**Project 3 - Final Computational Application**

Minerva University

CS110: Data Structures and Algorithms

Prof. B. Richard

December 26, 2022

**Decoding the relationships between genes**

In this assignment, I investigate the relationship between genes to establish the genealogical tree and estimate the probability of each possible change a string can undergo: insertion, deletion, or mutation.

1.  I created a function that calculated the longest common sequence between any 2 strings. I included the suggested test cases and my function passed all of them. Apart from the suggested test cases, I included three more. I examined the output of the function if the 2 strings given have no characters in common and if they are totally identical. I also examined the output of 2 empty strings and included them in my code guard clause to prevent the function of running in any case when at least of the strings is empty because we do not need it to run to know that there are no common characters.

I used the function to examine the set of strings given. As per my examination, the maximum number of combinations of possible longest common subsequences (lcs) is 13428744 between b and f, and the minimum number of combinations of lcs is 4 between e and f.

I created the matrix that shows the length of the longest common subsequence between any 2 strings. The matrix is represented in figure 1.

**2. Generate the matrix**

```
[[80 69 62 73 67 66 70]
 [69 88 62 76 62 61 72]
 [62 62 85 61 75 72 61]
 [73 76 61 84 62 63 79]
 [67 62 75 62 81 76 63]
 [66 61 72 63 76 87 65]
 [70 72 61 79 63 65 91]]
```

*Figure 1*: matrix of the longest common sequence between any 2 strings

Based on the matrix obtained above we can see that there are some stings that are more closely related to each other than others. We can also see that our matrix is mirrored. However, given that the strings have different lengths it does not necessarily reflect the real relationship between them. For example, if a string A and a string B have the longest common subsequence of 50 when A has a length of 50 and B has a length of 150 A would be much stronger related to B than B to A. There might be a third-string C that would have lengths of 300 and 90 in the longest common subsequence with B. In this case in the presented table, we will see that C is more strongly related to B than A. To consider the length of the strings I incorporated an additional part in my function to create the same table but with the information represented in percentages. As we can see in figure 2 if we take a look at the percentages the matrix would not be mirrored anymore and we will get a deeper insight in the relationship between the strings.

```
[[100  78  72  86  82  75  76]
 [ 86 100  72  90  76  70  79]
 [ 77  70 100  72  92  82  67]
 [ 91  86  71 100  76  72  86]
 [ 83  70  88  73 100  87  69]
 [ 82  69  84  75  93 100  71]
 [ 87  81  71  94  77  74 100]]
```

*Figure 2*: matrix represents the longest common sequence in percentages.

To examine the precise relationships between strings I would use the minimum edit distance which is also known as Levenshtein distance. This metric shows how many changes should be made to a string A so that it becomes a string B.

```
[[0, 22, 30, 14, 20, 28, 24],
 [22, 0, 36, 14, 35, 40, 23],
 [30, 36, 0, 34, 14, 23, 42],
 [14, 14, 34, 0, 31, 37, 14],
 [20, 35, 14, 31, 0, 14, 39],
 [28, 40, 23, 37, 14, 0, 40],
 [24, 23, 42, 14, 39, 40, 0]]
```

*Figure 3*: the figure represents the minimum edit distance.

I examined figure 3 and noted the two closest connections for each element. Examining only 1 relationship might not be enough to see a clear pattern and examining 3 connections of every node would result in too much noise in the data analyzed.

A is the most closely related to D and E.

A→ D, A→E

B is the most closely related to D and A

B→D, B→A

C is the most closely related to E and F

C→E, C→F

D is the most closely related to A and B and G. Here there are 3 connections since all of them have the same value.

D→A, D→G, D→B

E is the most closely related to C and F

E→C, E→F

F is the most closely related to E and C

F→E, F→C

G is the most closely related to D and B

G→D, G→B

I focus on detecting mutual connections. The set of nodes E C F are all the closest to each other (mutually connected) so we can assume they are a part of one subfamily. To understand who is the parent node and who are the children nodes we will take a look at the strength of the relationships. The longest distance will be between 2 brothers since they are connected with each other not directly but throught the common parent. Since 23 is the largest distance and it is between C and F they are brothers. Therefore E is a parent node of C and F.

Another set of 3 although not that distinct is B G D. The largest distance is between G and B so they are brothers. D has a strong connection with both of them so it is a parent node. A third strong connection the D node has is with A. So A is a parent node of D. A also has a strong connection with E, the previously found parent of another parent of the tree. Therefore we can say that A is a parent of both E and D. Therefore A is a grandparent of the tree.

**3. Examine the precise relationships between strings.**

**Local strategy:**

My local strategy will consist of creating all the possible trees and choosing the best one. I use the top-down approach. That means that I start from the top of the tree parent node and I

finish at the bottom level. I will take one step at a time and make a decision looking only at one row. I will use greedy properties by making a local optimal solution (choosing the best match from a single row) to reach global optimal solution.

I will choose each element as a parent and will build a tree for each of those parents. I will take a row that corresponds to the index of the parents. All the values in this row represent the minimum edit distance between the parent and every string. I will choose 2 elements with the smallest value. Since this is a root it cannot have parents so the closets element for it would be its children. After that, I will take each of those children as a parent and will find their children and continue until all the nodes are used.

When all the trees are constructed I will sum up all the edit distances between them to determine which tree has the lowest edit distance, thus is the best.

This strategy is very complex because we have to consider every combination possible and create every possible tree to find the one that would be optimal. It does not guarantee reaching the globally optimal solution since it does not consider the whole picture and makes only the locally optimal choice. For example, it might be possible that if the distance between the grandparent and parent is a bit bigger it would drastically decrease the distance between the parent and child which would result in the overall decrease of the edit distance of the tree. However, we will not be able to figure it out with this approach since we do not look at parent-child relationships when we find grandparent-parent relationships. Therefore it is not guaranteed to find the globally optimal solution with this strategy.

**Global strategy:**

My global strategy would be looking at the whole matrix instead of examining it row by row.  I will use the bottom-up approach by starting with the lower level (in this case it is grandchildren) and moving up to the parent node. I will find the highest edit distance on the matrix. The highest edit distance means that the elements are on opposite sides from each other. In our case of 7 elements they are located on the third level from different sides (have different parent nodes). When we found the combination of elements that produce the highest edit distance we can see which exact elements result in that number. Those elements are considered "cousins".  After that, we examine each of these "cousins" to find their "brother" nodes, the second node from the same parent. We know that the minimum edit value of the element is when we compare it to itself since no changes should be made. The second minimum value should be the connection with a parent or a connection with a child. Since we focus on the lowest level, the node does not have children so the second highest connection will be with its parent. The third highest connection will be with another child node of the same parent, its 'brother'. We find 'brothers' for each of these distant elements. As we have both children we can merge them to consider them as one entity. In those merged columns we find the highest element that will represent the most distant element to the children - "uncle". If we switch those elements we will get the set of parent nodes and children. After we finished those steps we have only 1 element left - the root node of the tree. We run loops to determine which element was not used. As we did so we will construct the tree in the form of lists. This strategy will help to get the globally optimal solution since we are looking at all relationships at the same time. It allowes us to make some sacrifices (less optimal solutions) on a local level to reach the best solution possible on a global level. The ability of this approach to reach a globally optimal solution was supported by

the evidence that the output of the global strategy function matched with the theoretically

hypothesized tree created above. On the other hand, the output from the function based on the

local approach did not produce the expected output.

The resulting tree from the global solution (an index-based version) is presented in figures 4 and
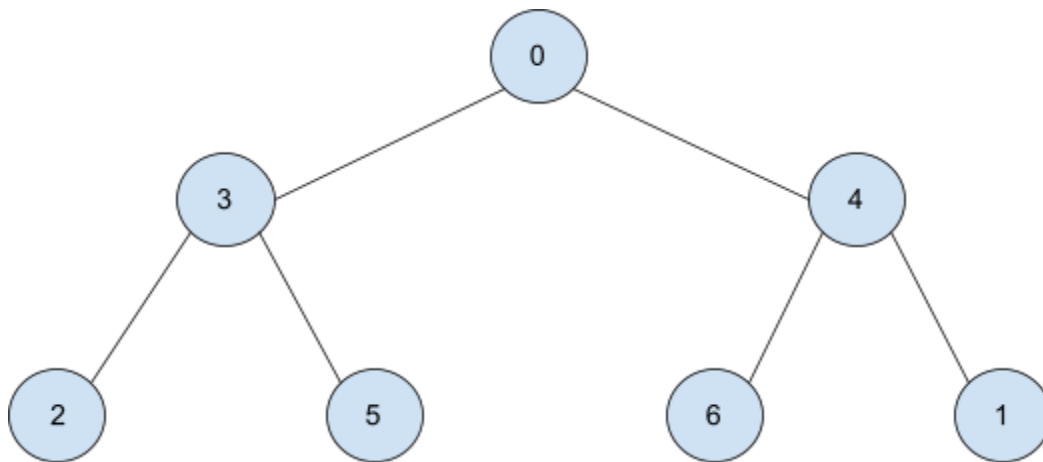
5.



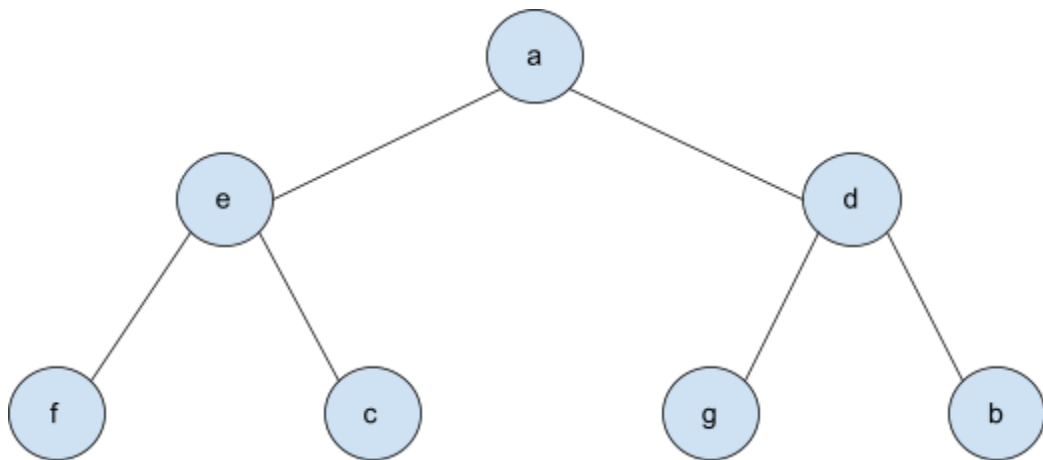*Figure 4*: index-based version of the global tree



*Figure 5*: letter-based version of the global tree,

However, the local approach is much easier to implement on a larger scale since the code is easy to scale. The global approach is very efficient but unfortunately, the scalability of my implementation is limited and it should be edited for a bigger scale.

**4. For each algorithmic approach, what is your solution's computational complexity to produce genealogy binary trees?**

Time complexities for both trees depend on $n$ - a number of genes and $m$ - the length of the genes. Both algorithms share in common the calculation of the Edit distance since it outputs the matrix both approaches work with. Edit distance has a time complexity of $O(n * m)$. Matrix is constructed using 2 loops so the time complexity for the matrix part will be $O(n^2)$. So the total time complexity of the edit distance calculation and matrix construction will be $O(n^2(n * m))$ (GeeksforGeeks, 2022b).

This is the only part where the algorithm's time complexity relies on the length of the string since after the matrix is calculated we rely solely on the integers in the matrix and depend on the number of the strings. As we see in the time complexity notation the number of the genes matters 3 times more than the length of the genes when it comes to calculating the time matrix which is the input for both trees.

     **Time complexity for the local approach:**

     Local approach will have comparatively worse time complexity because we have to deconstruct a tree starting with each parent node. It will scale as the number of genes increases because the more genes we add the more tree we will have to construct.

To calculate the tie complexity for the whole approach I will calculate the time complexity for each function consisting of:

- Find_min: it contains 2 for loops, 1 of them is based on the length of the row which is basically the number of genes n. Second repeats 2 times. So time complexity is about $O(2n)$

- Creating matrix time complexity was defined above and it is $O(n^2(n * m))$

- Cleaning table takes $O(n^2)$ because of the 2 for loops.

- Finding children will take about $O(n)$ because of the for loop. While loop although exists on a large scale the complexity it adds would not be significant.

- Greedy approach takes $O(n)$ because of the for loop

- Finding sum is $O(n)$ because of the for loop

Total time complexity is $O(2n + 3n + n^2 + n^2(n * m)) = O(5n + n^2 + n^2(n * m))$

Since 5n will not be that significant compared to exponential we can cancel it out.

$O(n^2 + n^3 + nm))$

So we will see exponential growth

Also, we will see more extreme growth if the number of genes increases and the length of genes stays the same the vice versa since n contributes to the time complexity much more than m.
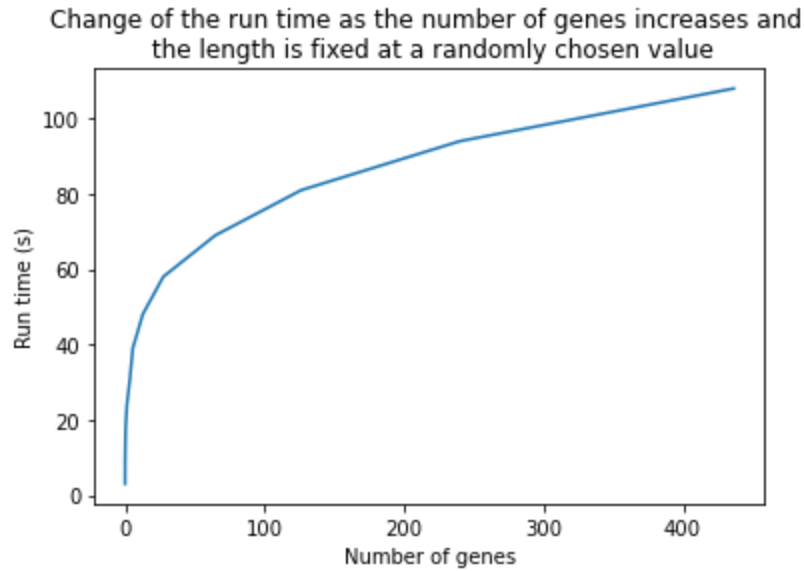
Change of the run time as the number of genes increases and
the length is fixed at a randomly chosen value



*Figure 6*: figure represents the growth of the run time of the greedy approach over time. The length of strings stays the same but the number of strings increases from 0 to 400.

Figure 6 shows that the dependency between the increase in the run time and the increase in the number of genes is exponential as it was hypothesized. It means that for each gene added to the tree the time to construct the tree will increase by 2+m times.

**For the global approach:**

Global approach would be more efficient than the local because we do not run the function for each node, so its time complexity is expected to be smaller that the one of the greedy.

- Constructing a matrix takes $O(n^2(n * m))$

- To find cousins we will need to do 1 for loop which will result in $O(n)$ time complexity.

- In finding the brother function we will need to do sorting 2 times. The time complexity of

the in-build sort function in Python in $O(2nlogn)$(GeeksforGeeks, 2022a).

- Finding parents have 2 loops one runs through all the rows and the second only through 2 elements. It results in the $O(2n)$ time complexity

- Merging will take $O(n^2)$ because of 2 loops

- Building the tree will take $O(2n)$

Total time complexity is

$$O(2n + 2n + n^2 + 2nlogn + n + n^2(n * m)) = 5n + 2nlogn + n^2(n * m))$$

If we follow the same principle as before and cancel out $5n$ we will end up with

$O(2nlogn + n^2(n * m))$ which is significantly smaller than the time complexity of the greedy approach.

So if we compare the scaling of both approaches for an increase of 1 gene greedy approach will increase by exponentially more time than the dynamic approach.

**5. How would you estimate the probabilities of insertions, deletions, and mutations?**

**a. Algorithmic approach for estimating the probabilities of insertion, deletion, and mutation**

We know the longest common subsequence for any two strings. It means that the rest that is not common was either added, deleted, or changed. My algorithmic strategy would focus on analyzing the part of the strings where they differ from each other. Let's assume that we know that string B was created as a result of the mutilation, insertion, and deletion of the elements into string A. (A is a parent string of B).

If string A is smaller than string B, the insertion would be found by simply finding the difference between the two strings since to get from A to B we need to insert additional elements. If we want to consider the opposite case, string B is smaller than string A we would assign the difference between the strings to the deletion since we need to remove some elements to get from B to A. The mutilation is equal to the length of the smaller string without the common values (LCS). Example:

String A has the length of 30 elements

String B has the length of 25 elements

The longest common subsequence of A and B is 10, so we subtract 10 from the length of both strings and will operate with the values we get.

$$30 - 10 = 20$$

$$25 - 10 = 15$$

We see that the string B is smaller than A. Since A is the ancestor of B the difference between the strings is the result of deletion.

$$20 - 15 = 5$$

As we subtract the remaining length of one element from another we get the number of deletions that were made.

The number of mutations is the part of A that was not deleted and does not match B. The number of mutations is equal to 15.

To find the probability we have to divide the number of all manipulations and divide the number of each manipulation by the total.

In this case:

$$total = 5 + 15 = 20$$

$$P_{deletion} = 5/20 = 0.25$$

$$P_{mutation} = 15/20 = 0.75$$

$$P_{insertion} = 0/20 = 0$$

The total number will be found by dividing the number of changes by the length of the string that was changed.

**b. Produce an estimation in Python that would take your algorithmic strategy into practice.**

I produced an estimation in Python. I obtained the following results:

$$P_{mutation} \approx 0.195$$

$$P_{deletion} \approx 0$$

$$P_{insertion} \approx 0.049$$

These probabilities seem reasonable since they are less than 5% . We also know that the minimum edit distance between the closest nodes in the tree is relatively small (14-20 edits) which serves as evidence supporting these results. The probability for deletion is 0 because the length of the tree increases by levels. This is supported by figure 1 in Appendix.

**6. List all the LOs and HCs you have exercised while working on this final assignment.**

a. LO Applications:

**#CodeReadability**

I added docstrings and comments where it was appropriate. All my docstrings were clean and formatted in the same way. I visually separated the description of the function and the input/output based on the feedback from the previous assignment. I assigned logical and self-explanatory names for my variables so that I did not have to write extra comments for more explanation. I also made sure that the code itself is clear and I avoided unnecessary steps that might be confusing. I kept my code as simple and succinct as possible.

**Word Count:** 91 words

**#Professionalism**

I did the explanation part of the assignment clearly. I have also proofread my assignment to avoid typos as I was advised in the previous assignment. I also proofread my comments and docstrings. I was also staying within the word count where it was specified. I formatted my assignment according to the APA style. I used professional language but at the same time avoided complicated expressions because my audience was a confused student and I did not want to confuse them even more.

**Word Count:** 83 words

#### #ComplexityAnalysis:

I applied this LO by analyzing each function both local and global approaches consist of and creating equations of the time complexity. I pointed out which parts of the functions contribute tot he time complexity and how they affect the run time. I explained the significance of time complexity in the real life - what it means for the time and change in the number of genes.

**Word Count:** 67 words

#### #ComputationalCritique

I made a strong application of this LO by providing efficiently connecting analytical and experimental analysis. I also ensured that I created appropriate experiments that provide useful insight into the efficiency of the code. I liked my reasoning on this LO to the LO #ComplexityAnalysis to make my application even stronger. I also explained which features of the global algorithm make it more efficient than the local algorithm.

**Word Count:** 68 words

#### #AlgoStratDataStrat

I explained my algorithmic strategy for both local and global algorithms in detail and went beyond a simple translation of the pseudocode to plain English. I exampled why I did what I did and how it improved my code. I discussed the usefulness of both local and global approaches. I pointed out the disadvantages of my implementation of the global strategy and an idea of how it could have been solved.

**Word Count:** 71 words

**#PythonProgramming**

All my code is working. I fixed all the bugs in my code. I wrote efficient code by using such strategies as guard clauses to prevent the function from running when it is not necessary. I avoided unnecessary loops that would negatively impact my time complexity. I also used multiple functions to separate pieces of the code and make it more clear.

**Word Count:** 62 words

b. HC Applications

**#dataviz:**

I applied this HC by producing high-quality data visualizations. I appropriately labeled the axis. I added an informative title and created captions for all the visualizations presented in this paper. I provided different types of data visualizations like graphs, line plots, and matrixes. I included multiple figures in this paper.

**Word Count:** 50 words

**#audience:**

I took into consideration my target audience of a confused student and explained everything in a detailed enough way. I avoided complex language and jargon. I also used examples to ease explanations. I used clear language and avoided vague statements.

**Word Count:** 40 words

**#deduction**:

I applied deductive reasoning when I was analytically analyzing a matrix to construct a tree. I defined premises and if the premises are true the conclusion should be true so I made the conclusion about the structure of the tree based on the proven premises about the parents. I concluded about the parents based on the distance in the interconnected triangles identified. I identified the connections in triangles by finding the mutually connected couples.

**Word Count:** 74 words

**#algorithms:**

I used this HC since developing the solution for the given problem will require strong algorithmic thinking. I accounted for all the potential edge cases and made sure that my code would be able to deal with them.  I identified an appropriate algorithmic strategy ( dynamical programming -  global strategy) and I justified my choice. Before writing the code I developed the pseudo code to have a between understanding of what I am aiming for. I was also drawing the execution of the code step by step while fixing bugs. I provided a thorough explanation of my code.

**Word Count:** 96 words

.

**Total Word Count:** 2823 words

Appendix:

```
1  from statistics import mean
2
3
4  def probabilities_calculation(set_strings, links):
5      mutation = []
6      deletion = []
7      insertion = []
8      strings = []
9      for i in set_strings:
10         strings.append(i[1])
11     for i in links:
12         str1 = strings[i[0]]
13         str2 = strings[i[1]]
14         print(len(str1)<len(str2))
15         p = probabilities(str1, str2, len(str1), len(str2))
16         mutation.append(p[0])
17         deletion.append(p[1])
18         insertion.append(p[2])
19     result = [mean(mutation), mean(deletion), mean(insertion)]
20     return result
21
22
23
24 probabilities_calculation(set_strings, links)
```

```
True
True
True
True
True
True
```

5]: [0.19491766580734352, 0.0, 0.049274625928333385]

*Figure 6*: This figure proves the plausibility of the 0% deletion probability since every other string in the

tree is longer than its parent.

Code:

```python
from Levenshtein import distance as lev
import numpy as np
"""
input: 2 strings
output: the length of the longest common sequence (integer)

it is based on 2 loops that check all the combination between every two elements of the strigs nd finds the highest
"""

def lcs(x , y):
    m = len(x)
    n = len(y)

    #create a table with dimensions corresponsing to the length of the strings
    L = [[None]*(n+1) for i in range(m+1)]

    #traverse throught all the combinations
    for i in range(m+1):
        for j in range(n+1):
            if i == 0 or j == 0 :
                L[i][j] = 0
            #if elements are the same add 1 to lcs and move diagonally
            elif x[i-1] == y[j-1]:
                L[i][j] = L[i-1][j-1]+1
            else:
            #assign the higest value between the one on the left and up.
                L[i][j] = max(L[i-1][j] , L[i][j-1])
    return L

def find_all_lcs(L, x, y, m, n):
    """
    input
    L: 2d list,
    x, y: string,
    m, n: integer (length of x, y)

    output:
    list of all possible longest common subsequences
    """
```

```python
39      """
40      #if one of them is empty
41      if m == 0 or n == 0:
42          return ['']
43      #if the elements are matching add it to the result
44      elif x[m-1] == y[n-1]:
45          result = [Z + x[m-1] for Z in find_all_lcs(L, x, y, m-1, n-1)]
46          return result
47      else:
48          result = []
49          #look at the elements up and left in the table
50          if L[m][n-1] >= L[m-1][n]:
51              result += find_all_lcs(L, x, y, m, n-1)
52          if L[m-1][n] >= L[m][n-1]:
53              result += find_all_lcs(L, x, y, m-1, n)
54          return result
55
56
57
58  def longest_common_subsequence(x, y):
59      """
60      input:
61      x, y: string
62
63      output:
64      tuple containing list of all the possible combinations of lcs, the length of lcs
65
66      produces the required format of the output
67      """
68      #set guard clause to prevent the unessesary run of the functnios.
69      if len(x) == 0 or len(y) == 0:
70          return (None, 0)
71
72      L = lcs(x , y)
73      result = find_all_lcs(L, x, y, len(x), len(y))
74      length = len(result[0])
75
76      #set the case for 0 common subsequencess
77      if '' in result:
78          result = None
79      return result, length
80
```

```
80
81  #the code is based on the code from CS110 class
82
```

```
1
2
3   x1, y1 = 'ABCBDAB', 'BDCABA'
4   x2, y2 = 'abc', ''
5   x3, y3 = 'abc', 'a'
6   x4, y4 = 'abc', 'ac'
7   x5, y5 = '', ''
8   x6, y6 = 'abc', 'efg'
9   x7, y7 = 'aaa', 'aaa'
10
11  """
12  I changed the order of the common subsequences in the first assertion statement
13  according to the appearance of letters in the second string
14  """
15  assert longest_common_subsequence(x1, y1) == (['BDAB', 'BCAB', 'BCBA'], 4)
16  assert longest_common_subsequence(x2, y2) == (None, 0)
17  assert longest_common_subsequence(x3, y3) == (['a'], 1)
18  assert longest_common_subsequence(x4, y4) == (['ac'], 2)
19  assert longest_common_subsequence(x5, y5) == (None, 0)
20  assert longest_common_subsequence(x6, y6) == (None, 0)
21  assert longest_common_subsequence(x7, y7) == (['aaa'], 3)
22
23
```

```
1
2   set_strings = [('a', 'ACAGCAAGCCATTCCTTAGAGAACGAAATTACGGCGACCGTCAGGGGCATAGCTCCGAGGCATACACTGACGTGTTGGGA'),
3   ('b', 'CAGCAGCACAGTCCTGAGATAGCGAAATGAACCGCGACCGTCAGGGGCCTTCGGCCTCCCGAGGTATTACGGCTGACAGTGGTTGGGA'),
4   ('c', 'GCACCAAGCCTTCACTAGAGGAACCCAGAACTAAGGGTCAGACCGTGTCGGGCCATGCATGCCGAGGTATTACGGACTAGTACGA'),
5   ('d', 'CAGCAGCACAGTCCTGAGATGACGAAATTAACCGGCGACCGTCAGGGGCCATAGCCTCCGGGGCATACGACTGACGTGGTGGGA'),
6   ('e', 'GCAGCAAGCCTTCACTTAGAGAACCGAAACTAGGGCAGACGGTCGGGCATGCATGCCGAGGATATACGGACTAGTACGGGA'),
7   ('f', 'GCTAGCACAAGCCTTCACTTAGATGAACCTAAACTAGGCAGACGGGTCGGCAAGCATAGTCCGAGGATATACAGCGACTAGTACGGG'),
8   ('g', 'CAGCAGACAGCTCCTTGTGATGAACGAAATCAACCGGCGACCATGCAGGCGGCCATAGCCTCCGGTGGCATCACGACTGAGCCGTGGTGGA')]
9
10
```

```python
11  def combinations(set_strings):
12      """
13      input:
14      set_strings: list of tuples with 2 strings
15
16      output:
17      string containing:
18      max_value, min_value integers,
19      min_key[0],min_key[1] max_key[0], max_key[1] strings
20      """
21
22      used = []
23      result = {}
24      for i in set_strings:
25          #do not compare strings twice
26          used.append(i)
27          for j in set_strings:
28              #do not compare 2 identical strings
29              if i == j:
30                  continue
31              if j not in used:
32                  length = len(longest_common_subsequence(i[1], j[1])[0])
33                  result[(i[0],j[0])] = length
34      max_value = max(result.values())
35      max_key = max(result, key = result.get)
36      min_value = min(result.values())
37      min_key = min(result, key = result.get)
38      output = "The maximum number of combinations of lcs is {} between {} and {}; the minimum number of combinations of l
39      return output
40
41
42
43
```

```python
1  combinations(set_strings)
```

```python
1  def matrix(set_strings,percentages = False):
2      """
3      input:
4      set_strings: list of tuples with 2 strings
5      percentages bool automatically set to "False"
6
7      output:
8      numpy matrix of the longest common sequences between every 2 strigs from the set_strings
9
10      program traverses throught every combination of strings and records the result for each row
11      there are 2 options for the output: lcs and lcs in percentages.
12      """
13      rows = []
14      maxi = []
15
16      for i in set_strings:
17          #row is updated after each string was compared to all the other in the inner loop
18          row = []
19          for j in set_strings:
20              row.append(lcs(i[1],j[1])[-1][-1])
21              #set the max value when the stright is compared to itself (100% match)
22              if j[0]==i[0]:
23                  maxi.append(lcs(i[1],j[1])[-1][-1])
24          # as the row is filled, add it to the main list
25          rows.append(row)
26      if percentages:
27          for row in rows:
28              for i in range(len(row)):
29                  #calculate the matrix using percentages
30                  row[i] = int(row[i]/maxi[i]*100)
31      return np.matrix(rows)
32
33
34  len_lcs_matrix = matrix(set_strings,0)
35
36  print(len_lcs_matrix)
```

```
1
2  #matrix for the local strategy
3  def matrix2(set_strings):
4      """
5      input:
6      set_strings: list of tuples with 2 strings
7
8      output:
9      numpy matrix of the minimum edit distance between every 2 strigs from the list
10
11     program traverses throught every combination of strings and records the result for each row
12     """
13     rows = []
14     for i in set_strings:
15         #row is updated after each string was compared to all the other in the inner loop
16         row = []
17         for j in set_strings:
18             if i == j:
19                 # to avoid bias set the comparison with itself to infinity
20                 row.append(float('inf'))
21             else:
22                 row.append(lev(i[1],j[1]))
23         rows.append(row)
24     return rows
25
26
27  matrix2(set_strings)
```

```
: [[inf, 22, 30, 14, 20, 28, 24],
  [22, inf, 36, 14, 35, 40, 23],
  [30, 36, inf, 34, 14, 23, 42],
  [14, 14, 34, inf, 31, 37, 14],
  [20, 35, 14, 31, inf, 14, 39],
  [28, 40, 23, 37, 14, inf, 40],
  [24, 23, 42, 14, 39, 40, inf]]
```

```
1  import math
2  #local greedy strategy
3
4  def find_min(row):
5      """
6      input
7      row: list with edit distace values(integers)
8
9      output
10     minimum: list
11     indexes: list
12
13     function finds minimum values in the row and records its indexes
14     """
15     #set starting values
16     minimum = [row[0], row[1]]
17     indexes = [0,1]
18
19     #skip the first 2 elements since they are starting values
20     for i in range(2,len(row)):
21         for j in range(len(minimum)):
22             #check if smaller
23             if row[i] < minimum[j]:
24                 e = minimum[j]
25                 ind = indexes[j]
26                 #remember to compare with the second element
27                 minimum[j] = row[i]
28                 indexes[j] = i
29
30                 if (j+1) < len(minimum) and e < minimum[j+1]:
31                     minimum[j+1] = e
32                     indexes[j+1] = ind
33                 break
34     return minimum, indexes
35
36
37
```

```python
38  def matrix2(set_strings):
39      """
40      input:
41      set_strings: list of tuples with 2 strings
42
43      output
44      rows: numpy matrix of the minimum edit distance
45
46      function traverses throught every combination of strings and records the result for each row
47      """
48      rows = []
49      for i in set_strings:
50          row = []
51          for j in set_strings:
52              #if the element is compared with itself use infinity to avoid bias
53              if i == j:
54                  row.append(float('inf'))
55              else:
56                  row.append(lev(i[1],j[1]))
57          rows.append(row)
58      return np.matrix(rows),rows
59
60  def clean_table(rows,indexes):
61      """
62      input
63      rows: numpy matrix
64      indexes: list with indexes of the colums to be cleaned.
65
66      output
67      numpy matrix
68      """
69      for row in rows:
70          for i in indexes:
71              row[i] = float('inf')
72
73
74  def find_children(rows,parent,greedy_tree):
75      """
76      input
77      rows: numpy matrix
78      parent: list
79      greedy tree: list
```

```python
74  def find_children(rows,parent,greedy_tree):
75      """
76      input
77      rows: numpy matrix
78      parent: list
79      greedy_tree: list
80
81      output
82      children: list
83      indexes: list
84
85      function find the children of a prenet based on its proximity
86      """
87      row = rows[parent]
88      #recors the value of the child and its location
89      children = find_min(row)[0]
90      indexes = find_min(row)[1]
91
92      #termination condition
93      if children.count(float('inf')) == 2:
94          return
95
96      #remove spare nodes from the tree
97      while children.count(float('inf')) > 0:
98          children.remove(float('inf'))
99
100     #add children to the tree
101     greedy_tree.append(children)
102     clean_table(rows,indexes)
103
104     for i in indexes:
105         find_children(rows,i,greedy_tree)
106
107     return children,indexes
108
```

```python
109 def build_tree(parent, set_strings):
110     """
111     input
112     parent: integer, parent node of the tree
113     set_strings: list of tuples with 2 strings
114
115     output
116     greedy_tree: list representing a tree
117
118     function builds the tree suing the defined above functinos
119     """
120     rows = matrix2(set_strings)[1]
121     clean_table(rows, [parent])
122     greedy_tree=[parent]
123     find_children(rows,parent,greedy_tree)
124     return greedy_tree
125
126
127 def find_sum(greedy_tree):
128     """
129     input
130     greedy_tree: list representing the tree
131
132     output
133     summ: integer
134
135     function calculated the sum of all the minmum edit distace of the tree
136     """
137     summ = 0
138     for i in greedy_tree[1:]:
139         summ+= sum(i)
140     node = greedy_tree[0]
141     return summ
142
```

```
143  def greedy_approach(set_strings, print_output = False):
144      """
145      input
146      set_strings: list of tuples with 2 strings
147      print_output: bollean automatically set to Flase
148
149      output
150      result: list
151
152      function finds the best tree from all the possible trees
153      """
154      max_sum = float('inf')
155      result = []
156      for i in range(len(set_strings)):
157          #build trees staring from each node as a parent node
158          tree = build_tree(i, set_strings)
159          if print_output:
160              print(tree)
161          #find the best tree based on the sum
162          if find_sum(tree) < max_sum:
163              max_sum = find_sum(tree)
164              result = [tree]
165          elif find_sum(tree) == max_sum:
166              result.append(tree)
167
168      if len(result) > 1 and print_output:
169          print('There is more than 1 solution')
170      return result
171
172  greedy_approach(set_strings, True)
173
```

```
1   def matrix_gloabl(set_strings):
2       """
3       input
4       set_strings: list of tuples with 2 strings
5
6       output
7       rows: numpy matrix
8
9       function traverses throught every combination of strings and records the result for each row
10      """
11      rows = []
12      set_strings = set_strings
13      for i in set_strings:
14          #row is updated after each string was compared to all the other in the inner loop
15          row = []
16          for j in set_strings:
17              row.append(lev(i[1],j[1]))
18          rows.append(row)
19      return rows
20
21
22  matrix_gloabl(set_strings)
```

```
[[0, 22, 30, 14, 20, 28, 24],
 [22, 0, 36, 14, 35, 40, 23],
 [30, 36, 0, 34, 14, 23, 42],
 [14, 14, 34, 0, 31, 37, 14],
 [20, 35, 14, 31, 0, 14, 39],
 [28, 40, 23, 37, 14, 0, 40],
 [24, 23, 42, 14, 39, 40, 0]]
```

```python
1   #global stategy
2
3   def find_cousins(rows):
4       """
5       input
6       rows: numpy matrix
7
8
9       output
10      cousins: list
11
12      function traverses the rows and finds the higest edit distance
13      """
14      maxi = []
15      for row in rows:
16          maxi.append(max(row))
17      highest = max(maxi)
18      #index of the element based on the colum
19      index1 = maxi.index(highest)
20      #index of the element based on the row
21      index2 = rows[index1].index(highest)
22      cousins = [index1,index2]
23      return cousins
24
```

```python
26  def find_bothers(rows):
27      """
28      input
29      rows: numpy matrix
30
31      output
32      result: list of lists with integers (indexes of brothers)
33
34      function copies rows where the "cousins" are located and finds third closes element (first is the element itself,
35      second is a parent element and thirs would be a "brother" element)
36      """
37      i1, i2 = find_cousins(rows)
38      r1 = sorted(rows[i1].copy())
39      r2 = sorted(rows[i2].copy())
40      bro1,bro2 = r1[2],r2[2]
41      b1,b2 = [rows[i1].index(bro1),rows[i2].index(bro2)]
42      result = [[i1,b1],[i2,b2]]
43      return result
44
45  def find_parents(rows):
46      """
47      input
48      rows: numpy matrix
49
50      output
51      parents: list of an integer and a list
52
53      function finds the highest edit distance in the merged colums - uncles of the merged elements
54      """
55      merged1 = rows[-1]
56      merged2 = rows[-2]
57      #clean the rows
58      #no need to clean colums since we consider only 2 merged columns
59      for row in rows:
60          for i in brothers:
61
62              row[i[0]] = 0
63              row[i[1]] = 0
64      parents = [merged1.index(max(merged1)),merged2.index(max(merged2))]
65      return parents
66
```

```python
68  def merge_colums(rows, indexes):
69      """
70      input
71      rows: numpy matrix
72
73      output
74      rows: numpy matrix
75
76
77      function merges 2 colums of the "brother" elements in 1 based on their averages
78      """
79      merged1,merged2 = [],[]
80      for row in rows:
81          for i in indexes:
82
83              #switch 0 to infinity to avoid confound
84              if row[i[0]] == 0:
85                  row[i[0]] = float('inf')
86              elif row[i[1]] == 0:
87                  row[i[1]] = float('inf')
88
89              #sort merged values in the coresponsing columns
90              if indexes.index(i) == 0:
91                  merged1.append((row[i[0]] + row[i[1]]) / 2)
92              else:
93                  merged2.append((row[i[0]] + row[i[1]]) / 2)
94
95      #add new colums to the matrix
96      rows.append(merged1)
97      rows.append(merged2)
98      return rows
99
100
```

```python
101  def build_global_tree(rows,brothers,parents):
102      """
103      input
104      rows: numpy matrix
105      brothers: list of lists with integers (indexes of brothers)
106      parents: list of an integer and a list
107
108
109      output
110      tree: list representing the structure of the tree
111
112      function finds grandparent - element that was not yet used and builds tree using the output of other functions.
113      """
114      n = len(rows[0])
115      grandparent = [i for i in range(n)]
116      #cross- assignment of parent because we were looking at the higest distace
117      for i in parents:
118          grandparent.remove(i)
119      for bro in brothers:
120          for i in bro:
121              grandparent.remove(i)
122      tree = [grandparent[0],[parents[1],brothers[0]],[parents[0],brothers[1]]]
123      return tree
124
125
126
127
128  rows = matrix_gloabl(set_strings)
129  brothers = find_bothers(rows)
130  rows = merge_colums(rows, brothers)
131  parents = find_parents(rows)
132  global_tree = build_global_tree(rows,brothers,parents)
133  print(global_tree)
```

[0, [3, [2, 5]], [4, [6, 1]]]

```python
#pribabilities estimation

def probabilities(str1, str2, len1, len2):
    """
    input
    str1, str2: strings
    len1, len2: integers (lengths of str1 and str2 respectively)

    output
    result: list of floats

    function the number of changed characters using lcs as a basis and calucates the percentage
    """
    l = lcs(str1 , str2)[-1][-1]
    #find the remaining of the string that does not include common characters
    rest1 = len1 - l
    rest2 = len2 - l
    insertion, deletion, mutation = 0,0,0
    # if the child is bigger there was no deletion
    if rest1 < rest2:
        insertion = rest2 - rest1
        mutation = rest1
    #if smaller - no insertion
    elif rest2 < rest1:
        deletion = rest1 - rest2
        mutation = rest2
    #in they are equal the characters mutated
    elif rest2 == rest1:
        mutation = rest1
    #since str2 is a child string
    result = mutation/len2, deletion/len2, insertion/len2
    return result
```

```python
def links(tree):
    """
    input
    tree: list of lists

    output
    links: list of lists. Each inner list contains 2 elements in the following form: [parent, child]

    function deconstructs the tree into parenth-children links
    """
    links = []
    for i in tree[1:]:
        #add grandparent - parent link
        links.append([tree[0],i[0]])
    #deconstruct each subtree and add parent - child links
    t1 = tree[1]
    for i in t1[1]:
        links.append([t1[0], i])
    t2 = tree[2]
    for i in t2[1]:
        links.append([t2[0], i])
    return links


links = links(global_tree)
print(links)
```

[[0, 3], [0, 4], [3, 2], [3, 5], [4, 6], [4, 1]]

```python
from statistics import mean


def probabilities_calculation(set_strings, links):
    """
    input
    set_strings: list of tuples with 2 strings
    links: list of lists


    output:
    result: list of floats (probabilities)

    function finction calculated the probabilities of change in all the strings in the tree
    """
    mutation = []
    deletion = []
    insertion = []
    strings = []
    #create list of strings only
    for i in set_strings:
        strings.append(i[1])
    for i in links:
        #find any 2 strings
        str1 = strings[i[0]]
        str2 = strings[i[1]]
        #find probabilities
        p = probabilities(str1, str2, len(str1), len(str2))
        mutation.append(p[0])
        deletion.append(p[1])
        insertion.append(p[2])
    result = [mean(mutation), mean(deletion), mean(insertion)]
    return result



probabilities_calculation(set_strings, links)
```

[0.19491766580734352, 0.0, 0.049274625928333385]

```python
1   import random
2   import time
3   import matplotlib.pyplot as plt
4   #tests for the greedy approach
5
6   #generate increasing length of the string
7   def increase_length_local():
8       """
9       input
10      None
11
12      output
13      local_time: list of floats
14      len_set: list of integers
15
16      function runs the test 15 times increasing the number of genes each time. the length of each gene is fixed at a rand
17      """
18      letters = 'AGCT'
19      local_time = []
20      str_set1 = []
21      len_set = []
22
23      M = random.randint(5,10) #fixed length of a gene
24      N = random.randint(3,5) #starting number of genes
25
26
27      for i in range(15): #number of tests
28          N = N+i
29          len_set.append(N)
30
31          for j in range(N):
32              str_set1.append(('n', ''.join([random.choice(letters) for i in range(M)])))
33
34          #start time
35          st = time.process_time()
36
37          # main program
38          greedy_approach(str_set1,0)
39
```
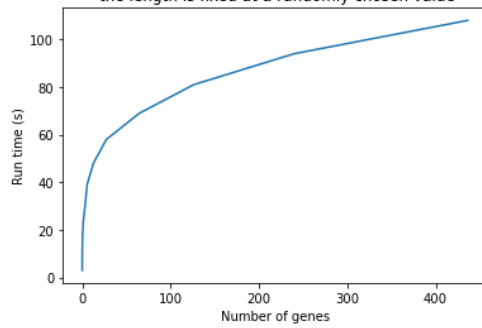
```python
31          for j in range(N):
32              str_set1.append(('n', ''.join([random.choice(letters) for i in range(M)])))
33
34          #start time
35          st = time.process_time()
36
37          # main program
38          greedy_approach(str_set1,0)
39
40          #end time
41          et = time.process_time()
42
43          local_time.append(et-st)
44
45
46      return local_time, len_set
47
48
49  lst = increase_length_local()
50
51  def build_graph(lst):
52      plt.plot(lst[0], lst[1])
53      plt.title('Change of the run time as the number of genes increases and \n the length is fixed at a randomly chosen va
54      plt.xlabel('Number of genes')
55      plt.ylabel('Run time (s)')
56      plt.show()
57
58
59      return plt.show()
60
61  build_graph(lst)
62
```

```
60
61  build_graph(lst)
62
```

Change of the run time as the number of genes increases and
the length is fixed at a randomly chosen value

**References**

*Algorithm Implementation/Strings/Longest common subsequence - Wikibooks, open books for an*

*open                               world*.                               (n.d.).

https://en.wikibooks.org/wiki/Algorithm_Implementation/Strings/Longest_common_subs

equence

Back To Back SWE. (2019, January 11). *Edit Distance Between 2 Strings - The Levenshtein*

*Distance      ("Edit      Distance"      on      LeetCode)*.      YouTube.

https://www.youtube.com/watch?v=MiqoA-yF-0M

*Levenshtein distance with substitution, deletion and insertion count*. (2020, May 13). Stack

Overflow.

https://stackoverflow.com/questions/61784300/levenshtein-distance-with-substitution-del

etion-and-insertion-count

GeeksforGeeks.      (2022,      December      10).      *Edit      Distance      |      DP-5*.

https://www.geeksforgeeks.org/edit-distance-dp-5/

GeeksforGeeks.      (2022a,      May      24).      *sort()      in      Python*.

https://www.geeksforgeeks.org/sort-in-python/

GeeksforGeeks. (2022b, December 14). *Print all longest common sub-sequences in*

*lexicographical                               order*.

https://www.geeksforgeeks.org/print-longest-common-sub-sequences-lexicographical-ord

er/

GeeksforGeeks.      (2022b,      December      10).      *Edit      Distance      |      DP-5*.

https://www.geeksforgeeks.org/edit-distance-dp-5/