u?

Objective:

- I want the first character of u to become the first character of y

↳ Possibility 1 I can delete the first character of u

↳ Possibility 2 I can replace the first character of u with the first character of y

↳ Possibility 3 I can insert the first character of y immediately before the first character of u

∴ We need to guess/try one of 3 possibilities

- replace $u[i] \rightarrow y[j]$
- insert $y[j]$
- delete $x[i]$

③ Recurrence

$$DP(i, j) = \min ($$

↙ ↘
suffix of x / suffix of y

Option 1: "cost of replace $x[i] \rightarrow y[j]$ " + $DP(i+1, j)$
Option 2: "cost of insert $y[j]$ " + $DP(i, j+1)$,
Option 3: "cost of delete $x[i]$ " + $DP(i+1, j)$

$$)$$

④ Recurse + memoize - nothing to be said here

⑤ $DP(0, 0)$

Question:

What is the running time?

#subproblems : $\Theta(|X| \times |Y|)$

time/subproblem : Each cost can be looked up on a table $\Rightarrow \Theta(1)$ time for the lookup.

$$\begin{aligned}
 \text{Total Time} &= \# \text{ subproblems} \times \text{time/subproblem} \\
 &= \Theta(|X| \times |Y|) \times \Theta(1) \\
 &= \Theta(|X| \times |Y|) \approx \Theta(n^2) \text{ if we assume } |X| = |Y| = n
 \end{aligned}$$

Knapsack Problem:

- List of items each with a size s_i and a value v_i
- The sizes are integers
- Knapsack of size S
- Objective: We want to maximize sum of values for a subset of items of total size $\leq S$

Question: How do we do this with DP?

Answer: With difficulty...

① S-b problems:

- We can look at this problem as a sequence of items, even though the order is not important
- This allows us to use the "subproblems for storage" strategies
 - suffix
 - prefix
 - subsetting
- Lets see with suffix i : of items. This is helpful because we can choose i items from the beginning
- What should I decide with the i th item
 - Should item i be included or not?
- We also need to tell the recursion how much available space we will have in the knapsack if we decide to take the i

- ∴ - suffix of items i :
- remaining capacity $x \leq S$

subproblems = $n \cdot S \Rightarrow \Theta(n \cdot S)$

- ② Guess / Try: is item i in subset or not?
- ③ Recurrence: # Guesses = 2 choices (yes/no) $\Rightarrow \Theta(1)$ time

$DP(i, x) = \max($

Option 1: "Don't include the item"

$DP(i+1, x),$

Option 2: "Include item i "

$DP(i+1, x - s_i) + v_i$

)

Total Time = # subproblems \times time/subproblem
 $= \Theta(n \cdot S) \times \Theta(1) = \Theta(n \cdot S)$

Question:
 Is this polynomial time?

↳ Not polynomial time!

↳ Polynomial time \Rightarrow polynomial of ~~the~~ the size of the input
 size of the input here is n and S

Question: Why? Any ideas?

- We would need $\log_2 S$ bits to encode all possible input sizes

- S is exponential in $\log S$ \Rightarrow exponential time algorithm !!
 1 bit $\Rightarrow 2^1$ combinations
 2 bits $\Rightarrow 2^2$ " " }

- The algorithm seems polynomial but it is not.
- Well if \underline{S} is small then the algorithm can be seen as polynomial
- Otherwise if \underline{S} is large then we have exponential time
- We call this ^{behaviour} pseudo polynomial time
- In an abuse of notation:

polynomial time < pseudo polynomial < exponential time
time

54 Longest Common Subsequence Problem [Cormen 2001]

Question: How "similar" two ~~strings~~ ^{sequences} are?

- Different measures of similarity:
 - Substring
 - Number of changes to turn one string into another
 - Subsequences
- Subsequence of a sequence is just the given sequence with zero or more elements left out. Formally:

Definition: Given a sequence $X = \langle x_1, x_2, \dots, x_m \rangle$ another sequence $Z = \langle z_1, z_2, \dots, z_k \rangle$ is a subsequence of X if there exists a strictly increasing sequence $\langle i_1, i_2, \dots, i_k \rangle$ of indices of X such that for all $j = 1, 2, \dots, k$ we have $x_{i_j} = z_j$.
- E.g. $Z = \langle B, C, D, B \rangle$ is a subsequence of $X = \langle A, B, C, \dots, B, D, A, B \rangle$ with corresponding index sequence $\langle 2, 3, 5, 7 \rangle$
- Given two sequences \underline{X} and \underline{Y} we say that a sequence \underline{Z} is a common subsequence of \underline{X} and \underline{Y} if \underline{Z} is a subsequence of both \underline{X} and \underline{Y} .

- E.g.:

- $X = \langle A, B, C, B, D, A, B \rangle$
 - $Y = \langle B, D, C, A, B, A \rangle$
- } A common subsequence: $\langle B, C, A \rangle$
 } Longest common subsequences:
- $\langle B, C, B, A \rangle$
 - $\langle B, D, A, B \rangle$

Theorem 15.1 - Optimal substructure of an LCS

- Let $X = \langle u_1, u_2, \dots, u_m \rangle$
 $Y = \langle y_1, y_2, \dots, y_n \rangle$

be sequences

- Let $Z = \langle z_1, z_2, \dots, z_k \rangle$ be an LCS of X and Y

- Then:

1. If $u_m = y_n$, then $z_k = u_m = y_n$ and z_{k-1} is an LCS of X_{m-1} and Y_{n-1}
2. If $u_m \neq y_n$, then $z_k \neq u_m$ implies that Z is an LCS of X_{m-1} and Y
3. If $u_m \neq y_n$, then $z_k \neq y_n$ implies that Z is an LCS of X and Y_{n-1}

Prefix

Notes:

- Given a sequence $X = \langle u_1, u_2, \dots, u_m \rangle$ we define the c th prefix of X, for $c = 0, 1, \dots, m$ as $X_c = \langle u_1, u_2, \dots, u_c \rangle$
- E.g. $X = \langle A, B, C, B, D, A, B \rangle$ then $X_4 = \langle A, B, C, B \rangle$ and X_0 is the empty subsequence

Proof:

- (1). If $z_k \neq u_m$ ~~then we could append $u_m = y_n$~~ then we could append $u_m = y_n$ to Z to obtain a common subsequence of X and Y of length $k+1$
- This would be a contradiction that Z is an LCS. Thus we must have that $z_k = u_m = y_n$
- Now, the prefix z_{k-1} is a length $(k-1)$ common subsequence of X_{m-1} and Y_{n-1} . We wish to show that it is an LCS.
- Suppose for the purpose of contradiction that there is a common subsequence W of X_{m-1} and Y_{n-1} with length greater than $k-1$
- Then appending $u_m = y_n$ to W produces a common subsequence of X and Y whose length is greater than k , which is a contradiction

$Y_{\{n-1\}}$

(2) If $z_k \neq u_m$ then z is a common subsequence of X_{m-1} and Y .

If there were a common subsequence w of X_{m-1} and Y with length greater than k then w would also be a common subsequence of X_m and Y , contradicting the assumption that z is an LCS of X and Y .

(3) The proof is symmetric to (2)

■ (end of proof)

15.4.2 - Recursive Solution

- Theorem 15.1 implies that there are either one or two subproblems to examine

- If $u_m = y_n$ we must find LCS of X_{m-1} and Y_{n-1} . Appending $u_m = y_n$ to this LCS yields an LCS of X and Y

- If $u_m \neq y_n$ then we must solve two subproblems:

Subproblem 1: Finding an LCS of X_{m-1} and Y

Subproblem 2: Finding an LCS of X and Y_{n-1}

Whichever of these ~~#~~ two LCS^s is longer is an LCS of X and Y .

- Define $C[i, j]$ to be the length of an LCS of sequences X_i and Y_j

If either $i = 0$ or $j = 0$ one of the sequences has length 0 so the LCS has length 0:

$$C[i, j] = \begin{cases} 0 & , \text{ if } i = 0 \text{ or } j = 0 \\ C[i-1, j-1] & , \text{ if } i, j > 0 \text{ and } u_i = y_j \\ \max(C[i, j-1], C[i-1, j]) & , \text{ if } i, j > 0 \text{ and } u_i \neq y_j \end{cases}$$

- # subproblems:

- We need to check each prefix of X - there are m

- We need to check each prefix of Y - there are n

- Total number of subproblems $m \cdot n$

- Total time = # subproblems \times time/subproblem

$= O(m \cdot n) \cdot O(1) = O(m \cdot n)$

15.4.3 Computing the length of an LCS

- For once lets see how we can calculate an iterative procedure

- LCS-Length (X, Y)

m = X.length

n = Y.length

let $\begin{cases} b[1..m, 1..n] \\ c[0..m, 0..n] \end{cases}$ be new tables

** Table b points to the optimal subproblem solution **

for i = 1 to m

 c[i, 0] = 0

 for j = 0 to n

 c[0, j] = 0

for i = 1 to m

 for j = 1 to n

 if $x_i == y_j$

 c[i, j] = c[i-1, j-1] + 1

 b[i, j] = "R"

 elseif c[i-1, j] > c[i, j-1]

 c[i, j] = c[i-1, j]

 b[i, j] = "A"

 else c[i, j] = c[i, j-1]

 b[i, j] = "L"

// c[m, n] contains the length of an LCS of X and Y

	J	0	1	2	3	4	5	6
i	y_j	B	D	C	A	B	A	
0	w_i	0	0	0	0	0	0	0
1	A	0	0 \uparrow	0 \uparrow	0 \uparrow	1 \leftarrow	1 \leftarrow	1 \leftarrow
2	B	0	1 \leftarrow	1 \leftarrow	1 \leftarrow	1 \uparrow	2 \leftarrow	2 \leftarrow
3	C	0	1 \uparrow	1 \uparrow	2 \leftarrow	2 \leftarrow	2 \uparrow	2 \uparrow
4	B	0	1 \leftarrow	1 \uparrow	2 \uparrow	2 \uparrow	3 \leftarrow	3 \leftarrow
5	D	0	1 \uparrow	2 \leftarrow	2 \uparrow	2 \uparrow	3 \uparrow	3 \uparrow
6	A	0	1 \uparrow	2 \uparrow	2 \uparrow	3 \leftarrow	3 \uparrow	4 \leftarrow
7	B	0	1 \leftarrow	2 \uparrow	2 \uparrow	3 \uparrow	4 \leftarrow	4 \uparrow

Example for strings:

$X = \langle A, B, C, B, D, A, B \rangle$

$Y = \langle B, D, C, A, B, A \rangle$

→ To reconstruct the elements of an LCS follow the $b[i, j]$ from the lower right corner. Each \leftarrow on the shaded sequence corresponds to an entry for which $w_i = y_j$ is a member of an LCS

An LCS = $\langle B, C, B, A \rangle$