


12 Edit distance

Summer Term 2011

Jan-Georg Smaus

Albert-Ludwigs-Universität Freiburg



Problem: similarity of strings

Edit distance

For two strings A and B , compute, as efficiently as possible, the edit distance $D(A,B)$ and a minimal sequence of edit operations which transforms A into B .

i n f - - o r m a t i k -
i n t e r p o l - a t i o n

06.07.2011 Theory 1 - Edit distance 2

Application of edit distance

Approximate string matching

For a given text T , a pattern P , and a distance d , find all substrings P' in T with $D(P,P') \leq d$

Sequence alignment

Find optimal alignments of DNA sequences

G A G C A - C T T G G A T T C T C G G
- - - C A C G T G G - - - - - - - -

06.07.2011 Theory 1 - Edit distance 3

Edit distance

Given: two strings $A = a_1 a_2 \dots a_m$ and $B = b_1 b_2 \dots b_n$

Goal: find minimal cost $D(A,B)$ for a sequence of edit operations to transform A into B .

Edit operations:

1. Replace a character in A by a character from B
2. Delete a character from A
3. Insert a character into A

06.07.2011 Theory 1 - Edit distance 4

Edit distance

Cost model:

$$c(a,b) = \begin{cases} 1 & \text{if } a \neq b \\ 0 & \text{if } a = b \end{cases}$$

$a = \epsilon, b = \epsilon$ possible

We assume the **triangle inequality** holds for c :

$$c(a,c) \leq c(a,b) + c(b,c)$$

→ Each character is changed at most once

06.07.2011 Theory 1 - Edit distance 5

Edit distances

Trace as representation of edit sequences

A = b a a c a a b c
B = a b a c b c a c

or using **indents**

A = - b a a c a - a b c
B = a b a - c b c a - c

Edit distance (cost): 5

An optimal trace can be divided into two optimal substraces.

06.07.2011 Theory 1 - Edit distance 6

Computation of the edit distance

Let $A_i = a_1 \dots a_i$ and $B_j = b_1 \dots b_j$

$$D_{ij} = D(A_i, B_j)$$

06.07.2011 Theory 1 - Edit distance 7

Computations of the edit distances

- Three possibilities of ending a trace:
 - a_m is replaced by b_n :

$$D_{m,n} = D_{m-1,n-1} + c(a_m, b_n)$$
 ($c(a_m, b_n)$ can be 0 or 1)
 - a_m is deleted: $D_{m,n} = D_{m-1,n} + 1$
 - b_n is inserted: $D_{m,n} = D_{m,n-1} + 1$

06.07.2011 Theory 1 - Edit distance 8

Computations of the edit distances

- Three possibilities of ending a trace:
 - a_m is replaced by b_n :

$$D_{m,n} = D_{m-1,n-1} + c(a_m, b_n)$$
 ($c(a_m, b_n)$ can be 0 or 1)
 - a_m is deleted: $D_{m,n} = D_{m-1,n} + 1$
 - b_n is inserted: $D_{m,n} = D_{m,n-1} + 1$

06.07.2011 Theory 1 - Edit distance 9

Computations of the edit distances

- Three possibilities of ending a trace:
 - a_m is replaced by b_n :

$$D_{m,n} = D_{m-1,n-1} + c(a_m, b_n)$$
 ($c(a_m, b_n)$ can be 0 or 1)
 - a_m is deleted: $D_{m,n} = D_{m-1,n} + 1$
 - b_n is inserted: $D_{m,n} = D_{m,n-1} + 1$

06.07.2011 Theory 1 - Edit distance 10

Computation of the edit distance

- Recurrence relation, if $m, n \geq 1$:

$$D_{m,n} = \min \left\{ \begin{array}{l} D_{m-1,n-1} + c(a_m, b_n) \\ D_{m-1,n} + 1 \\ D_{m,n-1} + 1 \end{array} \right\}$$
- Computation of all D_{ij} is required, $0 \leq i \leq m, 0 \leq j \leq n$.

06.07.2011 Theory 1 - Edit distance 11

Recurrence relation for the edit distance

Base cases:

$$D_{0,0} = D(\epsilon, \epsilon) = 0$$

$$D_{0,j} = D(\epsilon, B_j) = j$$

$$D_{i,0} = D(A_i, \epsilon) = i$$

Recurrence equation:

$$D_{i,j} = \min \left\{ \begin{array}{l} D_{i-1,j-1} + c(a_i, b_j) \\ D_{i-1,j} + 1 \\ D_{i,j-1} + 1 \end{array} \right\}$$

06.07.2011 Theory 1 - Edit distance 12

Order of computation for the edit distance

06.07.2011 Theory 1 - Edit distance 13

Algorithm for the edit distance

Algorithm edit_distance
Input: two strings $A = a_1 \dots a_m$ and $B = b_1 \dots b_n$
Output: the matrix $D = (D_{ij})$

```

1  $D[0,0] := 0$ 
2 for  $i := 1$  to  $m$  do  $D[i,0] = i$ 
3 for  $j := 1$  to  $n$  do  $D[0,j] = j$ 
4 for  $i := 1$  to  $m$  do
5   for  $j := 1$  to  $n$  do
6      $D[i,j] := \min( D[i-1,j] + 1,$ 
7                    $D[i,j-1] + 1,$ 
8                    $D[i-1,j-1] + c(a_i,b_j) )$ 

```

06.07.2011 Theory 1 - Edit distance 14

Example

			a	b	a	c
	0	1	2	3	4	
b	1					
a	2					
a	3					
c	4					

06.07.2011 Theory 1 - Edit distance 15

Example

	j	0	1	2	3	4
			a	b	a	c
i	0	0	1	2	3	4
1	b	1				
2	a	2				
3	a	3				
4	c	4				

06.07.2011 Theory 1 - Edit distance 16

Example

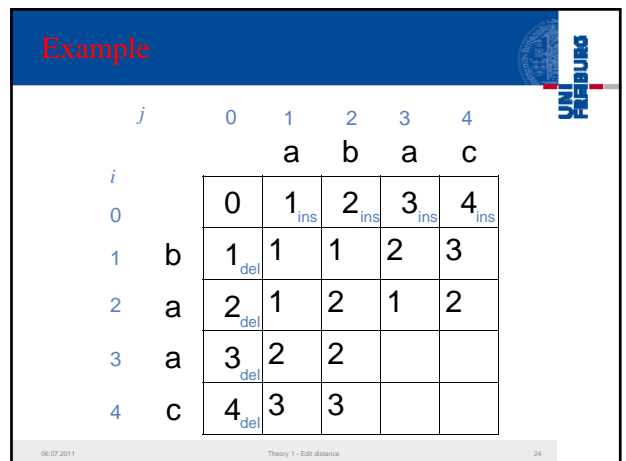
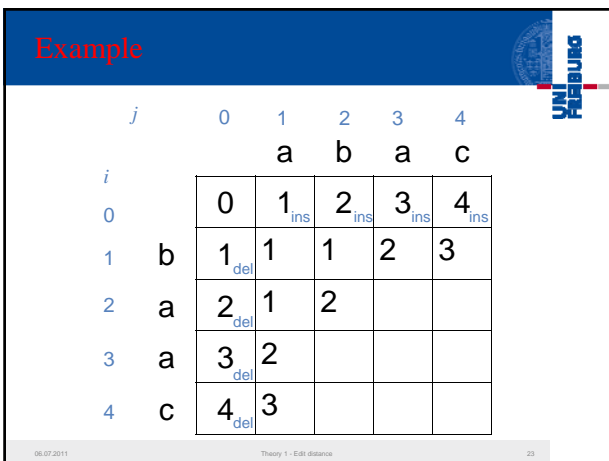
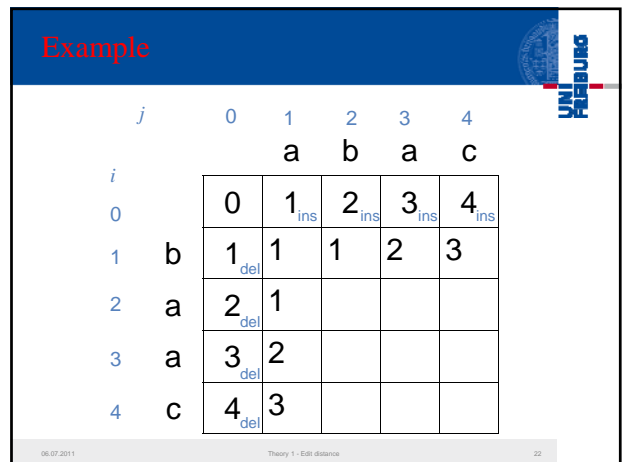
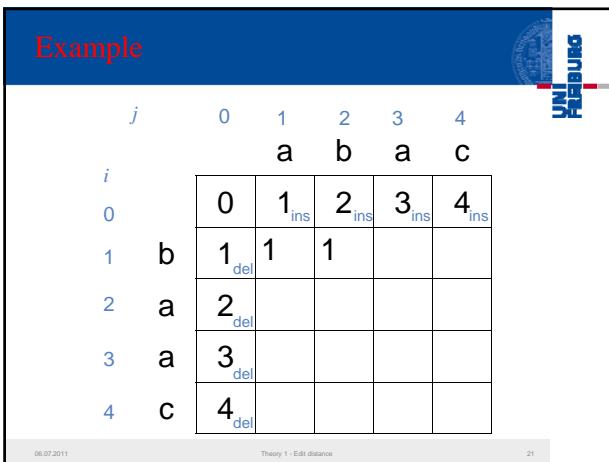
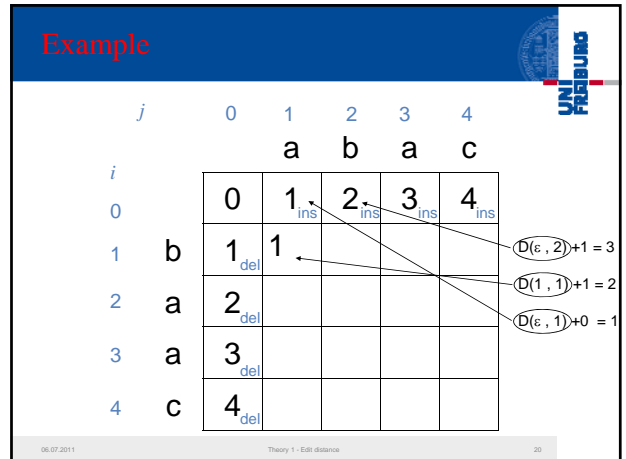
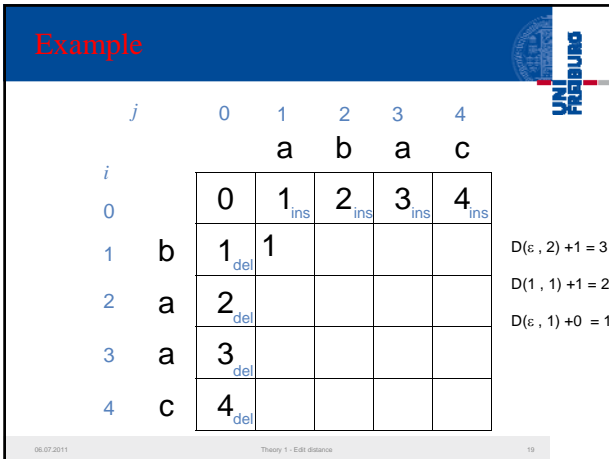
	j	0	1	2	3	4
			a	b	a	c
i	0	0	1 _{ins}	2 _{ins}	3 _{ins}	4 _{ins}
1	b	1 _{del}				
2	a	2 _{del}				
3	a	3 _{del}				
4	c	4 _{del}				

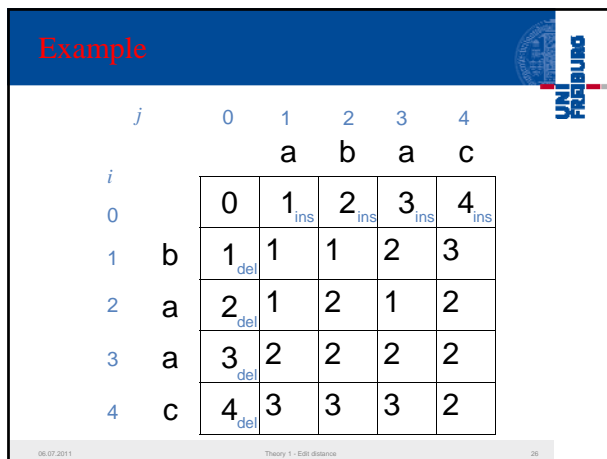
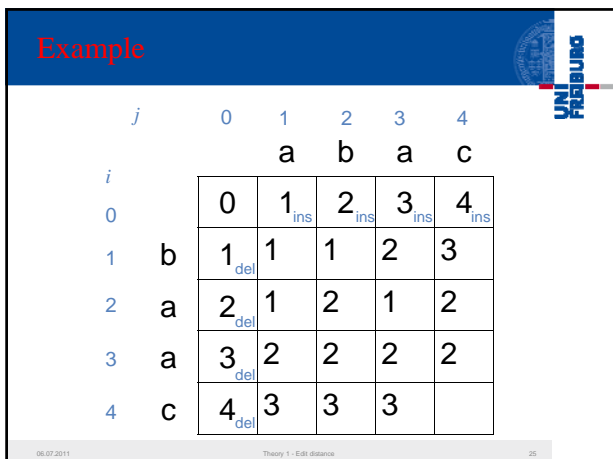
06.07.2011 Theory 1 - Edit distance 17

Example

	j	0	1	2	3	4
			a	b	a	c
i	0	0	1 _{ins}	2 _{ins}	3 _{ins}	4 _{ins}
1	b	1 _{del}	1			
2	a	2 _{del}				
3	a	3 _{del}				
4	c	4 _{del}				

06.07.2011 Theory 1 - Edit distance 18



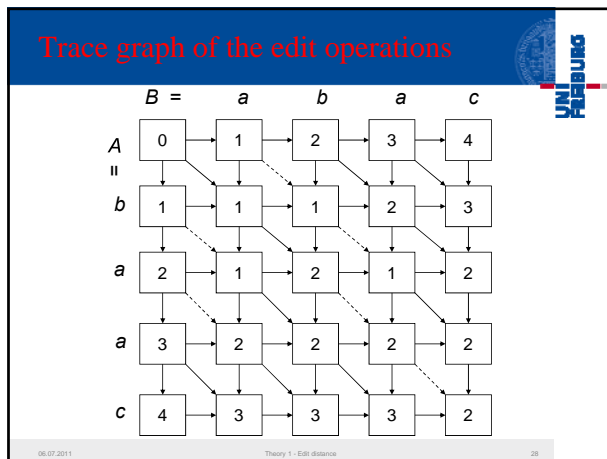


Computation of the edit operations

```

Algorithm edit_operations (i,j)
Input: matrix D (computed)
1 if i = 0 and j = 0 then return
2 if i ≠ 0 and D[i,j] = D[i - 1, j] + 1
3 then „delete a[j]“
4   edit_operations (i - 1, j)
5 else if j ≠ 0 and D[i,j] = D[i, j - 1] + 1
6 then „insert b[j]“
7   edit_operations (i, j - 1)
8 else
9   /* D[i,j] = D[i - 1, j - 1] + c(a[i], b[j]) */
10  „replace a[i] by b[j]“
    edit_operations (i - 1, j - 1)

Initial call: edit_operations(m,n)
    
```



Sub-graph of the edit operations

Trace graph: All possible traces which transform A into B, directed edges from vertex (*i, j*) to (*i + 1, j*), (*i, j + 1*) and (*i + 1, j + 1*).

Weights of the edges represent the edit costs.

Costs are monotonically increasing along an optimal path.

Each path from the upper left corner to the lower right corner, with monotonically increasing costs, represents an optimal trace.

Dynamic programming

- The algorithm just presented is an example of **dynamic programming**, an algorithm design technique often used for optimization problems.
- Generally usable for recursive approaches if the same partial solutions are required more than once
- Approach: store partial results in a table
- Advantage: better time complexity, often polynomial instead of exponential