

Design and Analysis of Algorithms

Manuel Lopez

100622700

Simple plagiarism detection: Enhance Levenshtein's Distance

The levenshtein distance between two documents find the minimum number of edits to transform the first document to the second. Levenshtein distance only allows for three types of edits Insert, Deletion, Substitution. To improve the levenshtein distance we preprocess the document to remove all the stop words and punctuation. Most search engine remove stop words because they give no real meaning to the algorithm. Levenshtein distance complexity is $O(n1 * n2)$.

How program levenshtein algorithm works:

Set **A** = length of document Original

Set **B** = length of document of Copy

Initialize a 2d array row size of **A** and Column size of **B** called distance

Observe the value of d2 (i from 1 to a).

Observe the value of d1 (j from 1 to b).

If the value at d2[i] is equals to value at d1[j], the cost becomes 0.

If the value at d2[i] does not equal d1[j], the cost becomes 2.

Set block of matrix Distance[d1, d2] of the matrix equal to the minimum of:

- the block immediately above add 1: Distance [d2-1, d1]
- the block immediately to the left add 1: Distance[d2, d1-1] +1.
- The block is diagonally above and to the left adds the cost: Distance[d2-1, d1-1] + cost.

To get similarity ratio:

$$\left(\left(\text{length}(\text{Document1}) + \text{length}(\text{Document2}) - \text{distance}[\text{row}][\text{col}] \right) / \left(\text{length}(\text{Document1}) + \text{length}(\text{Document2}) \right) \right)$$

$O(n1*n2)$

To preprocess Documents:

Memory complexity $O(n*n)$: All files are open as read byte then decoded to fix problems with utf-8 characters for comparison by pair then saved to an array.

We loop through the array removing all the Stop words and punctuation.

To remove punctuation I used the String library

After Documents are preprocessed we send the originalProcessed string and copyOriginal string to the levDistance function and get the similarity ratio between the two strings.

Loop through all the Documents while getting the ratio Save each ratio and then saving it to a csv File.

To run program:

levDistance folder contains **similarity.py** open terminal or command prompt and run Python3 similarity.py. All values are saved into a csv file called csv_out.csv

```
manuellopez@manuellopez-G551JW:~/Desktop/project$ python3 similarity.py
evaluation/1.txt vs abcde
[0.4308263695450325, 0.3492296404988995, 0.40065861690450055, 0.38959764474975467, 0.3487874465049929]
evaluation/2.txt vs abcde
[0.39710843373493976, 0.511522478277295, 0.3947217441193345, 0.3930576824910669, 0.3368807339449541]
evaluation/3.txt vs abcde
[0.40693430656934304, 0.36903039073806077, 0.8505376344086022, 0.4026974951830443, 0.3610133708655876]
evaluation/4.txt vs abcde
[0.37791527843883865, 0.35241301907968575, 0.4171848501978519, 0.46851385390428213, 0.3438749545619775]
evaluation/5.txt vs abcde
[0.38607888631090487, 0.3505683901723506, 0.3971475589687329, 0.38940657184894556, 0.4021390374331551]
```

References:

https://link.springer.com/chapter/10.1007/978-981-13-0755-3_6