

## String edit distance

Data Structures and Algorithms for Computational Linguistics III  
(ISCL-BA-07)

Çağrı Çöltekin  
ccoltekin@ufa.umi-tuebingen.de

University of Tübingen  
Seminar für Sprachwissenschaft

Winter Semester 2020/21

www.747320.0001.01.04

## Edit distance

- In many applications, we want to know how similar (or different) two strings are
  - Comparing two files (e.g., source code)
  - Comparing two DNA sequences
  - Spell checking
  - Approximate string matching
  - Determining similarity of two languages
  - Machine translation
- The solution is typically formulated as the (inverse) cost of obtaining one of the strings from the other through a number of *edit operations*
- Once we obtain the optimal edit operations, we may (depending on the edit operations) also be able to determine the optimal alignment between the strings

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 1 / 19

## Hamming distance

a simple distance metric between two sequences

- The Hamming distance measures number of different symbols in the corresponding positions

h	g	l	e	n	e	
h	g	l	e	n	e	
h	l	g	l	e	n	e
h	l	y	g	e	l	n

$0+1+0+0+0+0=1$        $0+1+1+1+0+1+1=5$

- Very easy/efficient calculation
- But cannot handle sequences of different lengths (consider *hygiene* – *hygiene*)

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 2 / 19

## A family of edit distance problems

- The same overall idea applies to a number of well-known problems/solutions that differ in the type of operations allowed
  - Hamming distance: only replacements
  - Longest common subsequence: (LCS) insertions and deletions
  - Levenshtein distance: insertions, deletions and substitutions
  - Levenshtein-Damerau distance: insertions, deletions and substitutions and transpositions (swap) of adjacent symbols
- Naive solutions to all (except Hamming distance) have exponential time complexity
- Polynomial-time solution can be obtained using *dynamic programming*

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 3 / 19

## Longest common subsequence (LCS)

Problem definition

- A subsequence is an order-preserving (but not necessarily continuous) sequence of symbols from a string (a version of the sequence where zero or more elements are removed)
  - hyg, gn, yrte, hnt, gme are subsequences of hygiene
- Note that a subsequence does not have to be a substring (substrings are continuous)
  - hyg, yrte, gme are substrings of hygiene
- The longest common subsequence (LCS) of two strings is the longest string that is a subsequence of both strings
  - LCS(hygiene, hygiene) = hygiene
  - LCS(hygiene, hygiene) = hygiene / hygiene
- LCS is exactly the problem solved by the UNIX *diff* utility
- It has wide-ranging applications from source-code comparison to bioinformatics (e.g., DNA sequencing)

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 4 / 19

## LCS: a naive solution

- A simple solution is:
  - Enumerate all subsequences of the first string
  - Check if it is also a subsequence of the second string
- There are exponential number of subsequences of a string
  - the string *abc* has 8 subsequences:
    - abc: nothing removed
    - ab, ac, bc: individual elements are removed
    - a, b, c: length-2 subsequences are removed
    - ε (empty string): all removed
  - For *abcd*, if we have subsequences of *abc* once with, and once without *d*
    - Each additional symbol doubles the number of subsequences
- For strings of size *n* and *m*, the complexity of the brute-force algorithm is  $O(2^n \cdot m)$

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 5 / 19

## LCS: recursive solution

demonstration

- Consider two strings  $Xx, Yy$  and their LCS  $Zz$  ( $X, Y, Z$  are possibly empty strings,  $x, y, z$  are characters)
- If  $x = y$ , then this character has to be part of the LCS,  $x = y = z$ , and  $Z$  must be the LCS of  $X$  and  $Y$
- If  $x \neq y$ , there are three cases
  - $x \neq y \neq z$ :  $Zz$  is also the LCS of  $X$  and  $Y$
  - $x = z \neq y$ :  $Zz$  is also the LCS of  $Xx$  and  $Y$
  - $y = z \neq x$ :  $Zz$  is also the LCS of  $X$  and  $Yy$
- This leads to following recursive definition:

$$\text{LCS}(Xx, Yy) = \begin{cases} \text{LCS}(X, Y) & \text{if } x = y \\ \max(\text{LCS}(Xx, Y), \text{LCS}(X, Yy)) & \text{otherwise} \end{cases}$$

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 6 / 19

## LCS: divide-and-conquer



- Note the **repeated computation**

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 7 / 19

## LCS: dynamic programming

general sketch

- For string indexes  $i$  and  $j$ , of strings  $X$  and  $Y$ , if we need  $\text{LCS}(X_{1..i}, Y_{1..j})$ ,  $\text{LCS}(X_{1..i-1}, Y_j)$ ,  $\text{LCS}(X_i, Y_{1..j-1})$
- In the standard algorithm, we do not store the LCS, but the length of the LCS,  $l_{i,j}$  for each  $i, j$
- Once we fill in the matrix, the  $l_{n,m}$  is the length of the LCS
- We can trace back and recover the LCS using the dynamic programming matrix

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 8 / 19

## LCS with dynamic programming

demonstration

		0	1	2	3	4	5	6	7	8
	c									
0	e	0	0	0	0	0	0	0	0	0
1	h	0	1	1	1	1	1	1	1	1
2	y	0	1	1	2	2	2	2	2	2
3	g	0	1	1	2	3	3	3	3	3
4	l	0	1	2	2	3	3	4	4	4
5	e	0	1	2	2	3	4	4	4	5
6	n	0	1	2	2	3	4	4	5	5
7	e	0	1	2	2	3	4	4	5	6

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 9 / 19

## Complexity of filling the LCS matrix

```
l = np.zeros(shape=(n+1, m+1))
for i in range(1, n+1):
    for j in range(1, m+1):
        if X[i] == Y[j]:
            l[i+1, j+1] = l[i, j] + 1
        else:
            l[i+1, j+1] = max(l[i+1, j], l[i, j+1])
```

- Two loops up to  $n$  and  $m$ , the time complexity is  $O(nm)$
- Similarly, the space complexity is also  $O(nm)$

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 10 / 19

## Recovering the LCS from the matrix

C. Çöltekin, IM | University of Tübingen

Winter Semester 2020/21 11 / 19

## Transforming one string to another

- The table (back arrows) also gives a set of edit operations to transform one string to another
- For LCS, operations are:
  - copy (diagonal arrows in the demonstration)
  - insert (left arrows in the demonstration – assuming original string is the vertical one)
  - delete (up arrows in the demonstration – assuming original string is the vertical one)
- These also form an alignment between two strings
- Different set of edit operations recovered will yield the same LCS, but different alignments

## LCS alignments

	0	1	2	3	4	5	6	7	8
0	ε	h	i	y	g	e	i	n	e
1	ε	0	1	1	1	1	1	1	1
2	h	0	1	1	2	2	2	2	2
3	y	0	1	1	2	3	3	3	3
4	g	0	1	2	2	3	3	4	4
5	i	0	1	2	2	3	4	4	5
6	e	0	1	2	2	3	4	4	5
7	n	0	1	2	2	3	4	4	5
8	e	0	1	2	2	3	4	4	5

Alignments:

h-y-i-ne  
 cicciccc  
 hiygi-ne  
 h-y-e-ne  
 cicciccc  
 hiy-wine

## LCS – some remarks

- We formulated the algorithm as optimizing the LCS
- Alternatively, we can consider costs associated with each operation:
  - copy = 0
  - delete = 1
  - insert = 1
- This is the typical application of LCS, as in *diff*
- In some applications we may want to have different costs for delete and insert (e.g., mapping lemmas to inflected forms of words)
- Similarly, we may want to assign different costs for different characters (e.g., higher cost to delete consonants in historical linguistics)

## Levenshtein distance

definition

- Levenshtein difference between two strings is the total cost of *insertions*, *deletions* and *substitutions*
- With cost of 1 for all operations

$$\text{lev}(Xx, Yy) = \begin{cases} \text{len}(X) & \text{if } \text{len}(Y) = 0 \\ \text{len}(Y) & \text{if } \text{len}(X) = 0 \\ \text{lev}(X, Y) & \text{if } x = y \\ 1 + \min \begin{cases} \text{lev}(X, Yy) \\ \text{lev}(Xx, Y) \end{cases} & \text{otherwise} \end{cases}$$

- Naive recursion (as defined above), again, is intractable
- But, the same dynamic programming method works

## Levenshtein distance

demonstration

	0	1	2	3	4	5	6	7	8
0	ε	h	i	y	g	e	i	n	e
1	h	0	1	1	2	3	4	5	6
2	y	2	1	1	1	2	3	4	5
3	g	3	2	2	2	1	2	3	4
4	i	4	3	2	3	2	2	3	4
5	e	5	4	3	3	3	2	3	3
6	n	6	5	4	4	4	3	3	4
7	e	7	6	5	5	5	4	4	4

## Levenshtein distance

edits and alignments

	0	1	2	3	4	5	6	7	8
0	ε	h	i	y	g	e	i	n	e
1	h	0	1	1	2	3	4	5	6
2	y	2	1	1	1	2	3	4	5
3	g	3	2	2	2	1	2	3	4
4	i	4	3	2	3	2	2	3	4
5	e	5	4	3	3	3	2	3	3
6	n	6	5	4	4	4	3	3	4
7	e	7	6	5	5	5	4	4	4

## Edit distance: extensions and variations

- Another possible operation we did not cover is *swap* (or transpose), which is useful for applications like spell checking
- In some applications (e.g., machine translation, OCR correction) we may want to have one-to-many or many-to-one alignments
- Additional requirements often introduce additional complexity
- It is sometimes useful to learn costs from data

## Summary

- Edit distance is an important problem in many fields including computational linguistics
- A number of related problems can be efficiently solved by dynamic programming
- Edit distance is also important for approximate string matching and alignment
- Reading suggestion: Goodrich, Tamassia, and Goldwasser (2013, chapter 13), Jurafsky and Martin (2009, section 3.11, or 2.5 in online draft)

Next:

- Algorithms on strings: tries
- Reading: Goodrich, Tamassia, and Goldwasser (2013, chapter 13),

## Acknowledgments, credits, references

- Goodrich, Michael T., Roberto Tamassia, and Michael H. Goldwasser (2013). *Data Structures and Algorithms in Python*. John Wiley & Sons, Incorporated. <https://doi.org/10.11197/9781118476734>.
- Jurafsky, Daniel and James H. Martin (2009). *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. second edition. Pearson Prentice Hall. <https://doi.org/10.1139/9780131304196-3>.

