



Algorithms Theory

14 – Dynamic Programming (4)

Edit distance

Approximate string matching

Sequence alignment

P.D. Dr. Alexander Souza

Dynamic programming

- Algorithm design technique, often applied to optimization problems
- Generally suitable for recursive approaches, when solutions to subproblems are required repeatedly.
- Approach: maintain a table of subproblem solutions
- Advantage: improved running time; often polynomial instead of exponential

Two different approaches

Bottom-up:

- + the table is maintained in an efficient way, time saving
- + subproblems are solved in a special, optimized order, space saving
- extensive rewriting of the original program code is necessary
- possibly, unnecessary subproblems are solved

Top-down: (memoization)

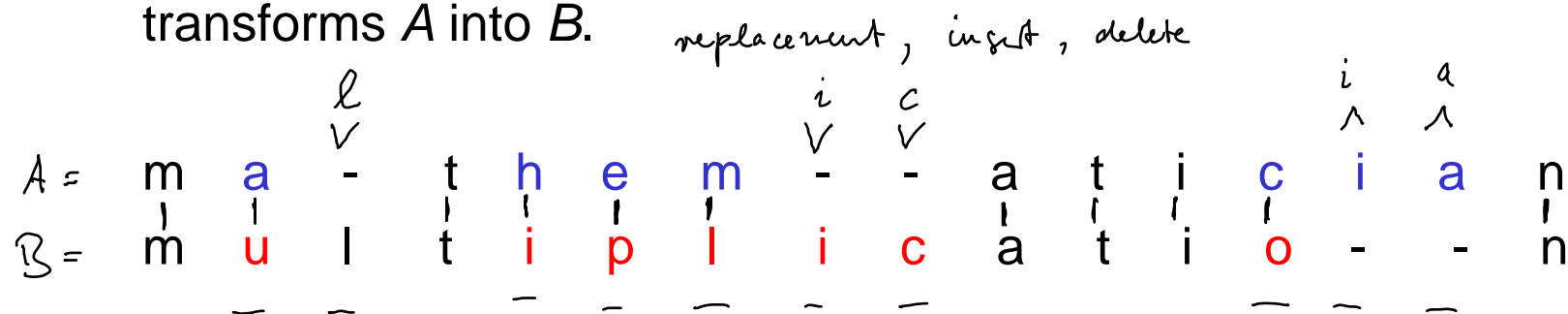
- + only slight modifications in the original program code are necessary
- + only those subproblems definitely required are solved
- separate table management is time consuming
- table size is often suboptimal

String matching problems

Here: Transform string A into string B

Edit distance: min. # operations.

For two given strings A and B , efficiently compute the edit distance $D(A,B)$ as well as a minimum sequence of edit operations that transforms A into B .



Cost of this sequence of operations ~~gives out~~: 10



String matching problems

Approximate string matching

For a given text T , a pattern P and a distance d , find all substrings P' of 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	G
-	-	-																	
-	-	-	C	A	C	G	T	G	G	-	-	-	-	-	-	-	-	-	-

align two strings such that the number of matching characters is as large as possible

Edit distance

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

Goal: Determine the minimum number $D(A,B)$ of edit operations required to transform A into B .

Edit operations:

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

$A =$ m a - t h e m - - a t i c i a n
 $B =$ m u l t i p l i c a t i o - - n

Edit distance

Unit-cost model:

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

$a = \varepsilon, b = \varepsilon$ possible

$a \neq \varepsilon, b \neq \varepsilon$: replacement

$a \neq \varepsilon, b = \varepsilon$: deletion

$a = \varepsilon, b \neq \varepsilon$: insertion

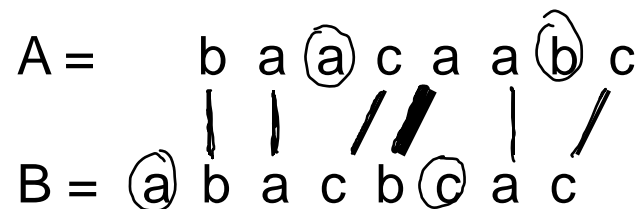
We assume the triangle inequality holds for c :

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

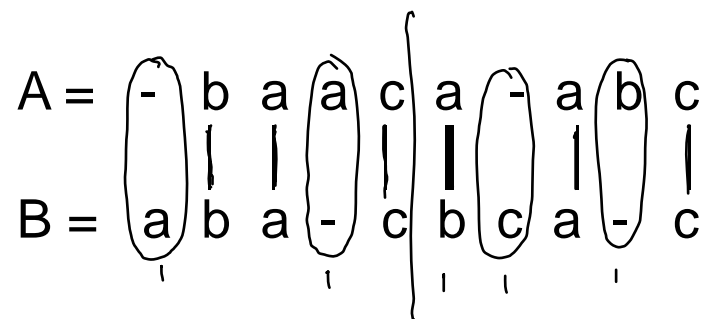
→ each character is changed at most once

Edit distance

Trace as representation of the sequence of edit operations:



or using indents:



Edit distance (costs) : 5

Trace is optimal only if two subtraces are also optimal.

Splitting an optimal trace yields two optimal subtraces
 → dynamic programming is suitable

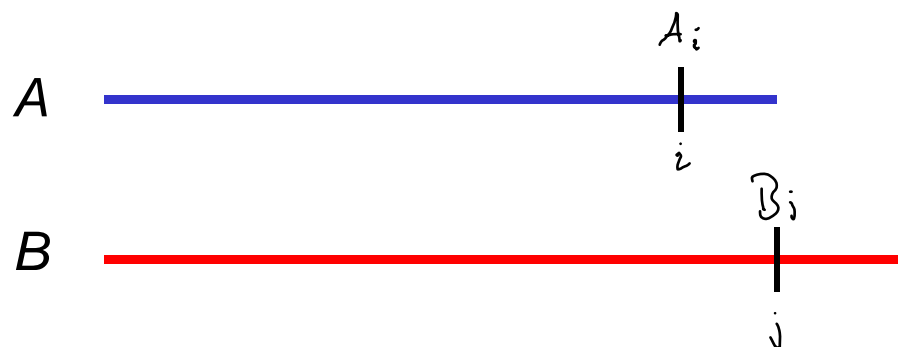
Computation of the edit distance

Prefixes

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

$$D_{i,j} = D(A_i, B_j)$$

We are interested in $D_{m,n}$



Computation of the edit distance

possible
↓

Three ways of ending a trace:

① a_m is replaced by b_n :

$$D_{m,n} = \underline{D_{m-1,n-1} + c(a_m, b_n)}$$

$$(1) \begin{array}{l} a_1 a_2 \dots a_{m-1} | a_m \\ b_1 b_2 \dots b_{n-1} | b_n \end{array}$$

② a_m is deleted: $D_{m,n} = \underline{D_{m-1,n} + 1}$

$$(2) \begin{array}{l} a_1 a_2 \dots a_{m-1} | a_m \\ b_1 b_2 \dots b_n | \text{---} \end{array}$$

③ b_n is inserted: $D_{m,n} = \underline{D_{m,n-1} + 1}$

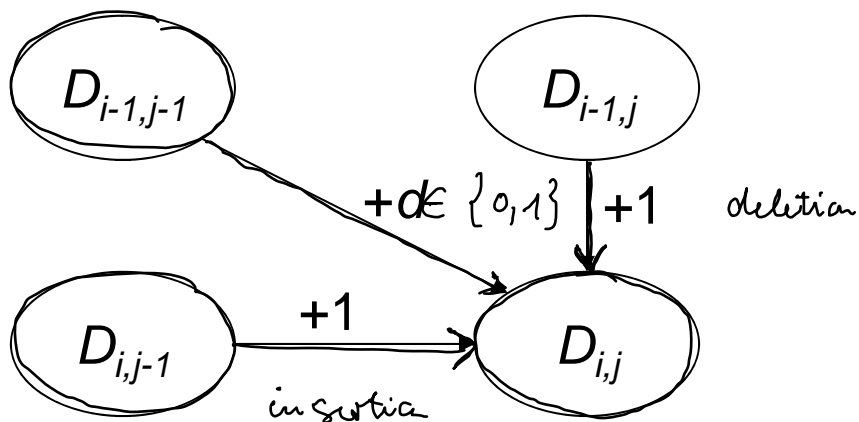
$$(3) \begin{array}{l} a_1 a_2 \dots a_m | \text{---} \\ b_1 b_2 \dots b_{n-1} | b_n \end{array}$$

Computation of the edit distance

Recurrence relation (for $m, n \geq 1$):

$$\underline{D_{m,n}} = \min \left\{ \begin{array}{l} D_{m-1,n-1} + c(a_m, b_n) \\ D_{m-1,n} + 1 \\ D_{m,n-1} + 1 \end{array} \right. \begin{array}{l} \text{replacement} \\ \text{deletion} \\ \text{insertion} \end{array}$$

→ computation of all $D_{i,j}$ required, $0 \leq i \leq m, 0 \leq j \leq n$.



Recurrence relation for the edit distance

Base cases:

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

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

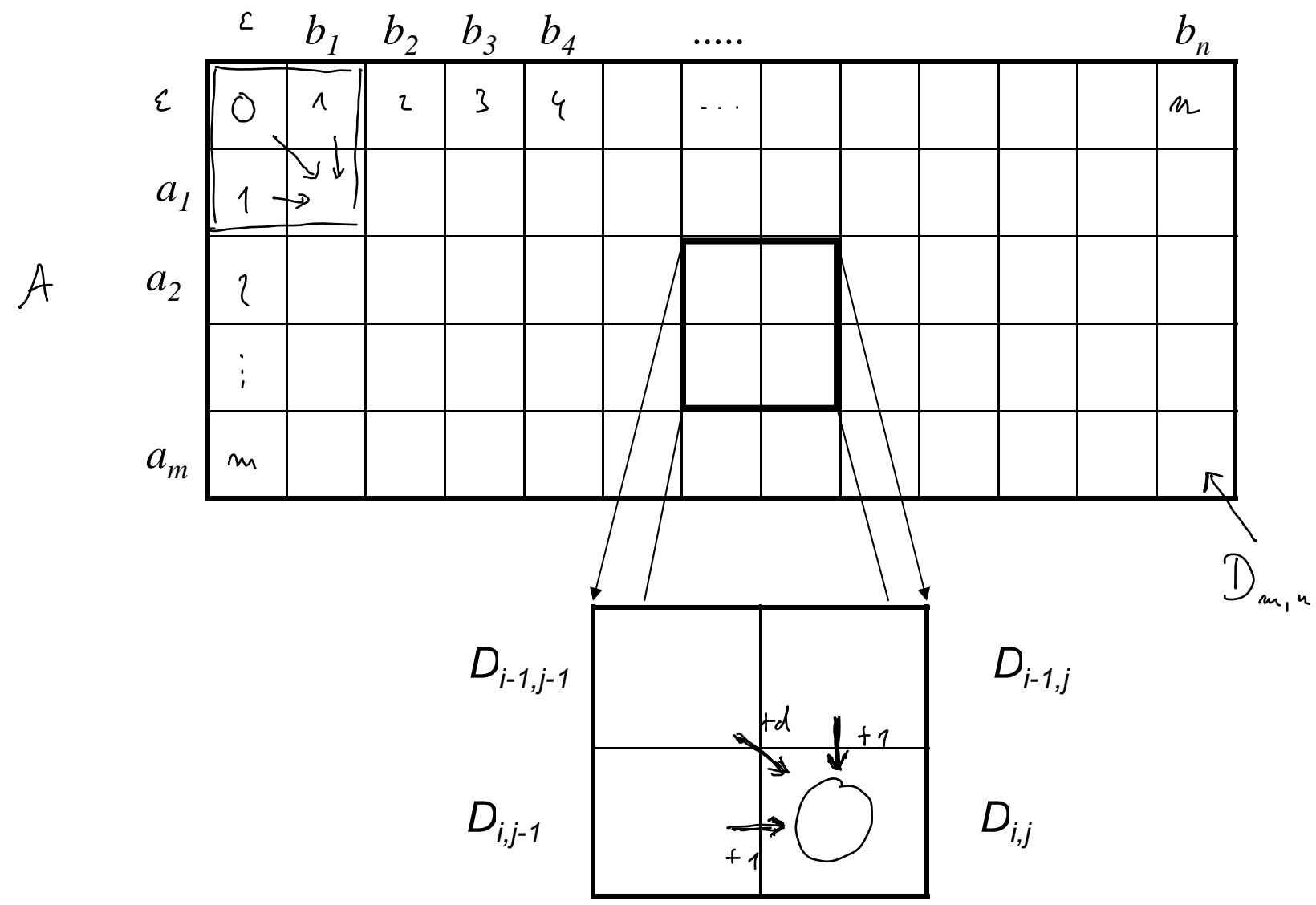
$$D_{i,0} = D(A_i, \varepsilon) = i$$

Recurrence relation:

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

Order of solving the subproblems

B



Algorithm for computing the edit distance

Algorithm *Edit-distance*

Input: two strings $A = a_1 \dots a_m$ and $B = b_1 \dots b_n$

Output: 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 $D[i-1, j-1] + \underline{c(a_i, b_j)})$

Running time $O(m + n + m \cdot n) = \underline{O(m \cdot n)}$

Example



3

a b a c

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

A

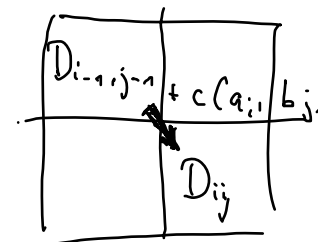
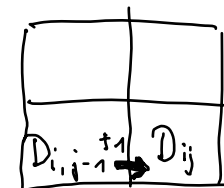
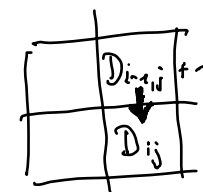
Computing the edit operations

Algorithm Edit-operations (i, j)

Input: matrix D (already computed)

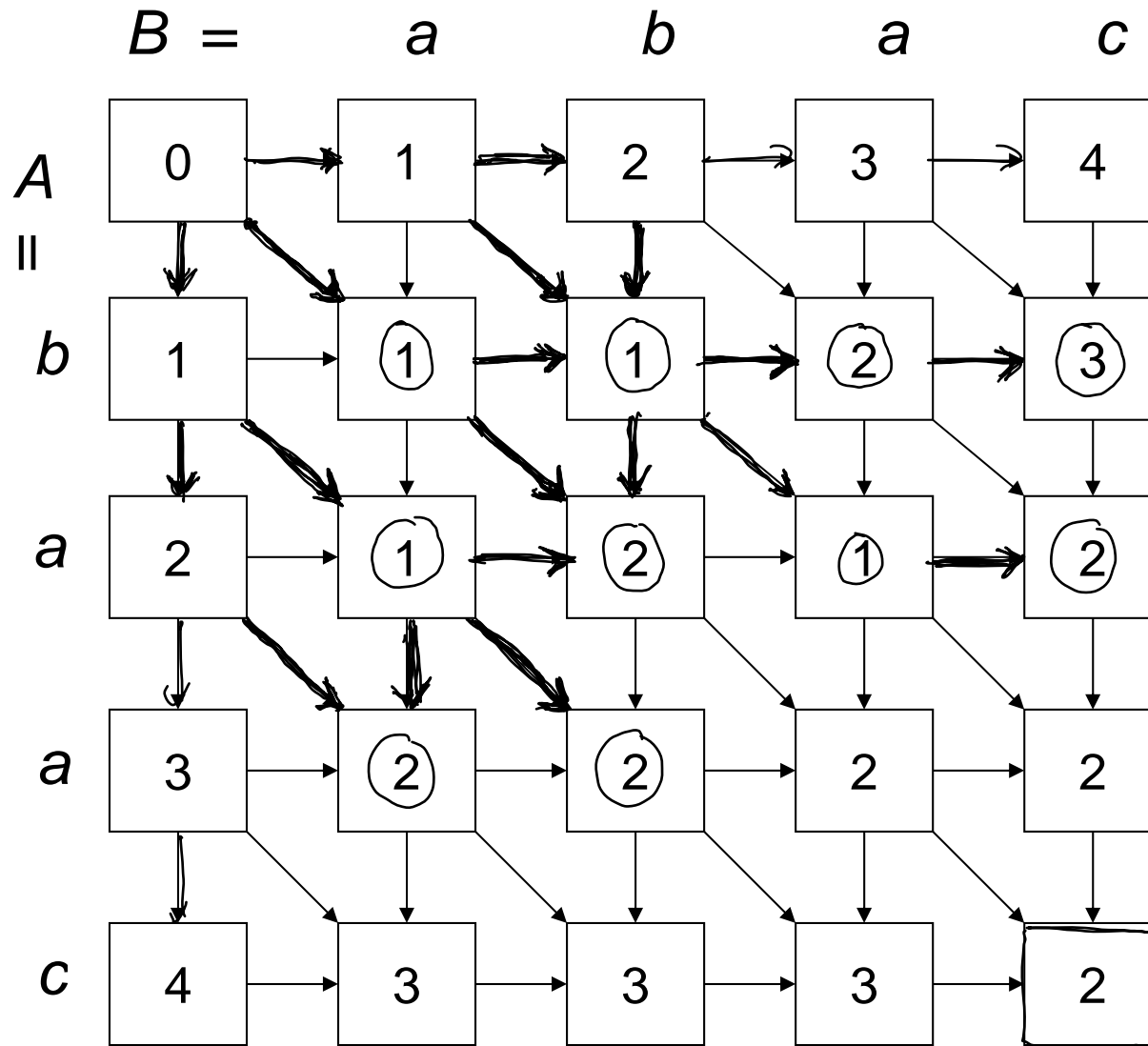
Output: sequence of edit operations

- 1 **if** $i = 0$ **and** $j = 0$ **then return**
- 2 **if** $i \neq 0$ **and** $D[i, j] = D[i - 1, j] + 1$
- 3 **then** Edit-operations ($i - 1, j$)
- 4 „delete $a[i]$ “
- 5 **else if** $j \neq 0$ **and** $D[i, j] = D[i, j - 1] + 1$
- 6 **then** Edit-operations ($i, j - 1$)
- 7 „insert $b[j]$ “
- 8 **else** /* $D[i, j] = D[i - 1, j - 1] + c(a[i], b[j])$ */
- 9 Edit-operations ($i - 1, j - 1$)
- 10 „replace $a[i]$ by $b[j]$ “

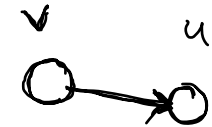


Initial call: Edit-operations(m, n)

Trace graph of the edit operations



feasible
transition

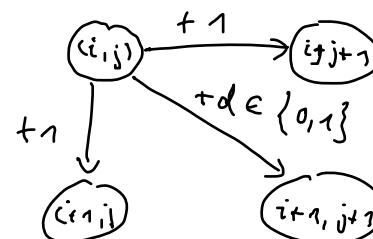


Trace graph of the edit operations

Trace graph:

Representation of all possible traces of operations that transform A into B . Directed edges from vertex (i, j) to vertices $(i + 1, j)$, $(i, j + 1)$ and $(i + 1, j + 1)$.

Edge weights represent the edit costs.



Along an optimal path, costs increase monotonically.

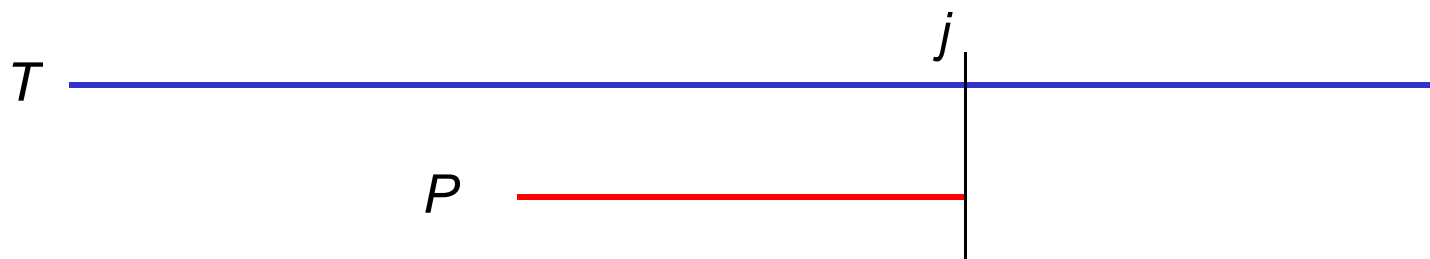
Each path from the upper left corner to the lower right corner with monotonically increasing costs represents an optimal trace.

Approximate string matching

Given: Two strings $T = t_1 t_2 \dots t_n$ (text) and $P = p_1 p_2 \dots p_m$ (pattern).

Goal: Find an interval $[j', j]$, $1 \leq j', j \leq n$, such that the substring $T_{j', j} = t_{j'} \dots t_j$ of T is the one with the highest similarity to the pattern P . Thus, for all other intervals $[k', k]$, $1 \leq k', k \leq n$:

$$\underline{\underline{D(P, T_{j', j})}} \leq \underline{\underline{D(P, T_{k', k})}}$$



Approximate string matching

Naive approach:

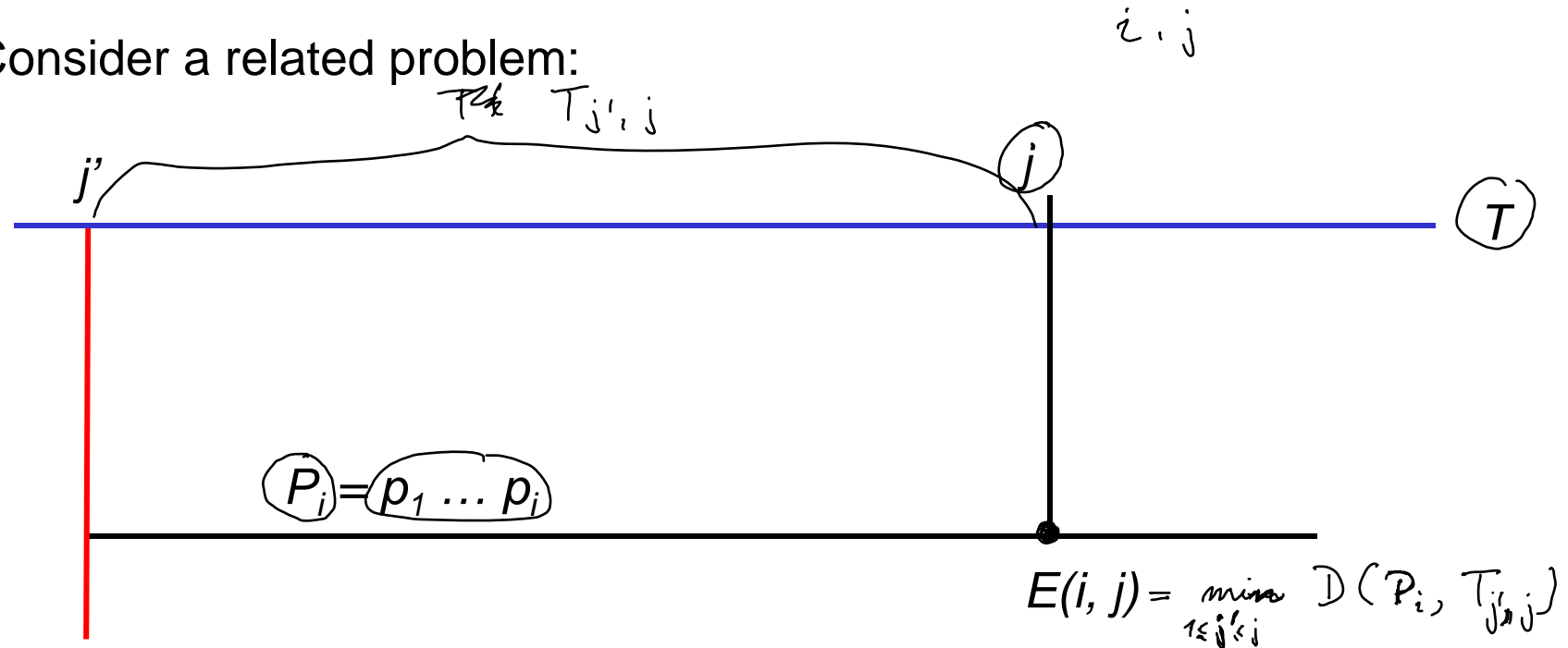
for all $1 \leq j', j \leq n$ **do**
 compute $D(P, T_{j', j})$
 choose the minimum

$$O\left(n^2 \cdot n \cdot m\right) = O\left(n^3 \cdot m\right)$$

Approximate string matching



Consider a related problem:



For each position j in the text and each position i in the pattern compute the minimum edit distance between P_i and any substring $T_{j',j}$ of T that ends at position j .

Approximate string matching

Method:

for all $1 \leq i \leq n$ do

determine j so that $D(P, T_{j',j})$ is minimized

For $1 \leq i \leq m$ and $1 \leq j \leq n$ let:

$$E_{i,j} = \min_{1 \leq j' \leq j+1} D(P_i, T_{j',j})$$

$j' = j+1$

$j' = j+1 : T_{j+1,j} = \epsilon$
empty string

Optimal trace:

$P_i =$	b	a	a	c	a	a	b	c
			//	//			/	/
$T_{j',j} =$	b	a	c	b	c	a	c	

Approximate string matching

Recurrence relation:

$$E_{i,j} = \min \left\{ \begin{array}{l} E_{i-1,j-1} + c(p_i, t_j), \\ E_{i-1,j} + 1, \\ E_{i,j-1} + 1 \end{array} \right. \begin{array}{l} \text{replacement} \\ \text{delete } p_i \\ \text{insert } p_i \end{array}$$

Remarks:

The index j' may differ for $E_{i-1,j-1}$, $E_{i-1,j}$ and $E_{i,j-1}$.

A subtrace of an optimal trace is an optimal subtrace.

*Compute the entries for increasing i
for each i , compute entries for increasing j*

Approximate string matching

Base cases:

$$E_{0,0} = E(\varepsilon, \varepsilon) = 0$$

$$E_{i,0} = E(P_i, \varepsilon) = i$$

whereas

$$E_{0,j} = E(\varepsilon, T_j) = 0$$

$$\min_{1 \leq j' \leq j+1} D(\varepsilon, T_{j',j}) = D(\varepsilon, \underbrace{T_{j+1,j}}_{\varepsilon}) = 0$$

Observation:

An optimal sequence of edit operations that transforms P into $T_{j',j}$ does not start with an insertion of character t_j .

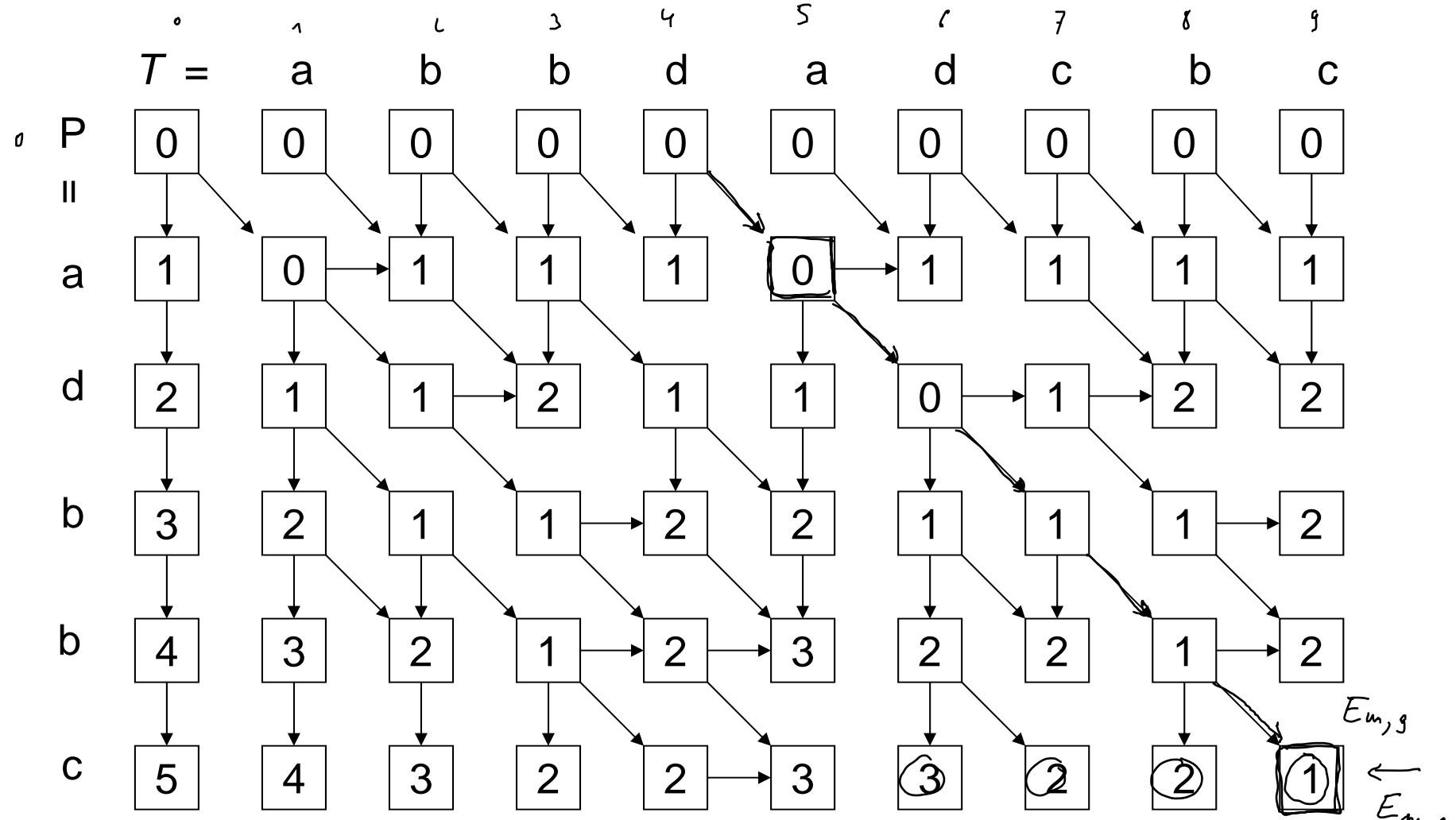
$$E_{m, \bullet} = E(P_m, T_{\bullet})$$



Approximate string matching

$T_{5,9} = edc b c$
 $P_m = a d b b c$

Dependency graph:



Find minimum entry in row m

Approximate string matching

Theorem:

If there is a path from $E_{0, j-1}$ to $E_{i, j}$ in the dependency graph, then $T_{j', j}$ is a substring of T that has the highest similarity to P_i , ending at position j and satisfying

$$D(P_i, T_{j', j}) \equiv E_{i, j}$$

We are interested in the entire pattern P_m , i.e., we are interested in the values $E_{m, \cdot}$. In particular $\min_j E_{m, j}$



Similarity of strings

Sequence alignment:

For two given DNA sequences, insert spaces (or dashes) such that, after placing the resulting strings one above the other, the number of matching characters is maximized.

G	A	⊖	C	G	G	A	T	T	A	G
G	A	T	C	G	G	A	A	T	A	G

Similarity of strings

Measuring the similarity of two characters:

example value	setting	in general
+ 1	for a <u>match</u>	} <u>s(a,b)</u>
- 1	for a <u>mismatch</u>	
- 2	for spaces	- c

Measuring the similarity of two sequences :

$$S(A, B) = \sum_{\substack{\text{pairs } (a_i, b_i) \\ a_i, b_i}} \text{similarity of } (a_i, b_i)$$

Goal: Find an alignment that maximizes the similarity.



Similarity of strings

Similarity $S(A,B)$ of two strings A and B

Operations:

1. Replacement of a character a by some character b :
Gain: $s(a,b)$
2. Deletion of a character from A , insertion of a character from B
Loss: $-c$

Goal:

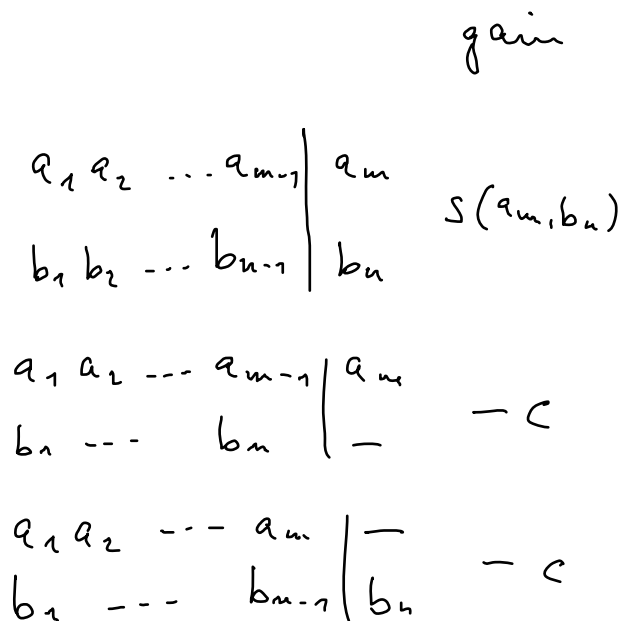
Find a sequence of operations that transforms A into B such that the total gain is maximized.

Similarity of strings

$$S_{i,j} = S(\underline{A_i}, \underline{B_j}), 0 \leq i \leq m, 0 \leq j \leq n$$

Recurrence relation:

$$S_{m,n} = \max \left(S_{m-1,n-1} + \overset{\text{replacement}}{s(a_m, b_n)}, S_{m-1,n} - c, S_{m,n-1} - c \right)$$



Base cases:

$$S_{0,0} = S(\varepsilon, \varepsilon) = 0$$

$$S_{0,j} = S(\varepsilon, B_j) = -jc$$

$$S_{i,0} = S(A_i, \varepsilon) = -ic$$

Most similar substrings

Given: Two strings $A = a_1 \dots a_m$ and $B = b_1 \dots b_n$.

Goal: Find two intervals $[i', j'] \subseteq [1, m]$ and $[k', l'] \subseteq [1, n]$ with

$$S(\underline{A_{i',j'}}, \underline{B_{k',l'}}) \geq S(\underline{A_{k',k'}}, \underline{B_{l',l'}}),$$

for all $[k', k'] \subseteq [1, m]$ and $[l', l'] \subseteq [1, n]$.

Naive approach:

for all $[i', j'] \subseteq [1, m]$ **and** $[k', l'] \subseteq [1, n]$ **do**
 compute $S(\underline{A_{i',j'}}, \underline{B_{k',l'}})$

Most similar substrings

Method:

for all $1 \leq i \leq m$, $1 \leq j \leq n$ do

compute i' and j' so that $S(\underline{A_{i',i}}, \underline{B_{j',j}})$ is maximized

For $0 \leq i \leq m$ and $0 \leq j \leq n$ let:

$$H_{i,j} = \max_{\substack{1 \leq i' \leq i+1, \\ 1 \leq j' \leq j+1}} S(A_{i',i}, B_{j',j})$$

$A_{i+1,i} = \varepsilon$
 $B_{j+1,j} = \varepsilon$

Optimal trace:

$$\begin{array}{r}
 A_{i',i} = \text{b a a c a - a b c} \\
 \quad \quad | \quad | \quad | \quad | \quad | \quad | \\
 B_{j',j} = \text{b a - c b c a - c}
 \end{array}$$

Most similar substrings

Recurrence relation:

$$H_{i,j} = \max \left\{ \begin{array}{l} H_{i-1,j-1} + s(a_i, b_j) \\ H_{i-1,j} - c \\ H_{i,j-1} - c \\ \textcircled{0} \end{array} \right\} \quad \text{empty substrings possible}$$

Base cases:

$$\begin{aligned} H_{0,0} &= H(\varepsilon, \varepsilon) = 0 \\ H_{i,0} &= H(A_i, \varepsilon) = 0 \\ H_{0,j} &= H(\varepsilon, B_j) = 0 \end{aligned}$$