

# Addendum to Lecture 3: Edit Distance (with animation)

Information Retrieval  
Computer Science Tripos Part II

Simone Teufel

Natural Language and Information Processing (NLIP) Group



**UNIVERSITY OF  
CAMBRIDGE**

`Simone.Teufel@cl.cam.ac.uk`

Lent 2014

- **Edit distance** between two strings  $s_1$  and  $s_2$  is the minimum number of basic operations that transform  $s_1$  into  $s_2$ .
- **Levenshtein distance:** Admissible operations are [insert](#), [delete](#) and [replace](#)

## Levenshtein distance

dog	–	do	1 (delete)
cat	–	cart	1 (insert)
cat	–	cut	1 (replace)
cat	–	act	2 (delete+insert)

# Levenshtein distance: Distance matrix

		s	n	o	w
	0	1	2	3	4
o	1	1	2	2	3
s	2	1	2	3	3
l	3	2	2	3	4
o	4	3	3	2	3

# Edit Distance: Four cells

		s	n	o	w
	<u>  </u> 0	<u>  1  </u> 1  1	<u>  2  </u> 2  2	<u>  3  </u> 3  3	<u>  4  </u> 4  4
o	<u>  1  </u> 1	<u>  1  2  </u> 2  1	<u>  2  3  </u> 2  2	<u>  2  4  </u> 3  2	<u>  4  5  </u> 3  3
s	<u>  2  </u> 2	<u>  1  2  </u> 3  1	<u>  2  3  </u> 2  2	<u>  3  3  </u> 3  3	<u>  3  4  </u> 4  3
l	<u>  3  </u> 3	<u>  3  2  </u> 4  2	<u>  2  3  </u> 3  2	<u>  3  4  </u> 3  3	<u>  4  4  </u> 4  4
o	<u>  4  </u> 4	<u>  4  3  </u> 5  3	<u>  3  3  </u> 4  3	<u>  2  4  </u> 4  2	<u>  4  5  </u> 3  3

Cormen et al:

- **Optimal substructure:** The optimal solution contains within it subsolutions, i.e, optimal solutions to subproblems
- **Overlapping subsolutions:** The subsolutions overlap and would be computed over and over again by a brute-force algorithm.

For edit distance:

- **Subproblem:** edit distance of two prefixes
- **Overlap:** most distances of prefixes are needed 3 times (when moving right, diagonally, down in the matrix)

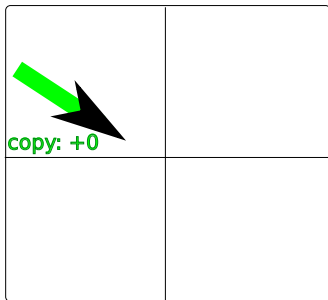
# Levenshtein distance: Algorithm

LEVENSHTEINDISTANCE( $s_1, s_2$ )

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

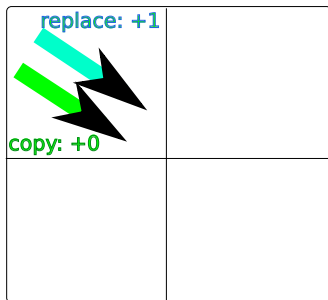
# Each cell of Levenshtein matrix


# Each cell of Levenshtein matrix

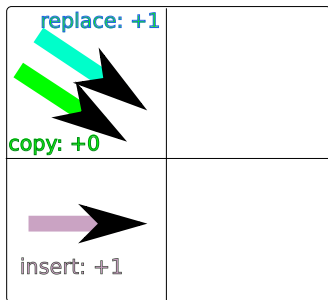




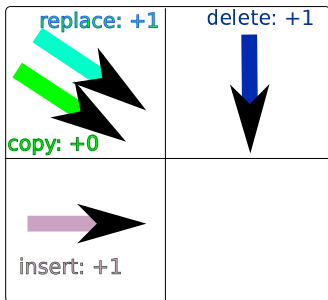
# Each cell of Levenshtein matrix



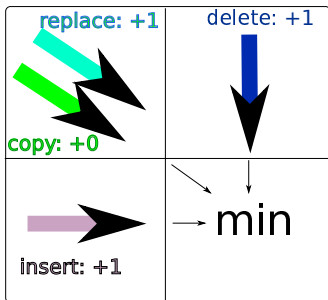
# Each cell of Levenshtein matrix



# Each cell of Levenshtein matrix



# Each cell of Levenshtein matrix



# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	$\frac{\quad}{\quad}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$				
s	$\frac{2}{2}$				
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2			
s	2 2				
l	3 3				
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	<u>  </u> 0	<u>  </u> 1 1	<u>  </u> 2 2	<u>  </u> 3 3	<u>  </u> 4 4
o	<u>  </u> 1 1	<u>  </u> 1 2 2 1			
s	<u>  </u> 2 2				
l	<u>  </u> 3 3				
o	<u>  </u> 4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2		
s	2 2				
l	3 3				
o	4 4				



# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	<u>  </u> 0	<u>  </u> 1 1	<u>  </u> 2 2	<u>  </u> 3 3	<u>  </u> 4 4
o	<u>  </u> 1 1	<u>  </u> 1 2 2 1	<u>  </u> 2 3 2 2		
s	<u>  </u> 2 2				
l	<u>  </u> 3 3				
o	<u>  </u> 4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	<u>  </u> 0	<u>  </u> 1 1	<u>  </u> 2 2	<u>  </u> 3 3	<u>  </u> 4 4
o	<u>  </u> 1 <u>  </u> 1	<u>  </u> 1 2 <u>  </u> 2 1	<u>  </u> 2 3 <u>  </u> 2 2	<u>  </u> 2 4 <u>  </u> 3	
s	<u>  </u> 2 <u>  </u> 2				
l	<u>  </u> 3 <u>  </u> 3				
o	<u>  </u> 4 <u>  </u> 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	
s	2 2				
l	3 3				
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3
s	2 2				
l	3 3				
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2				
l	3 3				
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3			
l	3 3				
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	<u>  </u> 0	<u>  </u> 1 1	<u>  </u> 2 2	<u>  </u> 3 3	<u>  </u> 4 4
o	<u>  </u> 1 <u>  </u> 1	<u>  </u> 1 2 <u>  </u> 2 1	<u>  </u> 2 3 <u>  </u> 2 2	<u>  </u> 2 4 <u>  </u> 3 2	<u>  </u> 4 5 <u>  </u> 3 3
s	<u>  </u> 2 <u>  </u> 2	<u>  </u> 1 2 <u>  </u> 3 1			
l	<u>  </u> 3 <u>  </u> 3				
o	<u>  </u> 4 <u>  </u> 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2		
l	3 3				
o	4 4				



# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	<u>  </u> 0	<u>  </u> 1 1	<u>  </u> 2 2	<u>  </u> 3 3	<u>  </u> 4 4
o	<u>  </u> 1 <u>  </u> 1	<u>  </u> 1 2 <u>  </u> 2 1	<u>  </u> 2 3 <u>  </u> 2 2	<u>  </u> 2 4 <u>  </u> 3 2	<u>  </u> 4 5 <u>  </u> 3 3
s	<u>  </u> 2 <u>  </u> 2	<u>  </u> 1 2 <u>  </u> 3 1	<u>  </u> 2 3 <u>  </u> 2 2		
l	<u>  </u> 3 <u>  </u> 3				
o	<u>  </u> 4 <u>  </u> 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3	
l	3 3				
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	
l	3 3				
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	<u>  </u> 0	<u>  </u> 1 1	<u>  </u> 2 2	<u>  </u> 3 3	<u>  </u> 4 4
o	<u>  </u> 1 <u>  </u> 1	<u>  </u> 1 2 <u>  </u> 2 1	<u>  </u> 2 3 <u>  </u> 2 2	<u>  </u> 2 4 <u>  </u> 3 2	<u>  </u> 4 5 <u>  </u> 3 3
s	<u>  </u> 2 <u>  </u> 2	<u>  </u> 1 2 <u>  </u> 3 1	<u>  </u> 2 3 <u>  </u> 2 2	<u>  </u> 3 3 <u>  </u> 3 3	<u>  </u> 3 4 <u>  </u> 4
l	<u>  </u> 3 <u>  </u> 3				
o	<u>  </u> 4 <u>  </u> 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3				
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4			
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2			
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3		
o	4 4				



# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2		
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3	
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4				

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5			

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3			

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4		



# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3		

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3	2 4 4	

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3	2 4 4 2	

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3	2 4 4 2	4 5 3

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3	2 4 4 2	4 5 3 3

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3	2 4 4 2	4 5 3 3

Edit distance OSLO–SNOW is 3!

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3	2 4 4 2	4 5 3 3

How do I read out the editing operations that transform OSLO into SNOW?

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3	2 4 4 2	4 5 3 3

cost | operation || input | output

1 | insert || \* | w



# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3	2 4 4 2	4 5 3 3

cost	operation	input	output
------	-----------	-------	--------

0	(copy)	o	o
1	insert	*	w

# Example: Edit Distance OSLO – SNOW

		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3	2 4 4 2	4 5 3 3

cost	operation	input	output
------	-----------	-------	--------

1	replace	l	n
0	(copy)	o	o
1	insert	*	w

# Example: Edit Distance OSLO – SNOW

			s	n	o	w
		0	1 1	2 2	3 3	4 4
o		1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s		2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l		3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o		4 4	4 3 5 3	3 3 4 3	2 4 4 2	4 5 3 3

cost	operation	input	output
0	(copy)	s	s
1	replace	l	n
0	(copy)	o	o
1	insert	*	w

# Example: Edit Distance OSLO – SNOW

			s	n	o	w
		0	1 1	2 2	3 3	4 4
o	1 1	1 2	2 3	3 4	4 5	
	1 1	2 1	2 2	3 2	3 3	
s	2 2	1 2	2 3	3 3	3 4	
	2 2	3 1	2 2	3 3	4 3	
l	3 3	3 2	2 3	3 4	4 4	
	3 3	4 2	3 2	3 3	4 4	
o	4 4	4 3	3 3	2 4	4 5	
	4 4	5 3	4 3	4 2	3 3	

cost	operation	input	output
1	delete	o	*
0	(copy)	s	s
1	replace	l	n
0	(copy)	o	o
1	insert	*	w