# Similarity Search The String Edit Distance

#### Nikolaus Augsten

nikolaus.augsten@plus.ac.at Department of Computer Science University of Salzburg



WS 2025/26

Version November 6, 2025

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

1/31

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

2/31

String Edit Distance Motivation and Definition

# Outline

- String Edit Distance
  - Motivation and Definition
  - Brute Force Algorithm
  - Dynamic Programming Algorithm
  - Edit Distance Variants
- 2 Conclusion

#### Outline

- String Edit Distance
  - Motivation and Definition
  - Brute Force Algorithm
  - Dynamic Programming Algorithm
  - Edit Distance Variants
- 2 Conclusion

String Edit Distance Motivation and Definition

# Motivation

How different are

Augsten (Univ. Salzburg)

- hello and hello?
- hello and hallo?
- hello and hell?
- hello and shell?

Augsten (Univ. Salzburg) Similarity Search WS 2025/26 3 / 31

Similarity Search

WS 2025/26

4/31

String Edit Distance Motivation and Definition

#### What is a String Distance Function?

#### Definition (String Distance Function)

Given a finite alphabet  $\Sigma$ , a string distance function,  $\delta_s$ , maps each pair of strings  $(x, y) \in \Sigma^* \times \Sigma^*$  to a positive real number (including zero).

$$\delta_s: \Sigma^* \times \Sigma^* \to \mathbb{R}_0^+$$

•  $\Sigma^*$  is the set of all strings over  $\Sigma$ , including the empty string  $\varepsilon$ .

Augsten (Univ. Salzburg)

Similarity Search String Edit Distance Brute Force Algorithm WS 2025/26

# Outline

String Edit Distance

- Motivation and Definition
- Brute Force Algorithm
- Dynamic Programming Algorithm
- Edit Distance Variants
- 2 Conclusion

String Edit Distance Motivation and Definition

# The String Edit Distance

#### Definition (String Edit Distance)

The *string edit distance* between two strings, ed(x, y), is the minimum number of character insertions, deletions and replacements that transforms x to y.

- Example:
  - hello-hallo: replace e by a
  - hello→hell: delete o
  - hello→shell: delete o, insert s
- Also called Levenshtein distance.<sup>1</sup>

<sup>1</sup>Levenshtein introduced this distance for signal processing in 1965 [Lev65].

Similarity Search

String Edit Distance Brute Force Algorithm

Augsten (Univ. Salzburg)

WS 2025/26

6/31

#### Gap Representation

- Gap representation of the string transformation  $x \to y$ : Place string x above string y
  - with a gap in x for every insertion,
  - with a gap in y for every deletion,
  - with different characters in x and y for every replacement.
- Any sequence of edit operations can be represented with gaps.
- Example:

```
hallo
shell
```

- insert s
- replace a by e
- delete o

Augsten (Univ. Salzburg) Similarity Search WS 2025/26 Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

#### Deriving the Recursive Formula: Optimal Substructure

- Recursive solution: is applicable only to problems with optimal substructure property.
- Optimal substructure property of a problem:
  - optimal solution to larger problem computable from the optimal solutions of subproblems
- Examples:
  - Shortest path has optimal substructure: If a is on shortest path P from a to  $b \Rightarrow$  the section  $a \rightarrow c$  on P is the shortest path between a and c.
  - Longest simple<sup>2</sup> path does *not* have optimal substructure. Counter example [CLRS09]: Consider longest path  $q \rightarrow t$  and subpath  $q \rightarrow r$ .



<sup>2</sup>i.e., the path has no cycles

Augsten (Univ. Salzburg)

WS 2025/26

String Edit Distance Brute Force Algorithm

# Deriving the Recursive Formula: Optimal Substructure

#### Proof: Optimal Substructure String Edit Distance (by contradiction)

• Last column contributes with c = 0 or c = 1 to cost of gap(x, y), thus:

$$cost(gap(x, y)) = cost(gap(x', y')) + c$$

• Assume gap(x', y') is not optimal, i.e., cost(gap(x', y')) > ed(x', y'). Let  $gap^*(x', y')$  be the respective gap representation:

$$cost(gap^*(x', y')) = ed(x', y') < cost(gap(x', y'))$$

• By extending  $gap^*(x', y')$  with the last column, we get a gap representation  $gap^*(x, y)$  with cost below ed(x, y), which contradicts the definition of the edit distance.

$$cost(gap^*(x,y)) = cost(gap^*(x',y')) + c$$
  
$$< cost(gap(x',y')) + c = ed(x,y)$$

# Deriving the Recursive Formula: Optimal Substructure

#### Lemma (Optimal Substructure of String Edit Distance Problem)

Given a gap representation, gap(x, y), between two strings x and y, such that the cost of gap(x, y) is the string edit distance ed(x, y). If we remove the last column of gap(x, y), then the gap representation of the remaining columns, gap(x', y'), has cost ed(x', y') between the resulting substrings, x' and y'.

- Example:
  - x = hallo, y = shell, cost(gap(x, y)) = ed(x, y) = 3
  - x' = hall, y = shell, cost(gap(x', y')) = ed(x', y') = 2

Augsten (Univ. Salzburg)

WS 2025/26

String Edit Distance Brute Force Algorithm

# Deriving the Recursive Formula

Example:

- Notation:
  - x[1...i] is the substring of the first i characters of x ( $x[1...0] = \varepsilon$ )
  - x[i] is the *i*-th character of x
- Recursive Formula:

Augsten (Univ. Salzburg)

$$\begin{array}{rcl} \operatorname{ed}(\varepsilon,\varepsilon) & = & 0 \\ \operatorname{ed}(x[1..i],\varepsilon]) & = & i \\ \operatorname{ed}(\varepsilon,y[1..j]) & = & j \\ \operatorname{ed}(x[1..i],y[1..j]) & = & \min(\operatorname{ed}(x[1..i-1],y[1..j-1])+c, \\ & & \operatorname{ed}(x[1..i-1],y[1..j])+1, \\ & & \operatorname{ed}(x[1..i],y[1..j-1])+1) \end{array}$$

where c = 0 if x[i] = y[j], otherwise c = 1.

String Edit Distance Brute Force Algorithm

#### Brute Force Algorithm

#### ed-bf(x, y)

$$\begin{split} m &= |x|, \ n = |y| \\ \text{if } m &= 0 \text{ then return } n \\ \text{if } n &= 0 \text{ then return } m \\ \text{if } x[m] &= y[n] \text{ then } c = 0 \text{ else } c = 1 \\ \text{return } \min(\text{ed-bf}(x,y[1\dots n-1]) + 1, \\ &= \text{ed-bf}(x[1\dots m-1],y) + 1, \\ &= \text{ed-bf}(x[1\dots m-1],y[1\dots n-1]) + c) \end{split}$$

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

String Edit Distance Dynamic Programming Algorithm

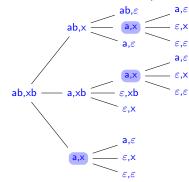
#### Outline

- String Edit Distance
  - Motivation and Definition
  - Brute Force Algorithm
  - Dynamic Programming Algorithm
  - Edit Distance Variants
- 2 Conclusion

String Edit Distance Brute Force Algorithm

#### Brute Force Algorithm

• Recursion tree for ed-bf(ab, xb):



- Exponential runtime in string length :-(
- Observation: Subproblems are computed repeatedly (e.g. ed-bf(a, x) is computed 3 times)
- Approach: Reuse previously computed results!

Augsten (Univ. Salzburg)

WS 2025/26

String Edit Distance Dynamic Programming Algorithm

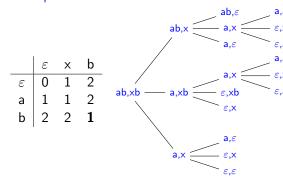
# Dynamic Programming Algorithm - Top Down

- Store distances between all prefixes of x and y
- Use matrix  $C_{0..m,0..n}$  with

$$C_{i,j} = \operatorname{ed}(x[1 \dots i], y[1 \dots j])$$

where  $x[1..0] = y[1..0] = \varepsilon$ .

• Example:



String Edit Distance Dynamic Programming Algorithm

#### Dynamic Programming Algorithm - Bottom Up

#### ed-dyn(x, y)

$$\begin{split} C: & \mathit{array}[0..|x|][0..|y|] \\ & \textbf{for } i = 0 \textbf{ to } |x| \textbf{ do } C[i,0] = i \\ & \textbf{for } j = 1 \textbf{ to } |y| \textbf{ do } C[0,j] = j \\ & \textbf{for } j = 1 \textbf{ to } |y| \textbf{ do } \\ & \textbf{ for } i = 1 \textbf{ to } |x| \textbf{ do } \\ & \textbf{ if } x[i] = y[j] \textbf{ then } c = 0 \textbf{ else } c = 1 \\ & C[i,j] = \min(C[i-1,j-1] + c, \\ & C[i-1,j] + 1, \\ & C[i,j-1] + 1) \end{split}$$

Augsten (Univ. Salzburg)

String Edit Distance Dynamic Programming Algorithm

ins  $\rightarrow$ 

WS 2025/26

#### Understanding the Solution

Exam

mple:			ren 📐	ε	m	0	n	d
			ε	0 <	—1 ÷	<b>−</b> 2 ←	—3 ÷	<b>-4</b>
	moon mond	del	m	î	0 +	-1 <	—2 ÷	-3
		$\downarrow$	o	2	1	0	-1 <	-2
			o	3	2	1	1 +	_2
			n	4	↑ 3	↑ <sup>1</sup>	1 +	_2

- Each arrow represents an edit operation with minimal cost
- Cost 2 in cell (n, d) can either result from replacing n by d (diagonal arrow) or by inserting d (horizontal arrow)
- Each path from bottom right to top left corner represents a valid set of edit operations

String Edit Distance Dynamic Programming Algorithm

# Dynamic Programming Algorithm - Properties

- Complexity:
  - O(mn) time (nested for-loop)
  - O(mn) space (the  $(m+1)\times(n+1)$ -matrix C)
- Improving space complexity (assume m < n):
  - we need only the previous column to compute the next column
  - we can forget all other columns
  - $\Rightarrow O(m)$  space complexity

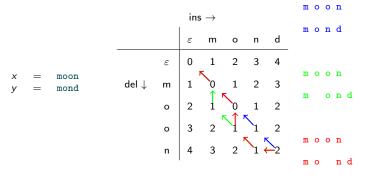
Augsten (Univ. Salzburg)

String Edit Distance Dynamic Programming Algorithm

WS 2025/26

#### Understanding the Solution

• Example:



- Solution 1: replace n by d and (second) o by n in x
- Solution 2: insert d after n and delete (first) o in x
- Solution 3: insert d after n and delete (second) o in x

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

19/31

Augsten (Univ. Salzburg)

String Edit Distance Dynamic Programming Algorithm

#### Dynamic Programming Algorithm

#### $ed-dyn^+(x,y)$

```
col_0: array[0..|x|]
col_1 : array[0..|x|]
for i = 0 to |x| do col_0[i] = i
for j = 1 to |y| do
   col_1[0] = j
   for i = 1 to |x| do
       if x[i] = y[j] then c = 0 else c = 1
       col_1[i] = \min(col_0[i-1] + c,
                      col_1[i-1]+1,
                      col_0[i] + 1
   col_0 = col_1
```

Augsten (Univ. Salzburg)

WS 2025/26

WS 2025/26

#### Distance Metric

#### Definition (Distance Metric)

A distance function  $\delta$  is a *distance metric* if and only if for any x, y, z the following hold:

- $\delta(x, y) = 0 \Leftrightarrow x = y$  (identity)
- $\delta(x, y) = \delta(y, x)$  (symmetric)
- $\delta(x, y) + \delta(y, z) \ge \delta(x, z)$  (triangle inequality)

#### **Examples:**

- the Euclidean distance is a metric
- d(a, b) = a b is not a metric (not symmetric)

Outline

- String Edit Distance
  - Motivation and Definition
  - Brute Force Algorithm
  - Dynamic Programming Algorithm
  - Edit Distance Variants
- 2 Conclusion

Augsten (Univ. Salzburg)

String Edit Distance Variants

# Introducing Weights

• Look at the edit operations as a set of rules with a cost:

$$lpha(arepsilon,b) = \omega_{ins} \qquad ext{(insert)} \ lpha(a,arepsilon) = \omega_{del} \qquad ext{(delete)} \ lpha(a,b) = \begin{cases} \omega_{rep} & \text{if } a \neq b \\ 0 & \text{if } a = b \end{cases} \qquad ext{(replace)}$$

where  $a, b \in \Sigma$ , and  $\omega_{ins}, \omega_{del}, \omega_{rep} \in \mathbb{R}_0^+$ .

- Edit script: sequence of rules that transform x to y
- Edit distance: edit script with minimum cost (adding up costs of single rules)
- Example: so far we assumed  $\omega_{ins} = \omega_{del} = \omega_{rep} = 1$ .

Augsten (Univ. Salzburg) Similarity Search WS 2025/26 23 / 31 Augsten (Univ. Salzburg) Similarity Search WS 2025/26 24 / 31

#### Weighted Edit Distance

• Recursive formula with weights:

$$C_{0,0} = 0$$
  
 $C_{i,j} = \min(C_{i-1,j-1} + \alpha(x[i], y[j]), C_{i-1,j} + \alpha(x[i], \varepsilon), C_{i,j-1} + \alpha(\varepsilon, y[j]))$ 

where  $\alpha(a, a) = 0$  for all  $a \in \Sigma$ , and  $C_{-1,i} = C_{i,-1} = \infty$ .

• We can easily adapt the dynamic programming algorithm.

Augsten (Univ. Salzburg)

WS 2025/26

WS 2025/26

# Allowing Transposition

- Transpositions
  - switch two adjacent characters
  - can be simulated by delete and insert
  - typos are often transpositions
- New rule for transposition

$$\alpha(ab,ba) = \omega_{trans}$$

allows us to assign a weight different from  $\omega_{ins} + \omega_{del}$ 

• Recursive formula that includes transposition:

$$\begin{array}{rcl} C_{0,0} & = & 0 \\ C_{i,j} & = & \min(C_{i-1,j-1} + \alpha(x[i], y[j]), \\ & & C_{i-1,j} + \alpha(x[i], \varepsilon), \\ & & C_{i,j-1} + \alpha(\varepsilon, y[j]), \\ & & C_{i-2,j-2} + \alpha(x[i-1]x[i], y[j-1]y[j])) \end{array}$$

where  $\alpha(ab, cd) = \infty$  if  $a \neq d$  or  $b \neq c$ ,  $\alpha(a, a) = 0$  for all  $a \in \Sigma$ , and  $C_{-1,j} = C_{i,-1} = C_{-2,j} = C_{i,-2} = \infty$ .

#### Variants of the Edit Distance

- Unit cost edit distance (what we did so far):
  - $\omega_{\textit{ins}} = \omega_{\textit{del}} = \omega_{\textit{rep}} = 1$
  - $0 < ed(x, y) < \max(|x|, |y|)$
  - distance metric
- Hamming distance [Ham50, SK83]:
  - called also "string matching with k mismatches"
  - allows only replacements
  - $\omega_{rep} = 1$ ,  $\omega_{ins} = \omega_{del} = \infty$
  - $0 \le d(x,y) \le |x|$  if |x| = |y|, otherwise  $d(x,y) = \infty$
  - distance metric
- Longest Common Subsequence (LCS) distance [NW70, AG87]:
  - allows only insertions and deletions
  - $\omega_{ins} = \omega_{del} = 1$ ,  $\omega_{rep} = \infty$
  - $0 \le d(x, y) \le |x| + |y|$
  - distance metric
  - LCS(x, y) = (|x| + |y| d(x, y))/2

Augsten (Univ. Salzburg)

#### Example: Edit Distance with Transposition

• Example: Compute distance between x = meal and y = mael using the edit distance with transposition ( $\omega_{ins} = \omega_{del} = \omega_{rep} = \omega_{trans} = 1$ )

Augsten (Univ. Salzburg)

• The value in red results from the transposition of ea to ae.

String Edit Distance Variants

# Text Searching

- Goal:
  - search pattern p in text t(|p| < |t|)
  - allow k errors
  - match may start at any position of the text
- Difference to distance computation:
  - $C_{0,j} = 0$  (instead of  $C_{0,j} = j$ , as text may start at any position)
  - result: all  $C_{m,j} \leq k$  are endpoints of matches

Augsten (Univ. Salzburg)

Similarity Search

Conclusion

WS 2025/26

Summary

- Edit distance between two strings: the minimum number of edit operations that transforms one string into the another
- Dynamic programming algorithm with O(mn) time and O(m) space complexity, where  $m \le n$  are the string lengths.
- Basic algorithm can easily be extended in order to:
  - weight edit operations differently,
  - support transposition,
  - simulate Hamming distance and LCS,
  - search pattern in text with k errors.

# Example: Text Searching

• Example:

• Solutions: 3 matching positions with  $k \le 2$  found.

```
survey
surge
survey
surger
surve
surgery
```

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

Alberto Apostolico and Zvi Galill.

The longest common subsequence problem revisited.

Algorithmica, 2(1):315-336, March 1987.

Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein.

Introduction to Algorithms, 3rd Edition.

MIT Press. 2009.

Richard W. Hamming.

Error detecting and error correcting codes.

Bell System Technical Journal, 26(2):147-160, 1950.

Vladimir I. Levenshtein.

Binary codes capable of correcting spurious insertions and deletions of

Problems of Information Transmission, 1:8–17, 1965.

Saul B. Needleman and Christian D. Wunsch.

Augsten (Univ. Salzburg) WS 2025/26 31/31 Similarity Search

Augsten (Univ. Salzburg) Similarity Search

WS 2025/26

A general method applicable to the search for similarities in the amino acid sequence of two proteins.

Journal of Molecular Biology, 48:443-453, 1970.



David Sankoff and Josef B. Kruskal, editors.

Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison.

Addison-Wesley, Reading, MA, 1983.

Augsten (Univ. Salzburg) WS 2025/26 Similarity Search 31/31