

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 demo – assuming original string is the vertical one)
 - delete (up arrows in the demo)
- 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	4
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-g-i-n-e
 cicciccc
 h-y-g-i-n-e
 cicciccc
 h-y-g-e-n-e
 cicciccc
 h-y-g-w-n-e

LCS – some remarks

- We formulated the algorithm as maximizing the LCS
- Alternatively, we can minimize the costs associated with each operation:
 - copy = 0
 - delete = 1
 - insert = 1
- The cost settings above are the typical, e.g., 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

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

- 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	1	0	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	1	0	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 *saup* (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: [goodrich2013](#), [jurafsky2009](#)
- Next:
- Algorithms on strings: tries
 - Reading: [goodrich2013](#),

Acknowledgments, credits, references

