

```
In [1]: 1 import numpy as np
        2 from IPython.display import display
```

Tracy Michaels

CSC 4740 Data Mining - Assignment 2

Implementing an algorithm for computing edit distance

Defining Edit Distance:

Edit distance is a way of quantifying how dissimilar two strings are to one another by counting the minimum number of operations required to transform one string into the other

Applications of Edit Distance:

- Natural Language Processing (Spell check, speech recognition)
- Computational Biology (DNA analysis)

Levenshtein Distance:

- For this assignment I will focus on Levenshtein distance, which is one type of edit distance
- This will calculate the minimum number of single-character edits (insertions, deletions, substitutions) required to transform one string into the other

3 operations for altering a string (each with a cost of 1 to perform):

- insert
- delete
- replace

last index in array will indicate the least number of operations to edit the strings to match

Algorithm:

```
In [2]: 1 def LevDis(s1: str, s2: str) -> int:
2         # construct array of size len(s1) by len(s2)
3         # fill row 0 and col 0 with incremental numbers from 0 to length of strings + 1 respectively
4         # adding one is done so that position [0][0] is a 0
5         lev_arr = np.zeros((len(s2) + 1, len(s1) + 1))
6
7         # filling initial values in row 0 and col 0
8         lev_arr[0] = [i for i in range(len(lev_arr[0]))]
9         lev_arr[:, 0] = [i for i in range(len(lev_arr[:, 0]))]
10
11
12        # iterate over array
13        for i in range(len(s2)):
14            for j in range(len(s1)):
15                # char matches, no incrementation
16                if s1[j] == s2[i]:
17                    # take the lowest value from 'backwards-adjacent' indecies
18                    # (i.e. [i-1][j-1], [i-1][j], or [i][j-1]) and place into target index
19                    # need +1 for target index as strings are 0-indexed but array will be
20                    # effectively 1-indexed since col[0] and row[0] are reserved for initial set-up
21                    lev_arr[i + 1, j + 1] = min(lev_arr[i+1, j], lev_arr[i, j+1], lev_arr[i, j])
22                else:
23                    # doesn't match, will need to perform insert, delete, or replace
24                    # increment lowest value from 'backwards-adjacent' indecies
25                    # (i.e. [i-1][j-1], [i-1][j], or [i][j-1]) then place into target index
26                    lev_arr[i + 1, j + 1] = min(lev_arr[i+1, j], lev_arr[i, j+1], lev_arr[i, j]) + 1
27
28        # displays array (might not work outside of jupyter as it depends on IPython.display)
29        display(lev_arr)
30        # return min edit distance
31        return lev_arr[-1, -1]
```

Example 1:

```
sString1 = "kitten"
sString2 = "sitting"
```

```
In [3]: 1 s1_1 = "kitten"
2         s2_1 = "sitting"
3
4         res_ex1 = LevDis(s1_1, s2_1)
5
6         print(f'Minimum edit distance between \'{s1_1}\'' and \'{s2_1}\'' is {res_ex1}')
```

```
array([[0., 1., 2., 3., 4., 5., 6.],
       [1., 1., 2., 3., 4., 5., 6.],
       [2., 2., 1., 2., 3., 4., 5.],
       [3., 3., 2., 1., 1., 2., 3.],
       [4., 4., 3., 1., 1., 2., 3.],
       [5., 5., 3., 2., 2., 2., 3.],
       [6., 6., 4., 3., 3., 3., 2.],
       [7., 7., 5., 4., 4., 4., 3.]])
```

Minimum edit distance between 'kitten' and 'sitting' is 3.0

Example 2:

```
sString1 = "GAMBOL"
sString2 = "GUMBO"
```

In [4]:

```
1 s1_2 = "GAMBOL"
2 s2_2 = "GUMBO"
3
4 res_ex2 = LevDis(s1_2, s2_2)
5
6 print(f'Minimum edit distance between \'{s1_2}\'' and \'{s2_2}\'' is {res_ex2}' )
```

```
array([[0., 1., 2., 3., 4., 5., 6.],
       [1., 0., 1., 2., 3., 4., 5.],
       [2., 1., 1., 2., 3., 4., 5.],
       [3., 2., 2., 1., 2., 3., 4.],
       [4., 3., 3., 2., 1., 2., 3.],
       [5., 4., 4., 3., 2., 1., 2.]])
```

Minimum edit distance between 'GAMBOL' and 'GUMBO' is 2.0

Observations:

One observation I found interesting was the relationship between the target index and those indices surrounding it, and how that relationship played a large role in the algorithm itself. The output makes sense as in each step the algorithm is selecting the minimum value from the select adjacent indices culminating in the last one being the one that would represent the least amount of steps required to transform one into the others.