# String Matching Assignment

KPR Lakmal

21001065

## Introduction

In this assignment I use horspool algorithm to implement the code.

In addition to that I have included four regular expression patterns and I have implemented four functions for each regular expression pattern.

**Regular expression pattern I have used.**

1. `.` (dot):

   - Pattern: `b.t`

   - Matches: "bat", "bet", "bit", "bot", "but", etc.

   - Explanation: The dot matches any single character except a newline.

2. `?` (question mark):

   - Pattern: `colou?r`

   - Matches: "color" and "colour."

   - Explanation: The question mark matches zero or one occurrence of the preceding character, in this case, the letter "u."

3.`^` are used to specify the beginning of a line or string.

4.`$` are used to specify the end of a line or string

# Why HorsPool.

Because Horspool algorithm matches from the end of the pattern and when there is a mismatch (bascharacter) found it's not going to search for the next right most occurrence of the bad character. Instead it's searching for the right most occurrence of the last character of the current text window in the pattern. There it skips lot of unwanted matches at once. In our case the text file contains a huge amount of texts. High number of skips decreases the time to search for a pattern in a text. Therefore Horspool algorithm is an efficient algorithm in this scenario.

```python
def horspool(text,pattern,text_file):

    alphabet = set(text + pattern)
    distlist = list(alphabet)
    m = len(pattern)
    n = len(text)

    dict={}
    for i in range(0,len(distlist)):
        dict[distlist[i]] = m
    for j in range(0,m-1):
        dict[pattern[j]] = m-j-1
    pos = 0
    count =0
    while pos<= n-m:
        j = m-1
        while j>=0 and text[pos+j] == pattern[j]:
            j = j-1
        if j==-1:
            count+=1
            linetwo = ("\n\tfound position is "+str(pos)+ " to " +str(pos+m-1)+" in the following text line")
            #print(linetwo)
            text_file.write(linetwo)
        pos = pos +dict[text[pos+m-1]]
    if count>0:
        linethree=("\n\t"+text)
        #print(linethree)
        text_file.write(linethree)

    return count
```

Above code describes how the horspool algorithm works. I have created a function named horspool and took text , pattern and output file as inputs.

And if there exist any matching in text then write matching line and matching position in output file.

**How it works:**

```python
alphabet = set(text + pattern)
distlist = list(alphabet)
m = len(pattern)
n = len(text)

dict={}
for i in range(0,len(distlist)):
    dict[distlist[i]] = m
```

As we know the horspool algorithm uses the last character of the current text window instead of bad character. In my solution, I've used a dictionary in python to keep track of the right most position of each character in the pattern so that it will be our HpBc table. For the HpBc table I've considered characters in both text and the pattern. Because there might be some characters in the text which cannot be found in the pattern anywhere and vice versa. So I've created a set and then converted it to a list so that I can access each character by an index. And for each character we need to instantiate a key in the dictionary. Initially all the keys have the same value which is the length of the pattern.

After that we need to calculate real right most position of the each character.

```python
for j in range(0,m-1):
    dict[pattern[j]] = m-j-1
```

Then we need two index variables to iterate thrplough the text and the pattern. For the text I've used variable "pos" which is initially equals to 0 and on each iteration of pos, we define another variable j which is initially equals to lentht(pattern)-1. Starting from the last character of the current text window. Value of j will be decreased if the two characters are equals or else value of pos will increase if there was a mismatch. To index "pos" to next suitable point we're using the pre-calculated HpBc table.

```
pos +=dict[text[pos+m-1]]
```

And the following part of the code check the pattern has any regular expression character. If there exist any regex character then its function call using if else statements.

```
for line in text:
    if pattern.find(".")!= -1:
        count+= horspool_dot(line,pattern,text_file)
    elif pattern.find("?")!= -1:
        count+= horspool_quest(line,pattern,text_file)
    elif pattern.find("^")!= -1:
        count+= horspool_begg(line,pattern,text_file)
    elif pattern.find("$")!= -1:
        count+= horspool_end(line,pattern,text_file)
    else:
        count+= horspool(line,pattern,text_file)
```

I have implemented four various functions for each regex character.

dot( . )                => call horspool_dot()

question mark( ? ) => call horspool_quest()

dollor sign ( $ )      => call horspool_begg()

^ mark ( ^ )           => call horspool_end()

These are some Test and results

# Test One

1. **Pattern.txt => " tr.ck "**
2. **Text.txt**

```
The word technology track and its uses have
immensely changed since the 20th century,
and with time, it has continued to evolve ever since tricky.
We are living in truck a world driven by technology.
trtck gThe advancement of technology has played an
important role in the development of human civilization,
along with cultural changes trick Technology provides
innovative ways of doing work through various smart and innovative means    trock.
hello evryone
```

**3.Output1.txt**

```
text1.txt Searching pattern is "tr.ck"
    found position is 20 to 24 in the following text line
    The word technology track and its uses have

    found position is 53 to 57 in the following text line
    and with time, it has continued to evolve ever since tricky.

    found position is 17 to 21 in the following text line
    We are living in truck a world driven by technology.

    found position is 0 to 4 in the following text line
    trtck gThe advancement of technology has played an

    found position is 28 to 32 in the following text line
    along with cultural changes trick Technology provides

    found position is 75 to 79 in the following text line
    innovative ways of doing work through various smart and innovative means    trock.

Number of matches found : 6
```

# Test two

1.Pattern.txt => " ^tree "

2.text.txt

```
1    tree have leaves
2    please leave me alone
3    I want grean leaves
4    trees are our lives
5    notrees can talk
6    trees produced oxigen treehello
7    hello treejack
```

3.output.txt

```
2    text3.txt Searching pattern is "^tree"
3        found position is 0 to 3 in the following text line
4        tree have leaves
5
6        found position is 0 to 3 in the following text line
7        trees are our lives
8
9        found position is 0 to 3 in the following text line
10       found position is 22 to 25 in the following text line
11       trees produced oxigen treehello
12
13       found position is 6 to 9 in the following text line
14       hello treejack
15   Number of matches found : 5
```

# Test Three

1. Pattern => "colou?r"
2. Text.txt

```
1    I have good colours
2    my favourite color is blue
3    we have colourfull books
4    what is your favourite
5    everything without colors are black and white
6    I want color pencils
7    my hobby is painting
8    Some colours can make us relax
```

3.output.txt

```
1
2    text4.txt Searching pattern is "colou?r"
3        found position is 12 to 17 in the following text line
4        I have good colours
5
6        found position is 13 to 17 in the following text line
7        my favourite color is blue
8
9        found position is 8 to 13 in the following text line
10       we have colourfull books
11
12       found position is 19 to 23 in the following text line
13       everything without colors are black and white
14
15       found position is 7 to 11 in the following text line
16       I want color pencils
17
18       found position is 5 to 10 in the following text line
19       Some colours can make us relax
20   Number of matches found : 6
```

# Test Four

1. Pattern => "riya$"
2. Text.txt

```
1   There is a saying that you can see the best of Sri Lanka's cultural monuments in the three ancient cities,
2   which are situated next to one another. These are Anuradhapura, Polonnaruwa and Sigiriya.
3   I visited the last two. Polonnaruwa was interesting, especially because it covers a huge place with many monuments,
4   from the smallest to the ones that are much bigger. But Sigiriya was the highlight of the Sri Lankan culture for me.
5   Sigiriya (or Sinhagiri) is an ancient rock fortress located near the town of Dambulla. The name Sigiriya refers to
6   a site of historical and archaeological significance that is nearly 200 metres high. If you ask anyone in Sri Lanka about Sigiriya,
7   they will proudly say it's the 8th wonder of the world. It is a massive monument, settled in the middle of an enormous fortress.
8   If you want to get to the top, you must conquer very steep steps, which spiral upwards. The climb looks more daunting than it actually is.
9   The citadel was built just before 500 BC, and is surrounded by a double moat, the outer part of which is now dry. The water gardens,
10  largely dry now, are a restful place to pass time before or after your climb. temple
```

## 3.output.txt

```
1
2   text7.txt Searching pattern is "riya$"
3       found position is 84 to 87 in the following text line
4       which are situated next to one another. These are Anuradhapura, Polonnaruwa and Sigiriya.
5
6       found position is 60 to 63 in the following text line
7       from the smallest to the ones that are much bigger. But Sigiriya was the highlight of the Sri Lankan culture for me.
8
9       found position is 4 to 7 in the following text line
10      found position is 100 to 103 in the following text line
11      Sigiriya (or Sinhagiri) is an ancient rock fortress located near the town of Dambulla. The name Sigiriya refers to
12
13      found position is 126 to 129 in the following text line
14      a site of historical and archaeological significance that is nearly 200 metres high. If you ask anyone in Sri Lanka about Sigiriya,
15
16  Number of matches found : 5
```