



String Matching Algorithms

Algorithms Project

14012402-4

Ruba Balubaid 442017269
Yara Al simlan 442011215
Noor Alhashmi 442007783
Shaden Anagreh 442017295
Hadel Alnasiri 442004645

Under Supervision:

Dr. Manal Alharbi

- **Project objectives:**

This project aims to present 3 different data structures to implement 3 different string matching algorithms (KMP, Naïve, and Longest common subsequence) on the same data, for analyzing and comparing the overall performances of each algorithm.

- **Programming Language :**

Python.

- **Project specifications:**

For this project, we implement the following 3 String Matching Algorithms:

- **Naive string matching algorithm:**

- algorithm pseudocode:***

- Input: two strings txt, pat
 - Output: the index in which a pattern (pat) match has been found in the text (txt).

Naïve_String_Matching(txt, pat)

- 1) $n = \text{length}(\text{txt})$
- 2) $m = \text{length}(\text{pat})$
- 3) for $i = 0$ to $(n-m)$
- 4) if $\text{pat}[1\dots m] = \text{txt}[i+1\dots i+m]$;
- 5) print "Match found at " i

- Longest Common Subsequence string matching algorithm:

algorithm pseudocode:

LCS-LENGTH (X,Y)

- 1) $m = X.length$
- 2) $n = Y.length$
- 3) let $b[1 \dots m, 1 \dots n]$ and $c[0..m, 0..n]$ be new tables
- 4) for $I = 1$ to m
- 5) $c[I, 0] = 0$
- 6) for $j = 0$ to n
- 7) $c[I, 0] = 0$
- 8) for $I = 1$ to m
- 9) for $j = 1$ to n
- 10) if $X == Y$
- 11) $c[I, j] = c[i-1, j-1] + 1$
- 12) $b[I, j] = c[i-1, j-1] + 1$
- 13) elseif $c[i-1, j] > c[I, j-1]$
- 14) $c[I, j] = c[i-1, j]$
- 15) $b[I, j] = c[i-1, j]$
- 16) else $c[I, j] = c[I, j-1]$
- 17) $b[I, j] = c[I, j-1]$
- 18) return c and b

- KMP string matching algorithm:

algorithm pseudocode:

COMPUTE- PREFIX- FUNCTION (P)

1. $m \leftarrow \text{length}[P]$ // 'p' pattern to be matched
2. $\Pi[1] \leftarrow 0$
3. $k \leftarrow 0$
4. for $q \leftarrow 2$ to m
5. do while $k > 0$ and $P[k + 1] \neq P[q]$
6. do $k \leftarrow \Pi[k]$
7. If $P[k + 1] = P[q]$
8. then $k \leftarrow k + 1$
9. $\Pi[q] \leftarrow k$
10. Return Π

KMP-MATCHER (T, P)

1. $n \leftarrow \text{length}[T]$
2. $m \leftarrow \text{length}[P]$
3. $\Pi \leftarrow \text{COMPUTE-PREFIX-FUNCTION}(P)$
4. $q \leftarrow 0$ // numbers of characters matched
5. for $i \leftarrow 1$ to n // scan S from left to right
6. do while $q > 0$ and $P[q + 1] \neq T[i]$
7. do $q \leftarrow \Pi[q]$ // next character does not match
8. If $P[q + 1] = T[i]$
9. then $q \leftarrow q + 1$ // next character matches
10. If $q = m$ // is all of p matched?
11. then print "Pattern occurs with shift" $i - m$
12. $q \leftarrow \Pi[q]$

- **Project specifications(cont.):**

We present the following 3 different data structures to implement each of the algorithms:

3. Hash table: We used this data structure because they allow for fast lookups and can store large amounts of data. They also provide constant time access to elements, which is important for string matching algorithms that need to quickly find matches.

2. List: Lists are a good choice for string matching algorithms because they are easy to traverse and can store large amounts of data. Additionally, lists can be sorted, which is useful for finding matches quickly.

3. Linked List: Linked lists are a good choice for string matching algorithms because they allow for fast insertion and deletion of elements, which is important when dealing with large datasets. Additionally, linked lists can be traversed quickly, which is important when searching for matches.

- **README :**

- *Required Libraries:*

Time, pandas

- *Data's file:*

We have three file contains data. We use "both_covid_data" file.

- **README(cont.) :**

- ***Project Files:***

The main file (StringMatchingAlgo), which contains as well as the data file, and the Project main report, the following three files:

1. Naïve_Algo:

This file contains 3 files, each present a different data structure to implement the Naive String Matching algorithm (Same data used):

1. Naïve Algo – List: (.py)

In this file, data will be read from the given file and stored in a data frame. A method named naive_string_match() will search for a pattern match in the data stored in the data frame, and return the indices for each pattern match has been found in a List named matches.

2. Naïve Algo – Linked List: (.py)

This file contains 2 classes (Node, Linked List), to create a linked list that contains nodes that hold data elements from the data frame, and a pointer to the next element. Linked List class contains 3 methods: append(), to add new nodes, printList(), to print the linked list content, and method extract(), which aims to extract the data from each node to be send to the method named naive_string_match(), which aims to search -in each data element - for a pattern matching and return the index in which it has been found.

3. Naïve Algo – Hash Table: (.py)

In this file the data will be stored from the data frame in a Hash Table, and it contains 2 methods: createHashTable() which create a hash table of all the substrings of the given pattern, and searchPattern() method which searches for all occurrences of the pattern (Naïve algorithm) in the given data. Method accepts each data element stored in the hash table and returns the index in which a pattern match was found.

- **README(cont.) :**

2. KMP algo:

This file contains 3 files, each present a different data structure to implement the KMP String Matching algorithm (Same data used) :

1. *KMPList: (.py)*

In this file there are some methods:

- 1) readlist: for read a data but we do not use
- 2,3) KMP and computeLPSArray: algorithm's method
- 4) read_csv: we use it for a read data

2. *KMPLinkedList: (.py)*

In this file there are a Node class and KMPSearch method for implement KMP algorithm. For read data we use a read_csv method which is convert a csv file to data frame and we use a pandas library here.

3. *KMPHashTable: (.py)*

This file contains same methods in KMPList file, but here we use a hash table data structure.

3. Longest algo:

This file contains 3 files, each present a different data structure to implement the longest String Matching algorithm (Same data used):

1. *Longest_List: (.py)*

in the code of the longest list, first, I created a data frame to read the data from a file, then I created a code to join all the elements to the row, then I defined two variables to read the lengths and defined an array to store the dp values, then I applied the following steps build $L[m+1][n+1]$ in bottom up fashion. Thus, the letters will be stored one by one, in the event that the next letter in x and y is the same, meaning the current letter is part of the longest, but if it is not, it will search for the largest number of the two, then it will move to the direction of the largest value.

- **README(cont.) :**

3. Longest algo (cont.):

2. Longest_Linkedlist: (.py)

In the longest LinkedList code, also, I created a data frame to read data from a file, then I created code to join all the elements to the row

Next, create a node class, define the node, and enter the data

In linked list class contents a Nood object


I also added a function to initialize the header, so that it prints the contents of the list starting from the header, after that I created a method to return the length of the list, and I also created a method to return the longest string between the two strings using LinkedList

3. Longest_Hashtable: (.py)


In the longest HashTable code, also, I created a data frame to read data from a file, then I created code to join all the elements to the row.

First, I created a table to store the results of the sub-problems, then it would store the lcs string from bottom to top, and it would find the length of the LCs. Then, create a character array to store the lcs string

If the current letter in X and Y is the same, then the current letter is part of the LCS, and if it is not, search for the larger number of the two and go to the direction of the larger value.



Python Codes

- Screenshots was taken from each python code has been written shown below.
- 

Longest List code:

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Sat Feb  4 10:13:00 2023
4
5  @author: shade
6  """
7  #importing pandas library
8  import pandas as pd
9  import time
10
11
12  df = pd.read_csv("both_covid_data.csv")
13  df = df.astype(str)
14
15
16  df["joined_row"] = df.apply(lambda row: ''.join(row), axis=1)
17
18
19
20  start_time = time.perf_counter()
21
22  def lcs(X, Y):
23
24      m = len(X)
25      n = len(Y)
26
27
28      L = [[None]*(n+1) for i in range(m+1)]
29
30
31      for i in range(m+1):
32          for j in range(n+1):
33              if i == 0 or j == 0 :
34                  L[i][j] = 0
35              elif X[i-1] == Y[j-1]:
36                  L[i][j] = L[i-1][j-1]+1
37              else:
38                  L[i][j] = max(L[i-1][j], L[i][j-1])
39
40
41
42      index = L[m][n]
43
```

```
41      index = L[m][n]
42
43
44
45      lcs="" * (index+1)
46
47      lcs[index] = ""
48
49
50      i=m; j=n;
51
52      while i > 0 and j > 0:
53
54          if X[i - 1] == Y[j - 1]:
55              lcs[index - 1]=X[i - 1];
56              i -= 1; j -= 1; index -= 1;
57
58          elif L [i - 1][j] > L [i][j - 1]:
59              i -= 1;
60
61          else:
62              j -= 1;
63
64      print ("LCS of " + X + " and " + Y + " is " + ''.join(lcs))
65
66  end_time = time.perf_counter()
67
68  print("Time taken for implementing (Longest String Matching algorithm) using Array:",end_time-start_time)
69
70
71
72
73
74
75
76
```

Longest HashTable code:

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Sat Feb  4 20:33:26 2023
4
5  @author: shade
6  """
7
8  import pandas as pd
9  import time
10
11  df = pd.read_csv("both_covid_data.csv")
12  df = df.astype(str)
13
14
15  df["joined_row"] = df.apply(lambda row: ''.join(row), axis=1)
16
17
18  start_time = time.perf_counter()
19
20
21  def longestCommonSubsequence(str1, str2):
22      m, n = len(str1), len(str2)
23      dp = [[0] * (n + 1) for _ in range(m + 1)]
24
25      for i in range(1, m + 1):
26          for j in range(1, n + 1):
27              if str1[i - 1] == str2[j - 1]:
28                  dp[i][j] = dp[i - 1][j - 1] + 1
29              else:
30                  dp[i][j] = max(dp[i - 1][j], dp[i][j - 1])
31
```

```
32      index = dp[m][n]
33      lcs = [""] * index
34
35      i, j = m, n
36      while i > 0 and j > 0:
37          if str1[i - 1] == str2[j - 1]:
38              lcs[index - 1] = str1[i - 1]
39              i -= 1
40              j -= 1
41              index -= 1
42          elif dp[i - 1][j] > dp[i][j - 1]:
43              i -= 1
44          else:
45              j -= 1
46
47      return ''.join(lcs)
48
49  end_time = time.perf_counter()
50
51  print("Time taken for implementing (Longest String Matching algorithm) using Hash Table:", end_time - start_time)
```

Longest linked list code:

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Sun Feb  5 12:34:29 2023
4
5  @author: shade
6  """
7
8  # -*- coding: utf-8 -*-
9  """
10 Created on Wed Feb  1 06:56:04 2023
11
12 @author: shade
13 """
14 import pandas as pd
15 import time
16
17 df = pd.read_csv("both_covid_data.csv")
18 df = df.astype(str)
19
20
21 df["joined_row"] = df.apply(lambda row: ''.join(row), axis=1)
22
23
24
25
26 start_time = time.perf_counter()
27
28 class Node:
29
30
31     def __init__(self, data):
32         self.data = data
33         self.next = None
34
35
36 class LinkedList:
37
38     def __init__(self):
39         self.head = None
40
41     def printList(self):
42
43         temp = self.head
```

```
44
45         while (temp):
46             print (temp.data)
47             temp = temp.next
48
49
50
51     def getlength(self):
52
53         temp = self.head
54
55         count = 0;
56
57         while (temp):
58             count += 1;
59             temp = temp.next;
60
61         return count;
62
63
64     def lcs(self, X, Y, m, n):
65
66         if m == 0 or n == 0:
67             return 0;
68
69         elif X[m-1] == Y[n-1]:
70             node = Node(X[m-1])
71             if self.head is None:
72                 self.head = node
73             else:
74                 current_node = self.head
75                 while current_node.next is not None:
```

```
82         current_node = self.head
83
84         while current_node.next is not None:
85             current_node=current_node.next
86             current_node=current_node.next
87             current_node.next=node
88             return 1 + self.lcs(X, Y, m-1, n-1);
89         else:
90             return max(self.lcs(X, Y, m, n-1), self.lcs(X, Y, m-1, n));
91
92
93
94
95
96 end_time = time.perf_counter()
97
98 print("Time taken for implementing (Longest String Matching algorithm) using Linked List:",end_time-start_time)
```

KMP linked list code:

```
53 for i in range(M):
54
55     # create a new node for every character of pattern
56     curr = Node(pat[i])
57
58     # append the new node at the end of linked list
59     prev.next = curr
60
61     # update previous pointer to current node
62     prev = curr
63
64 # create linked list of text characters
65 head2 = Node(None) # dummy node to mark beginning
66 prev = head2
67
68 for i in range(N):
69
70     # create a new node for every character of text
71     curr = Node(txt[i])
72
73     # append the new node at the end of linked list
74     prev.next = curr
75
76     # update previous pointer to current node
77     prev = curr
78
79 i, j, count = 0, 0, 0
80
81 while i < N:
82     # if characters match then move ahead in both lists
83     if (head2.data == head1.data):
84         count += 1
85         i += 1
86         j += 1
87         head2 = head2.next
88         head1 = head1.next
89     else:
90         if (j != 0):
91             j = 0
92         else:
93             i += 1
94             head2 = head2.next
95         if (count == M):
96             return True
97         return False
98 count += 1
```

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Wed Jan 11 07:24:49 2023
4
5 @author: ruba balubaid
6 """
7
8 """
9 Definition for variables and methods in code:
10
11 - myFile: is a file to read
12 - text: array to store information from a file
13 - readlist(): to read a txt file
14 - KMP(text, pattern): to take array and pattern
15 - n,m: to calculate length
16 - lps: will hold the longest prefix suffix values for pattern
17 """
18
19 # linked list data structure code
20 import time
21 import pandas as pd;
22
23 start_time = time.perf_counter()
24
25 """
26 def readlist():
27     # opening the file in read mode
28     myFile = open("both_covid_data.txt", "r")
29     text = []
30     for i in myFile:
31         text.append(str(i))
32     myFile.close()
33     return text
34
35 text = readlist()
36 """
37 # Python program for KMP Algorithm using Linked List
38
39 # A linked list node
40 class Node:
41     def __init__(self, data):
42         self.data = data
43         self.next = None
44
45 def KMPSearch(pat, txt):
46     M = len(pat)
47     N = len(txt)
48
49     # create linked list of pattern characters
50     head1 = Node(None) # dummy node to mark beginning
51     prev = head1
52
53     for i in range(M):
54
```

```
100 def read_csv():
101     # Read the CSV file into a DataFrame
102     df = pd.read_csv("both_covid_data.csv")
103     df = df.astype(str)
104     # print(len(df))
105     # Join all the elements in each row with no separator
106     df["joined_row"] = df.apply(lambda row: ''.join(row), axis=1)
107
108     # Define the target string
109     target = "0010"
110
111     # Initialize a variable to store the result
112     result = -1
113
114     # Iterate through the joined_row column of the DataFrame
115     for i, s in enumerate(df["joined_row"]):
116
117         col = KMPSearch(s, target)
118         if col != -1:
119             # If a match is found, store the index in the result variable
120             result = i
121             break
122
123     # Print the result
124     if result != -1:
125         print("Match found at index", result, "at col", col)
126     else:
127         print("No match found.")
128
129 read_csv()
130
131
132 end_time = time.perf_counter()
133
134 print("Time taken in KMP Linked List:", end_time - start_time)
135 print()
```

KMP list code:

```
48 def KMP(pat, txt):
49     M = len(pat)
50     N = len(txt)
51
52     # create lps[] that will hold the longest prefix suffix
53     # values for pattern
54     lps = [0]*M
55     j = 0 # index for pat[]
56
57     # Preprocess the pattern (calculate lps[] array)
58     computeLPSArray(pat, M, lps)
59
60     i = 0 # index for txt[]
61     while (N - i) >= (M - j):
62         if pat[j] == txt[i]:
63             i += 1
64             j += 1
65
66         if j == M:
67             print("Found pattern at index " + str(i-j))
68             j = lps[j-1]
69
70         # mismatch after j matches
71         elif i < N and pat[j] != txt[i]:
72             # Do not match lps[0..lps[j-1]] characters,
73             # they will match anyway
74             if j != 0:
75                 j = lps[j-1]
76             else:
77                 i += 1
78
79
80 def computeLPSArray(pat, M, lps):
81     len = 0 # length of the previous longest prefix suffix
82
83     lps[0] = 0 # lps[0] is always 0
84     i = 1
85
86     # the loop calculates lps[i] for i = 1 to M-1
87     while i < M:
88         if pat[i] == pat[len]:
89             len += 1
90             lps[i] = len
```

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Tue Jan 10 21:30:34 2023
4
5 @author: ruba balubaid
6 """
7
8
9 """
10 Definition for variables and methods in code:
11
12 - myFile: is a file to read
13 - text: array to store information from a file
14 - readlist(): to read a txt file
15 - KMP(text, pattern): to take array and pattern
16 - n,m: to calculate length
17 - lps: longest proper prefix will hold the longest prefix suffix values for pattern
18 - read_csv(): to convert csv file to data frame and read it
19
20 """
21
22 # we use array data structure in KMP algorithm
23
24 import time;
25 import pandas as pd;
26
27
28 start_time = time.perf_counter()
29
30 # Code to be timed
31
32
33 """
34 def readlist():
35     #opening the file in read mode
36     myFile = open("both_covid_data.txt","r")
37     text = []
38     for i in myFile:
39         text.append(str(i))
40     myFile.close()
41     return text
42
43 text = readlist()
44
45 """
```

```
91     i += 1
92 else:
93     # This is tricky. Consider the example.
94     # AAACAAAA and i = 7. The idea is similar
95     # to search step.
96     if len != 0:
97         len = lps[len-1]
98
99     # Also, note that we do not increment i here
100 else:
101     lps[i] = 0
102     i += 1
103
104
105 def read_csv():
106     # Read the CSV file into a DataFrame
107     df = pd.read_csv("both_covid_data.csv")
108     df = df.astype(str)
109     #print(len(df))
110     # Join all the elements in each row with no separator
111     df["joined_row"] = df.apply(lambda row: ''.join(row), axis=1)
112
113     # Define the target string
114     target = "0010"
115
116     # Initialize a variable to store the result
117     result = -1
118
119     # Iterate through the joined_row column of the DataFrame
120     for i, s in enumerate(df["joined_row"]):
121
122         col=KMP(target, s)
123         if col!=-1:
124             # If a match is found, store the index in the result variable
125             result = i
126             break
127
128     # Print the result
129     if result != -1:
130         print("Match found at index", result," at col",col)
131     else:
132         print("No match found.")
133
134 read_csv()
135
136
137 end_time = time.perf_counter()
138
139 print("Time taken in KMP List:", end_time - start_time)
140
141
```

KMP Hash table code:

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Sat Jan 28 18:24:23 2023
4
5  @author: ruba balubaid
6  """
7
8
9  import time
10 import pandas as pd;
11
12 #start = time.process_time_ns()
13 start_time = time.perf_counter()
14
15
16 """
17 def readlist():
18     #opening the file in read mode
19     myFile = open("both_covid_data.txt", "r")
20     text = []
21     for i in myFile:
22         text.append(str(i))
23     myFile.close()
24     return text
25
26 text = readlist()
27 """
28
29
30 def KMP(text, pattern):
31     n = len(text)
32     m = len(pattern)
33
34     # create a hash table to store the index of pattern
35     hash_table = [0] * 256
36
37     # fill the hash table
38     for i in range(m):
39         hash_table[ord(pattern[i])] = i + 1
40
41     # initialize lps array and variables
42     lps = [0] * m
43     j, i = 0, 0
```

```
45     while i < n:
46
47         # if characters match, increment both pointers
48         if text[i] == pattern[j]:
49             i += 1
50             j += 1
51
52         # if j is equal to length of pattern, we have found a match.
53         if j == m:
54             print("Pattern found at index " + str(i-j))
55
56             # reset j to lps value
57             j = lps[j-1]
58
59         # mismatch after j matches
60         elif i < n and text[i] != pattern[j]:
61
62             # Do not match lps[0..lps[j-1]] characters, they will match anyway
63             if j != 0:
64                 j = lps[j-1]
65
66             else:
67                 i += 1
68
69         return -1
70
71
72 def read_csv():
73     # Read the CSV file into a DataFrame
74     df = pd.read_csv("both_covid_data.csv")
75     df = df.astype(str)
76     #print(len(df))
77     # Join all the elements in each row with no separator
78     df["joined_row"] = df.apply(lambda row: ''.join(row), axis=1)
79
80     # Define the target string
81     target = "0010"
82
83     # Initialize a variable to store the result
84     result = -1
```

```
86     # Iterate through the joined_row column of the DataFrame
87     for i, s in enumerate(df["joined_row"]):
88
89         col=KMP(s,target)
90         if col!=-1:
91             # If a match is found, store the index in the result variable
92             result = i
93             break
94
95     # Print the result
96     if result != -1:
97         print("Match found at index", result, " at col", col)
98     else:
99         print("No match found.")
100
101 read_csv()
102
103
104
105
106 end_time = time.perf_counter()
107
108 print("Time taken in KMP hash table:", end_time - start_time)
109
110
```


Naive list code:

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Mon Jan 30 07:32:07 2023
4
5  @author: itzha
6  """
7
8  #importing the necessary libraries
9  import time
10 import pandas as pd
11
12 #creating a dataframe that read data from the given file.
13 df = pd.read_csv("both_covid_data.csv")
14 df = df.astype(str)
15 # Join all the elements in each row with no separator.
16 df["joined_row"] = df.apply(lambda row: ''.join(row), axis=1)
17
18
19
20 #Naive string matching algorithm to find patterns in data stored in a frame.
21 start_time = time.perf_counter()
22 def naive_string_match(dataframe, pattern):
23
24     # Initialize the matches list
25     matches = []
26
```

```
27     # Iterate over each row of the dataframe
28     for index, row in dataframe.iterrows():
29
30         # Iterate over each column of the row
31         for col in row:
32
33             # If pattern is found in the column value add it to matches list.
34             if pattern in str(col):
35
36                 matches.append(index)
37
38     return matches
39
40 #Pattern:
41 pattern = '0010'
42
43 matches = naive_string_match(df, pattern)
44 print("Matches found at positions:", matches) #Print the matches indices.
45
46 end_time = time.perf_counter()
47
48 #Calculating and printing the run time needed to search for a pattern match.
49 print("Time taken for implementing (Naive String Matching algorithm) using Array: ", end_time - start_time)
50
```


Naive linked list code:

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Tue Jan 31 22:27:54 2023
4
5  @author: itzha
6  """
7
8  import pandas as pd
9  import time
10
11  #Pattern :
12  pattern = '0010'
13
14  #Class to create nodes that hold the data element, and a pointer to the next node.
15  class Node:
16      def __init__(self, data):
17          self.data = data
18          self.next = None
19
20  # Linked List class contains a node object, to create a linked list.
21  class LinkedList:
22      def __init__(self):
23          self.head = None
24
25      # method to add new node to the linked list
26      def append(self, new_data):
27
28          # create a new node using the given data and assign it to the head of the linked list if it is empty.
29          if self.head is None:
30              self.head = Node(new_data)
31
32          # else, traverse till the last node and insert the new_node there.
33          else:
34
35              last = self.head
36
37              while (last.next):
38
39                  last = last.next
40
41              last.next = Node(new_data)
42
```

```
43      # method to print all nodes of the linked list
44      def printlist(self):
45
46          temp = self.head
47
48          while (temp):
49
50              print (temp.data)
51
52              temp = temp.next
53      #method to traverse data element in each node.
54      def extract(self):
55
56          temp = self.head
57
58          while (temp):
59              #each data element in the linked list will be send as an argument to the naive string match method, to search for a pattern match.
60              naive_string_match(pattern, temp.data)
61
62              #next element..
63              temp = temp.next
64
65      start_time = time.perf_counter()
66
67      # method to apply Naive String Matching algorithm on the data stored in Linked List
68      def naive_string_match(pat, txt):
69          M = len(pat)
70          N = len(txt)
71
72          #A loop to slide pat[] one by one.
73          for i in range(N - M + 1):
74              j = 0
75
76              # For current index i, check for pattern match.
77              while(j < M):
78                  if (txt[i + j] != pat[j]):
79                      break
80                  j += 1
81
82              if (j == M):
83                  print("Pattern found at index ", i) #Print the 1st index in which a pattern match has been found.
```

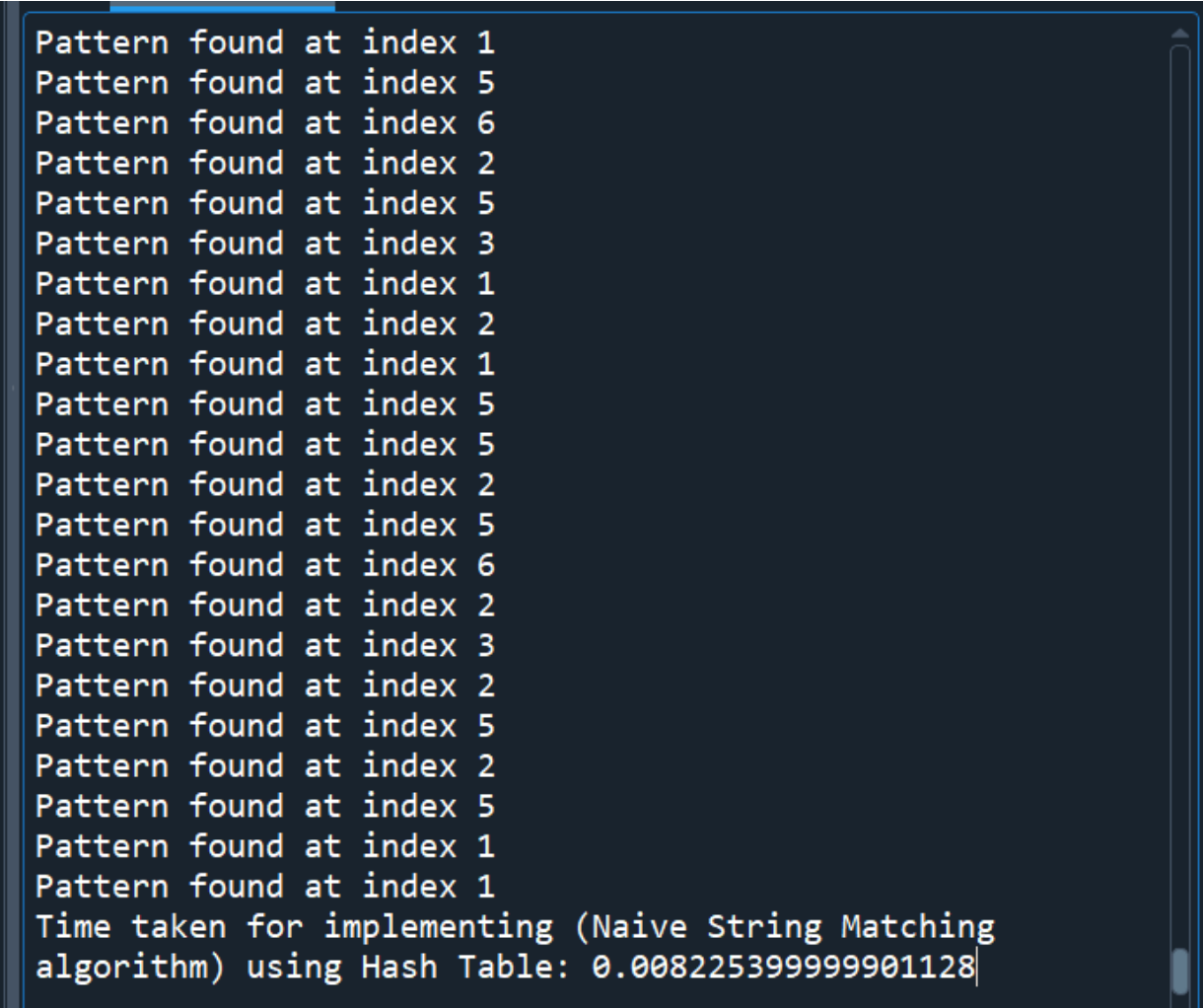
```
88  mylist=LinkedList()
89
90  #creating a dataframe that read data from the given file.
91  df = pd.read_csv("both_covid_data.csv")
92  df = df.astype(str)
93  #Join all the elements in each row with no separator
94  df["joined_row"] = df.apply(lambda row: ''.join(row), axis=1)
95  for i,s in enumerate(df["joined_row"]):
96      #Store the data in the linked list.
97      mylist.append(s)
98
99
100  #search for a pattern match for each linked list element.
101  mylist.extract()
102  end_time = time.perf_counter()
103
104  #Calculating and printing the run time needed to search for a pattern match.
105  print("Time taken for implementing (Naive String Matching algorithm) using Linked List:", end_time - start_time)
106
107
```

Naive hash table code:

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Tue Jan 31 22:45:24 2023
4
5  @author: itzha
6  """
7
8  #importing pandas library
9  import pandas as pd
10 import time
11
12 #creating a dataframe that read data from the given file.
13 df = pd.read_csv("both_covid_data.csv")
14 df = df.astype(str)
15 #Join all the elements in each row with no separator
16 df["joined_row"] = df.apply(lambda row: ''.join(row), axis=1)
17
18
19
20 #creating a hash table, and storing the data from the dataframe in.
21 hash_table = {df.index[i] : df.joined_row[i] for i in range(len(df))}
22
23
24
25 #creatin a hash table (for the pattern.)
26 table = {}
27
28 #method to create a hash table of all the substrings of the pattern.
29 def createHashTable(pattern):
30     m = len(pattern)
31     for i in range(m):
32         table[pattern[i:i+1]] = i
33
34
35 start_time = time.perf_counter()
36
37 #method to search for all occurrences of the pattern in the given text (data in the hash table) using Naive String Matching Algorithm.
38 def searchPattern(text, pattern):
39     n = len(text)
40     m = len(pattern)
```

```
41
42     #Create a hash table for pattern
43     createHashTable(pattern)
44
45     # Traverse through the text (data from the hash table)
46     for i in range(n-m+1):
47
48         # For current index i, check for pattern match
49         j = 0;
50
51         while j < m:
52
53             # If characters don't match, break the loop and move to next index in text.
54             if text[i + j] != pattern[j]:
55                 break;
56
57             j += 1;
58
59             # If *all* characters are matched, print the index at which it is found.
60             if j == m:
61                 print("Pattern found at index", i);
62
63
64
65
66
67 #Pattern:
68 pattern = '0010'
69
70 #A loop to traverse the data elements in the hash table, and search in each for a pattern matching.
71 for key, value in hash_table.items():
72     searchPattern(value, pattern)
73
74
75
76 end_time = time.perf_counter()
77
78 #Calculating and printing the run time needed to search for a pattern match.
79 print("Time taken for implementing (Naive String Matching algorithm) using Hash Table:", end_time - start_time)
80
81
```

- *Sample Output:*




```
Pattern found at index 1
Pattern found at index 5
Pattern found at index 6
Pattern found at index 2
Pattern found at index 5
Pattern found at index 3
Pattern found at index 1
Pattern found at index 2
Pattern found at index 1
Pattern found at index 5
Pattern found at index 5
Pattern found at index 2
Pattern found at index 5
Pattern found at index 6
Pattern found at index 2
Pattern found at index 3
Pattern found at index 2
Pattern found at index 5
Pattern found at index 2
Pattern found at index 5
Pattern found at index 1
Pattern found at index 1
Time taken for implementing (Naive String Matching
algorithm) using Hash Table: 0.008225399999901128
```

This screenshot is a part of the output from running (Naive Algo – Hash Table) python file.

as shown, the output will be the index in which a pattern match has been found, and that's for every data element (column by column) in the hash table.



Results & Analysis

- Comparison and analysis for the overall performances for each algorithm as well as the data structures, shown below.
- 

- **Device specifications:**

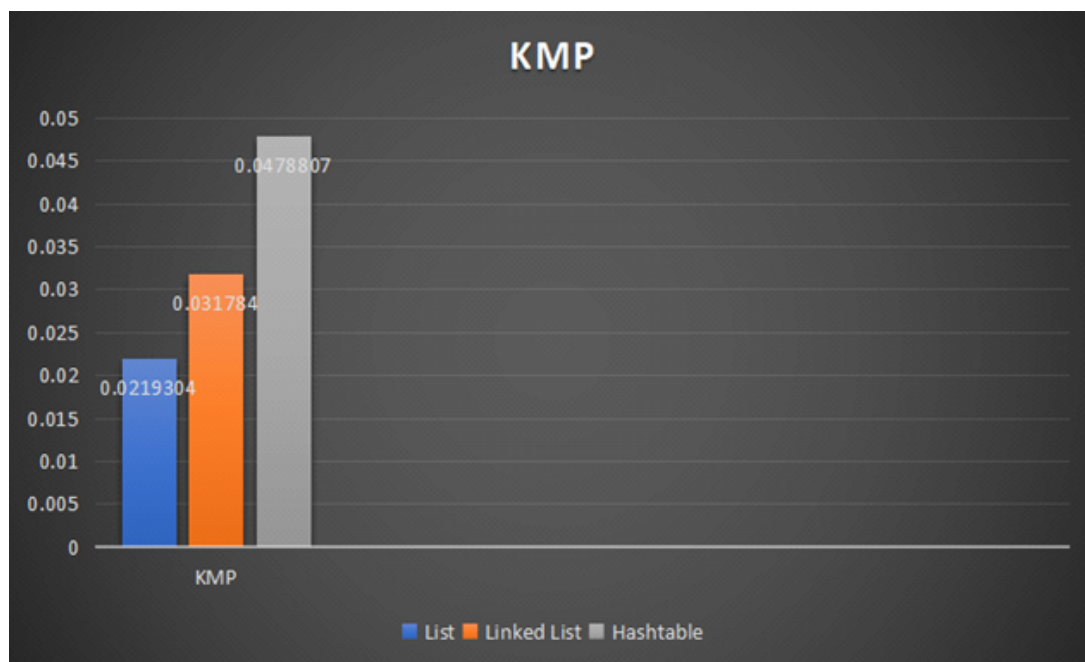
MacBook Pro 2020.

- M1 chip memory
- 16 gigabyte storage
- SSD disk 256 gigabyte
- Mac OS operating system

- **Data size:**

Data from both_covid_data.csv file, its size equals 2647.

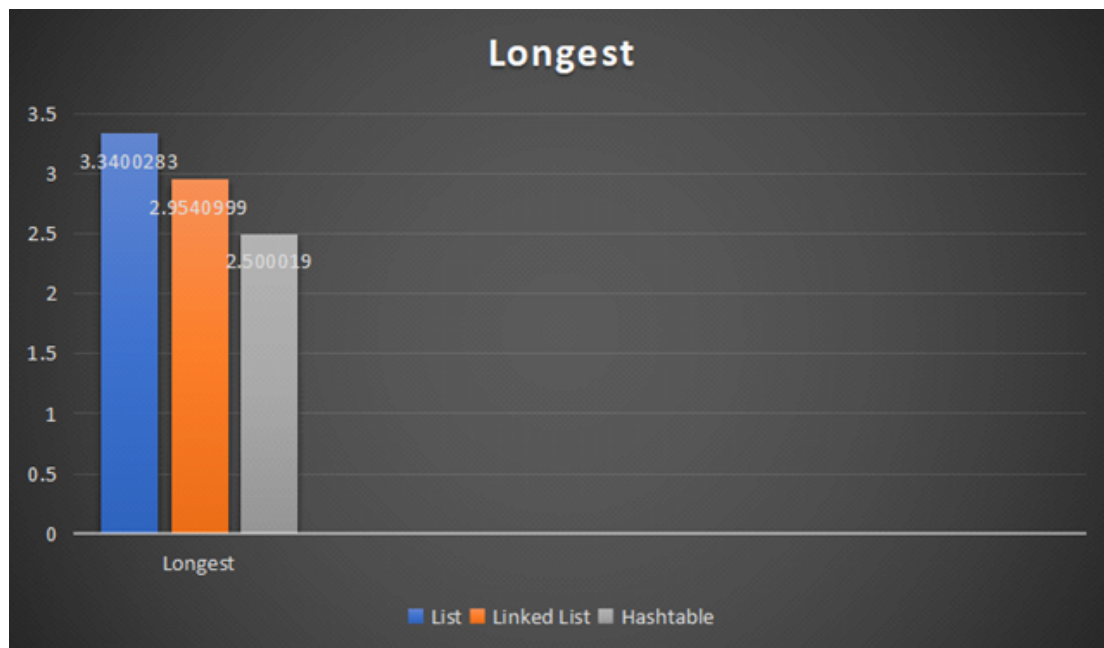
- **Results and analysis:**



List is the fastest in KMP because it took less time in the implementation process 0.0219304. and the list faster than the linked list of data structures because they are more efficient at searching for patterns, List data structures use an array of elements that can be searched in a linear fashion, while linked list data structures require traversing the entire list to find a pattern, The linear search of the list data structure is much faster than the traversal of the linked list structure.

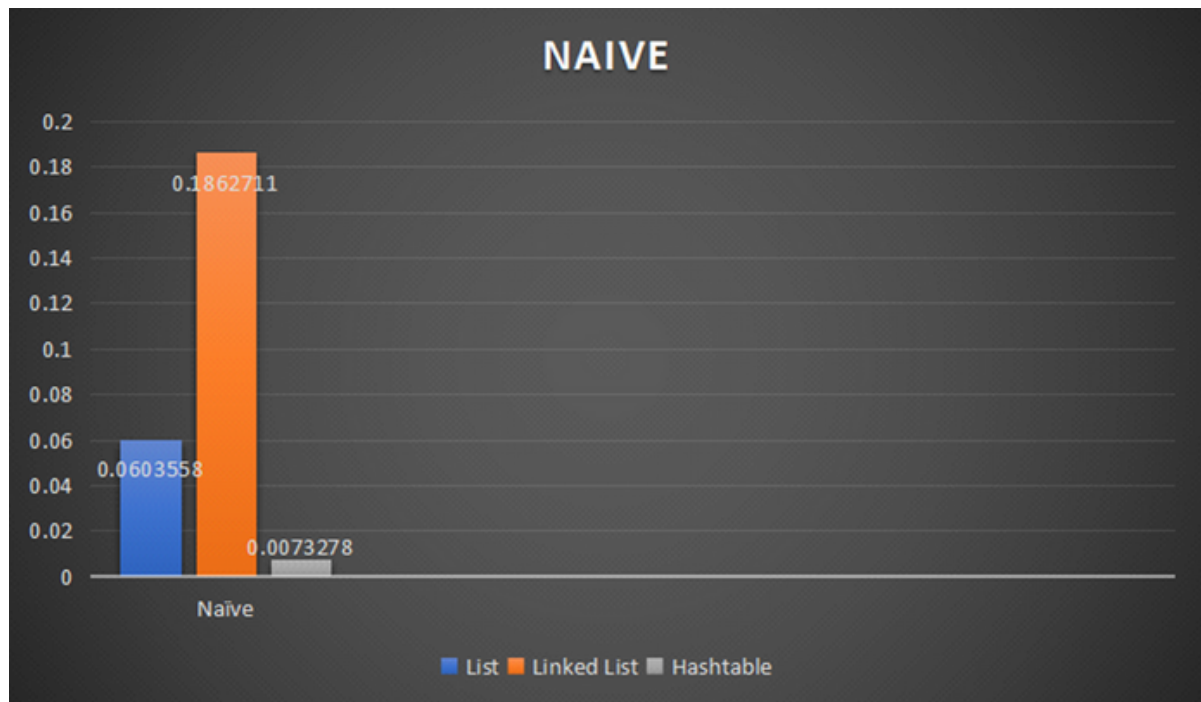
and faster than hash table because they can quickly traverse through the list of characters in the pattern to find a match, Hash tables require more time to search for a particular character in the pattern, as they must hash each character and then compare it to the stored value. List data structures can quickly traverse through the list of characters and find a match, making them faster than hash tables for KMP string matching.

- **Results and analysis (cont.):**



hash table fastest in longest because it took less time in the implementation process 2.500019. Hash tables are faster than lists because they use a hashing algorithm to quickly locate the desired string. This means that instead of having to search through an entire list, the hash table can quickly locate the desired string by using its hash value. This makes it much faster than searching through a list, which requires linear time complexity, and faster than linked list because they allow for constant time lookups. This means that when searching for a string, the algorithm can quickly determine if the string is present in the hash table without having to traverse through a linked list. Additionally, hash tables can store more data than linked lists, which allows for more efficient storage and retrieval of data.

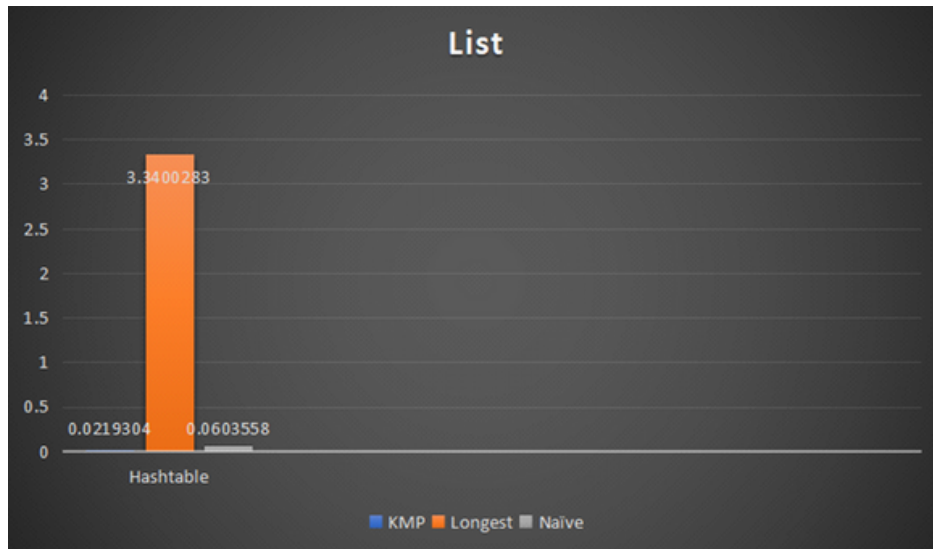
- **Results and analysis (cont.):**



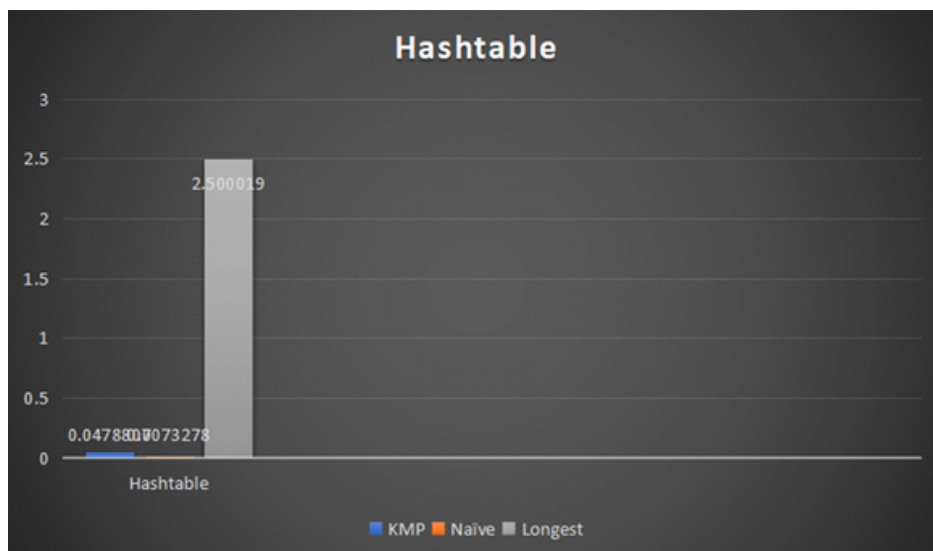
hash table fastest in naive because it took less time in the implementation process 0.0073278. Hash tables are faster than linked lists in naive string matching algorithms because they allow for faster lookups. Hash tables use a hashing function to map keys to values, which makes it easier and faster to find the desired value. Linked lists, on the other hand, require linear search time, which can be slow if the list is long , and faster than lists in naive string matching algorithms because they allow for faster lookups. Hash tables use a hashing function to map a key to an index in the table, which allows for constant time lookup of the value associated with the key. Lists, on the other hand, require linear time lookup since each element must be searched through sequentially. This makes hash tables much more efficient when searching for a specific value.

- **Results and analysis (cont.):**

(comparing between the same data structure used to implement each algorithm)



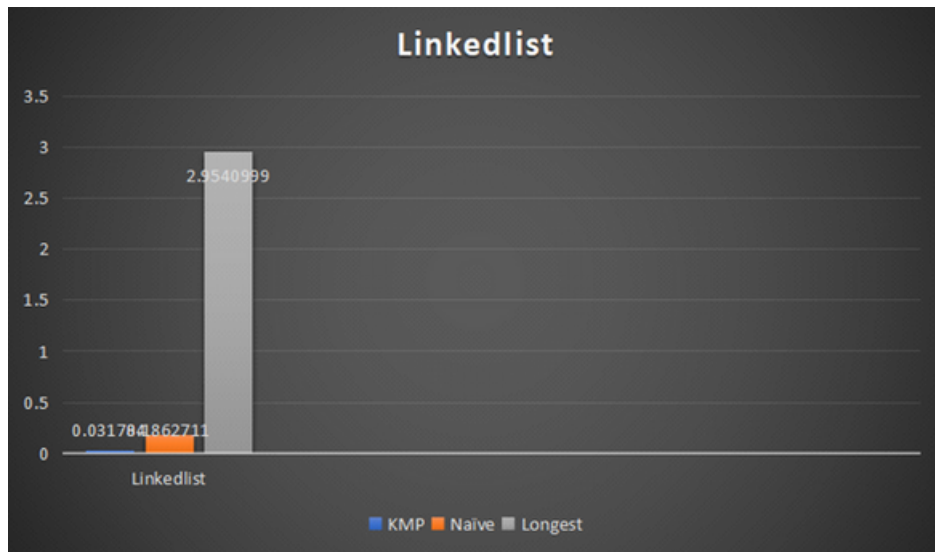
KMP is faster in the list data structure because it takes 0.0219304 for execute and The time complexity of the KMP string matching algorithm using a List is $O(n)$. This is because the algorithm only needs to iterate through the list once to find a match, and it does not need to perform any additional operations.



naive is faster in the hash table data structure because it takes 0.0073278 for execute and The time complexity of a naive string matching algorithm using a hash table is $O(n)$. This is because the algorithm must iterate through each character in the string to generate a hash value, and then compare the generated hash value with the stored hash values in the table.

- **Results and analysis (cont.):**

(comparing between the same data structure used to implement each algorithm)



KMP is faster in the linked list data structure because it takes 0.031784 for execute and the time complexity of the KMP string matching algorithm using a linked list is $O(n+m)$, where n is the length of the string and m is the length of the pattern.

• *Group Work Report:*

	Ruba Balubaid	Hadeel Alnasiri	Yara Al simlan	Noor Alhashmi	Shaden Anagreh
Naive codes (3 data structures)		✓			
KMP codes (3 data structures)	✓				
LCS codes (3 data structures)					✓
Comparisons			✓		
Project Report	✓	✓		✓	
Project presentation	✓				