

Report: String Similarity Matching (Q5)

Author: Syed Rahmath Ullah Hussaini

Guide: Dr. Ahmed Rimaz Faizabadi

Institution: Delloyd R&D

1. Introduction

String similarity matching is a fundamental problem in computer science, with applications in text comparison, error correction, license plate recognition, and data validation.

This project focuses on:

- Accepting two strings of **6–10 characters**.
 - Calculating **percentage similarity** between them.
 - Generating a **match report** showing which characters match and which do not.
 - Aligning strings when necessary to improve comparison.
-

2. Problem Statement

Given two input strings, the program must:

1. Compute a **similarity percentage** based on character differences using the **Levenshtein distance** algorithm.
 2. Align the strings to visually show insertions, deletions, and substitutions.
 3. Identify **matching** and **non-matching** characters.
 4. Output a **comprehensive match report** with aligned strings and character-level details.
-

3. Methodology

3.1 Levenshtein Distance

The Levenshtein distance between two strings is the **minimum number of single-character edits** (insertions, deletions, or substitutions) required to change one string into the other.

- A dynamic programming matrix is initialized with dimensions $(\text{len}(\text{str1})+1) \times (\text{len}(\text{str2})+1)$.
- The matrix is filled iteratively to calculate minimum edits.
- Similarity is then computed as:

$$\text{Similarity (\%)} = \frac{\text{max_length} - \text{distance}}{\text{max_length}} \times 100$$

$$\text{Similarity (\%)} = \frac{\text{max_length} - \text{distance}}{\text{max_length}} \times 100$$

3.2 String Alignment

To improve comparison, the program aligns strings by tracing back the dynamic programming matrix:

- **Diagonal move** → characters match
- **Up move** → deletion
- **Left move** → insertion
- **Other** → substitution

Aligned strings are reversed and used to generate the match report.

3.3 Match Report

The match report includes:

1. **Aligned String 1**
2. **Aligned String 2**
3. **Matching characters**
4. **Non-matching characters** (shown as char1 vs char2)

This allows visual and quantitative assessment of string similarity.

4. Implementation

Main Functions:

1. `calculate_similarity(str1, str2)`
 - Computes Levenshtein distance and similarity percentage.
 - Returns (similarity, matrix) for alignment.
 2. `align_strings(str1, str2, matrix)`
 - Performs alignment using the distance matrix.
 - Generates the match report including matches and mismatches.
 3. `main()`
 - Handles interactive input or command-line arguments.
 - Displays similarity percentage and match report.
-

5. Sample Output

Enter the first string (6-10 characters): HELLO123

Enter the second string (6-10 characters): H3LLO12

Similarity Percentage: 87.50%

Aligned String 1: HELLO123

Aligned String 2: H-3LLO12

Matching characters: H, L, L, O, 1, 2

Non-matching characters: E vs 3, 3 vs -, 3 vs 3

Explanation:

- Characters H, L, L, O, 1, 2 match exactly.
 - E and 3 differ, and an insertion/deletion is represented by -.
-

6. Observations

- The program correctly computes similarity for any strings between 6–10 characters.
 - Alignment improves clarity by showing edits required to transform one string into another.
 - The match report can be used in applications such as:
 - License plate comparison
 - Spell checking
 - OCR error correction
 - DNA sequence comparison (with slight modification)
-

7. Conclusion

The string similarity program successfully:

- Computes **percentage similarity** using Levenshtein distance.
- Aligns strings to improve visual comparison.
- Generates a detailed **match report** highlighting matches and mismatches.

This implementation provides a robust foundation for any application requiring string comparison and error analysis.

8. References

1. Python sys module documentation: <https://docs.python.org/3/library/sys.html>
2. Levenshtein distance algorithm: https://en.wikipedia.org/wiki/Levenshtein_distance
3. Dynamic programming in string comparison: Cormen et al., *Introduction to Algorithms*, 3rd Edition.