

UNIVERSITÉ LIBRE DE BRUXELLES

DATA STRUCTURES AND ALGORITHMS

INFO-F-413

---

# String Matching

---

*Author:*

Charlotte NACHTEGAEL

*Supervisor:*

Pr. Jean CARDINAL

December 2016



# 1 Introduction

Looking for patterns in a text is a problem encountered in a lot of situations, such as text-editing programs, in bioinformatics when searching for pattern in the DNA or even when searching for relevant queries with search engines on Internet.

We formalize the string matching problem as follows : we consider a pattern  $P$  as an array  $P[1...m]$  of length  $m$  and a text  $T$  as an array  $T[1...n]$  of length  $n$ , with  $m \leq n$ . The elements of the arrays are drawn from a finite alphabet  $\Sigma$  with a size  $|\Sigma|$ . We define that a pattern  $P$  occurs with a shift  $s$  in the text  $T$  if  $0 \leq s \leq n - m$  and  $P[1...m] = T[s+1...s+m]$ . If  $P$  occurs in  $T$  with a shift  $s$ , this shift is qualified **valid**, otherwise it is qualified **invalid**. The string-matching problem is to find all the valid shifts for the pattern  $P$  in the text  $T$ .

We will present during this report three different strategies for string matching: the Naive one, the Knuth-Morris-Pratt (KMP) algorithm [Knuth, 1977] and the Boyer-Moore (BM) algorithm [Boyer and Moore, 1977]. The notations used are from the book *Introduction to Algorithms, third edition* [Cormen et al., 2009] (CLRS). The aim of this report is to present the worst cases analysis and some average case analysis of the number of comparisons of characters made with these different strategies, proving it by demonstration and back these affirmations with experimental results of original implementations.

# 2 Necessary terminology and notation

The set of all finite strings is noted  $\Sigma^*$ . The **length** of a string  $x$  is denoted  $|x|$ . A string  $w$  is a **prefix** of the string  $x$  if  $x = wy$  for any  $y \in \Sigma^*$  and is noted  $w \sqsubset x$ .  $w$  is a **suffix** of  $x$  if  $x = yw$  for any  $y \in \Sigma^*$  and is noted  $w \sqsupset x$ . In both cases,  $|w| \leq |x|$ .

The  $k^{th}$  prefix of the pattern  $P[1...m]$ ,  $P[1...k]$ , is noted as  $P_k$ . Similarly the  $k^{th}$  prefix of the text  $T$  is noted  $T_k$ . Thusly, we can note the string matching problem as finding all the shifts  $s$  such as  $P \sqsubset T_{s+m}$ , with  $0 \leq s \leq n - m$ .

### 3 Naive algorithm

The naive algorithm is based on the idea of matching all the characters of the pattern for all the possible shifts, of the number of  $n-m+1$ , until there is mismatch, then the pattern shift from 1 to the right to begin to check the next shift. The computational time is  $O(m(n-m))$ . The maximum number of comparisons is  $m(n-m)$ , because you compare all your characters at each possible shift. This could happen when you have all your text with "a" and you pattern composed of at least  $(m-1)*"a"$ . This means that you would have to check all your pattern ( $m$  characters) before finding if your patterns occurs with shift  $s$  or not.

We implemented a worst case scenario where the text was "a"\*43 and the pattern was "ab" and we obtained a number of comparisons equal to  $m(n-m)$ .

In the CLRS book [Cormen et al., 2009], they proposed an upper bound value for the number of comparisons of characters, for randomly chosen pattern  $P$  and text  $T$  from an alphabet  $\Sigma_d = \{0, 1, \dots, d-1\}$  with  $d \geq 2$ , where the comparison stops when a mismatch occurs, as follows:

$$(n-m+1) \frac{1-d^{-m}}{1-d^{-1}} \leq 2(n-m+1)$$

This affirmation is based on the expected number of trials before success :

$$E[X] = \frac{1 - (1-p)^n}{p},$$

where  $p$  is the probability of success and  $E[X]$  the sum of the probabilities of the success after the  $i$ th trial. In this case, the success is defined as the fact that, after  $i-1$  matches, the  $i$ th character does not, with a probability of  $1 - \frac{1}{d}$ . The final success case is when all the characters match, so all  $m$  comparisons were made. Thusly, we obtain :

$$E[X] = \frac{1 - \left(1 - \left(1 - \frac{1}{d}\right)\right)^m}{1 - \left(\frac{1}{d}\right)} \quad (1)$$

$$E[X] = \frac{1 - d^{-m}}{1 - d^{-1}} \quad (2)$$

This expected number of comparisons is used for all possible shift, so  $(n - m + 1)$  times, this is why the total expected number is  $(n - m + 1) \frac{1 - d^{-m}}{1 - d^{-1}}$ . To complete the proof, we must, prove that  $\frac{1 - d^{-m}}{1 - d^{-1}}$  is less or equal to 2.

As said in the beginning,  $d \geq 2$ , so this means that if  $d = 2$ ,  $1 - d^{-1} = 0.5$  and as  $d$  increases, so does  $1 - d^{-1}$ ; and  $1 - d^{-m}$  will always be less than 1. Consequently  $\frac{1 - d^{-m}}{1 - d^{-1}} \leq 2$  ( $1/0.5$  for  $d = 2$  and always decreasing as  $d$  increases so the denominator increases and the numerator is asymptotically close to 1), confirming the affirmation that  $(n - m + 1) \frac{1 - d^{-m}}{1 - d^{-1}} \leq 2(n - m + 1)$  for the average number of comparisons for a random text and pattern.

We implemented the naive algorithm with text of different lengths from 1,000 to 10,000 and patterns of length 100, created randomly from a alphabet  $\Sigma = \{A, G, T, C\}$ . We averaged the results over 10 patterns of lengths 100. The results are plotted with the upper bound. We can see that the upper bound is largely overestimated when the text length increases (Fig 1). However, we can distinguish that when the text length is small, the upper bound is quite close to the real number of comparisons.

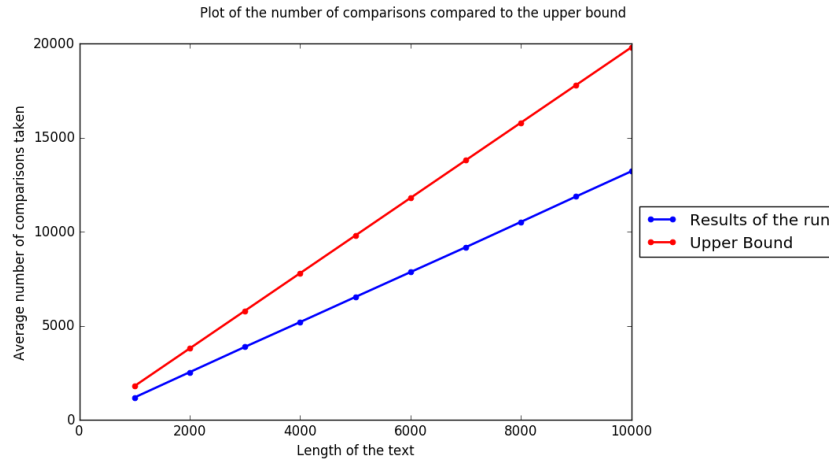


Figure 1: Average number of comparisons for different lengths of text averaged over 10 patterns of length 100 for the naive algorithm

## 4 KMP algorithm

The Knuth-Morris-Pratt algorithm is based on the construction of a prefix table designated as  $\pi$ , with  $\pi[i] = \max\{k : k < i \text{ and } P_k \sqsupseteq P_i\}$ . This value  $\pi[i]$  is used to shift after mismatch towards a place where a string of characters equal to the prefix of the pattern  $P$  were already matched, so as to avoid in-between invalid shifts. The worst case scenario is to obtain  $2n$  comparisons. This is due to the fact that at each iteration of the for loop (line 6), we either increase  $q$  or slide the pattern right. These two events can occur at most  $n$  times, so the loop takes at most  $2n$  times [Knuth, 1977]. This worst case happens in the same condition as described earlier for the naive algorithm, with a text composed only of "a" characters and a pattern equal to "ab". We confirmed with an implementation that the number of comparisons was indeed equal to  $2n$ .

**Algorithm 1** KMP string matching algorithm

---

```

1: procedure KMP_MATCHER( $T, P$ )
2:    $n \leftarrow T.length$ 
3:    $m \leftarrow P.length$ 
4:    $\pi \leftarrow COMPUTE - PREFIX - FUNCTION(P)$ 
5:    $q \leftarrow 0$ 
6:   for  $i \leftarrow 0$  to  $n$  do
7:     while  $q > 0$  and  $P[q + 1] \neq T[i]$  do
8:        $q \leftarrow \pi[q]$ 
9:     end while
10:    if  $P[q + 1] == T[i]$  then
11:       $q \leftarrow q + 1$ 
12:    end if
13:    if  $q == m$  then
14:      print "Pattern occurs with shift"  $i - m$ 
15:       $q \leftarrow \pi[q]$ 
16:    end if
17:  end for
18: end procedure

```

---

The implementation was done as previously described for the naive algorithm to prove the correctness of the upper bound. We can see that, again, the upper bound is much bigger than the concrete number of comparisons (Fig 2). This is due here to the fact that, despite the randomness of both the text and the pattern, the prefix of the pattern can be found along the rest of the pattern, allowing a more efficient shift along the text. The difference between the upper bound and the results is more marked than with the naive algorithm. This is the proof of the efficiency of the KMP algorithm compared to the naive method. Indeed, we can analyze the complexity of the KMP-matcher with an amortized analysis. The potential function is in function of  $q$ . The initial value of  $q$  is 0 (line 5). The value of  $q$  decreases at line 8 and 15 since  $q < \pi[q]$  and since  $\pi[q] \geq 0$ , the condition that the potential function is always greater or equal than 0 is respected. The value of  $q$  is incremented by 1 at line 11 at most once by for loop. The cost of the while loop at line 7-9 is paid by the decrease of the value of  $q$ . The cost of the for loop is paid by the incrementation at line 11, so that the for loop operation cost constant time  $O(1)$ . As the for loop is repeated  $n$  times, the cost of only the searching part of the KMP-matcher is  $O(n)$ .

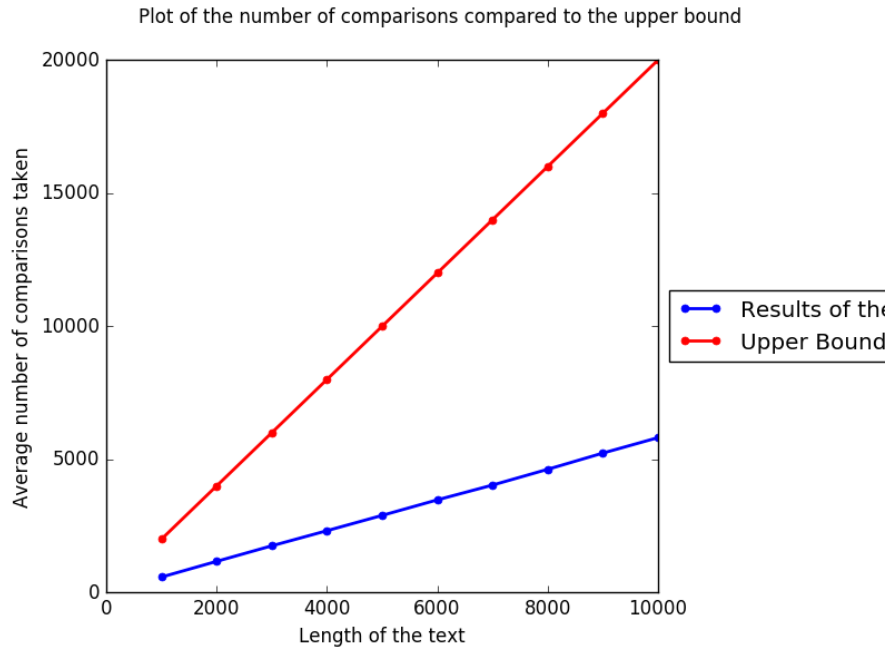


Figure 2: Average number of comparisons for different lengths of text averaged over 10 patterns of length 100 for the KMP algorithm

## 5 BM algorithm

The Boyer-Moore algorithm is based on the idea of matching from right to left with a shift calculated either with a suffix rule (inspired by the prefix rule in the KMP algorithm) or the bad character rule, where you look if the last bad match character can be found in the pattern at the most left position. The shift chosen is the biggest one. The worst case analysis by Knuth [Knuth, 1977] bounded at  $7n$  comparisons. However, Colussi found another upper bound of  $4n$  [Colussi, 1991]. While Guibas conjectures for  $2n$  comparisons [Guibas, 1980], the current used upper bound number of comparisons is  $3n$  for a non-periodic pattern [Cole, 1994].

We implemented as previously described to check the veracity of the upper bound. Again, we constated that the upper bound is largely overestimated (Fig 3). However, the difference here is even more marked than with the KMP algorithm. This means that the efficiency of the BM algorithm with random text and pattern is quite high. This result is due to the fact that BM algorithm shifts according to two different rules, making it able to shift as much as possible.

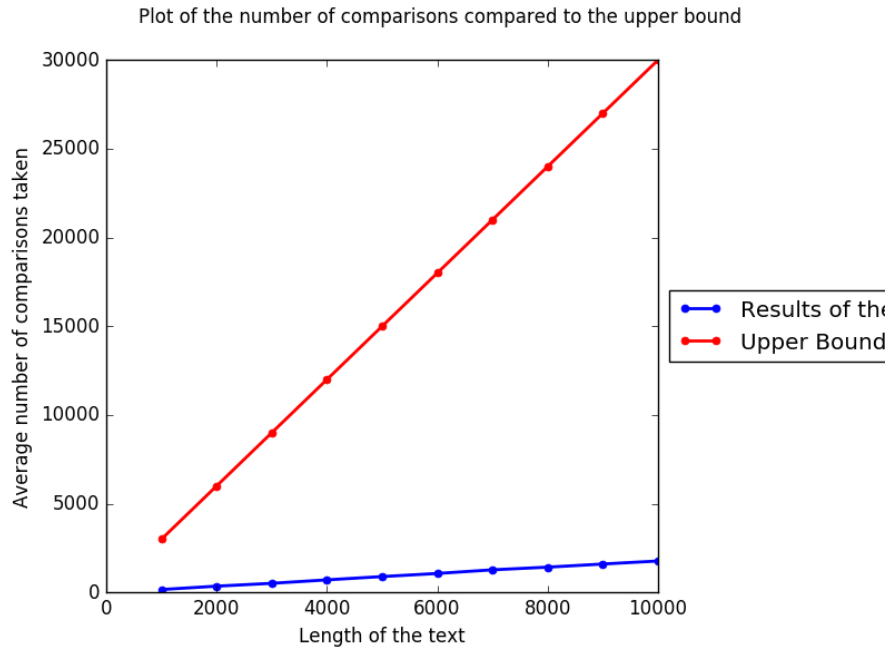


Figure 3: Average number of comparisons for different lengths of text averaged over 10 patterns of length 100 for the BM algorithm

## 6 Comparison of running time

We compared the time taken by different string matching algorithm. The text was of a length 10,000 and we modified the lengths of the pattern, from 10 to 500. We ran 10 different runs with 10 random words for each length of pattern, all created from an alphabet  $\Sigma = \{0, 1\}$ .

We compared the naive, the Rabin-Karp, the Knuth-Morris-Pratt and the Boyer-Moore algorithms, respectively with average complexity of  $O(mn)$ ,  $O(m + n)$ ,  $O(m + n)$  and  $O(m/n)$  in the best case. The Rabin-Karp algorithm is based on a rolling hash function. The pattern and the current part of the text of the same length as the pattern is hashed and compared as two numbers. A comparison of characters is only done when the two numbers are equal. The rolling hash function allows to compute the hash of the next text window in a constant time.



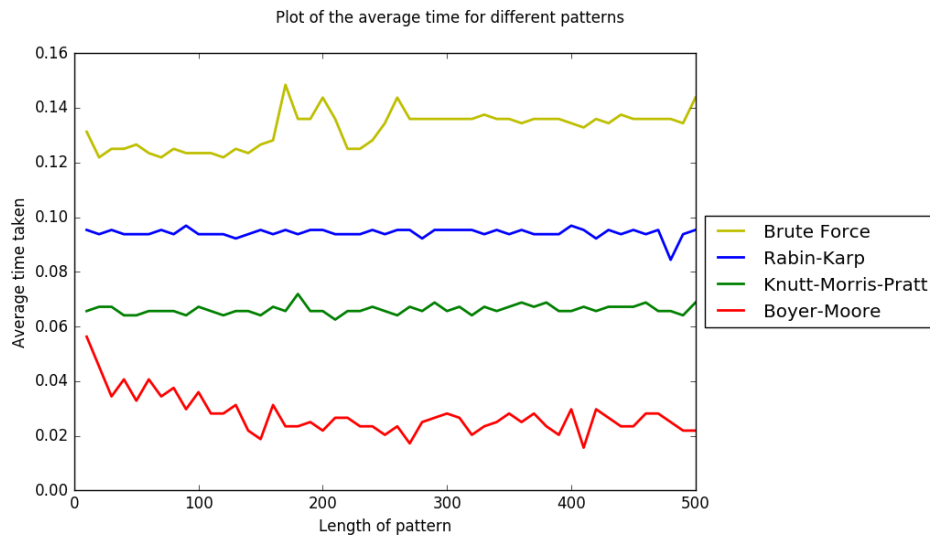


Figure 4: Average time taken by the different string matching algorithms for different lengths of pattern. The results were obtained over 10 runs of 10 random words for each length of pattern confronted to a random 10,000 text, created from a binary alphabet.

We can indeed observe that the naive algorithm performs the worst and that the Boyer-Moore performs the best (Fig 4). Even though the Rabin-Karp and the KMP algorithms have the same time complexity, in practice, we observe that the KMP performed better. This could be due to the fact that, even if the Rabin-Karp compares the pattern and the text only when a match is found with the two strings hashed. However, we could have what is called a **spurious hit**, a false positive hit obtained because of a collision of the hash function. This situation can be avoided by choosing the  $q$  parameter of the Rabin-Karp algorithm (a prime number usually chosen to be greater than the length of the pattern) is chosen randomly.

Rabin-Karp and KMP are also fairly linear, whereas we can see that the complexity of the naive algorithm increase with the length of the pattern and the BM complexity decreases with the length of the pattern. These results concords with the fact that the BM algorithm is the most used string matching algorithm in general compared to the other ones.

## References

- Robert S. Boyer and J. Strother Moore. A fast string searching algorithm. *Commun. ACM*, 20(10):762–772, October 1977. ISSN 0001-0782. doi: 10.1145/359842.359859.
- Richard Cole. Tight bounds on the complexity of the boyermoore string matching algorithm. *SIAM Journal on Computing*, 23, 10 1994. doi: 10.1137/s0097539791195543. URL <http://gen.lib.rus.ec/scimag/index.php?s=10.1137/s0097539791195543>.
- Livio Colussi. Correctness and efficiency of pattern matching algorithms. *Information and Computation*, 95(2):225 – 251, 1991. ISSN 0890-5401. doi: [http://dx.doi.org/10.1016/0890-5401\(91\)90046-5](http://dx.doi.org/10.1016/0890-5401(91)90046-5). URL <http://www.sciencedirect.com/science/article/pii/0890540191900465>.
- Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms, Third Edition*. The MIT Press, 3rd edition, 2009. ISBN 0262033844, 9780262033848.
- Andrew M. Guibas, Leo J.; Odlyzko. A new proof of the linearity of the boyer-moore string searching algorithm. *SIAM Journal on Computing*, 9, 11 1980. doi: 10.1137/0209051. URL <http://gen.lib.rus.ec/scimag/index.php?s=10.1137/0209051>.
- Jr. James H.; Pratt Vaughan R. Knuth, Donald E.; Morris. Fast pattern matching in strings. *SIAM Journal on Computing*, 6, 06 1977. doi: 10.1137/0206024.

## Code

```
1 """
2 INFO-F413 : Data Structures and Algorithms
3 All Algorithms for the Project 2016
4 Theme : String Matching
5 Naive, Rabin-Karp, Knutt-Morris-Pratt and Boyer-Moore
6 """
7
8 import time
9 import numpy as np
10 import matplotlib.pyplot as plt
11
12
13 def naive(text, pattern):
14     nb_comp = 0
15     for i in range(len(text)-len(pattern)):
16         for j in range(len(pattern)):
17             nb_comp += 1
18             if text[i+j] != pattern[j]:
19                 break
20             elif j == len(pattern) - 1:
21                 print("Pattern occurs with shift", i - 1)
22     return nb_comp
23
24
25 def Rabin_Karp(text, pattern, d, q):
26     """
27     :param text: string where you look for the pattern
28     :param pattern: string you are looking for
29     :param d: number, radix, generally the number of characters used
30     :param q: prime number
31     """
32     n = len(text)
33     m = len(pattern)
34     h = d ** (m - 1) % q
35     p = 0
36     t = 0
37
38     for i in range(m):
39         p = (d * p + ord(pattern[i])) % q
40         t = (d * t + ord(text[i])) % q
41
42     for s in range(n - m + 1):
```

```
43     # hit
44     if p == t:
45         # check if it is a spurious hit or not
46         if pattern == text[s:]:
47             print("The pattern occurs at shift", s-1)
48     if s < n - m:
49         t = (d*(t - ord(text[s])* h) + ord(text[s + m])) % q
50
51
52 def Knutt_Morris_Pratt(text, pattern):
53     n = len(text)
54     m = len(pattern)
55     prefix = prefix_function(pattern)
56
57     nb_comp = 0
58     q = -1 # number of character matched
59     for i in range(n):
60         # no match for next character when a match was already found
61         while q > -1 and pattern[q + 1] != text[i]:
62             nb_comp += 1
63             q = prefix[q]
64         # match with next character
65         if pattern[q + 1] == text[i]:
66             nb_comp += 1
67             q = q + 1
68         # complete match
69         if q == m - 1:
70             # print("Pattern occurs with shift", i - m + 1)
71             q = prefix[q]
72     nb_comp += 1
73     return nb_comp
74
75
76 def Boyer_Moore(text, pattern, alphabet):
77     n = len(text)
78     m = len(pattern)
79     delta = last_occurrence(pattern, alphabet)
80     gamma = good_suffix(pattern)
81
82     s = 0
83     nb_comp = 0
84
85     while s <= n - m:
86         j = m - 1
```

```

87     # good match
88     while j > 0 and pattern[j] == text[s + j]:
89         j = j - 1
90         nb_comp += 1
91     #complete march
92     if j == 0:
93         print("Pattern occurs at shift", s)
94         s = s + gamma[0]
95     else:
96         nb_comp += 1
97         # if gamma[j] > j - delta[text[s + j]]:
98         # print("Shift good suffix")
99         # else:
100        # print("Shift bad character")
101        s = s + max(gamma[j], j - delta[text[s + j]])
102    return nb_comp
103
104
105 def prefix_function(pattern):
106     """
107     Create a array where the ith entry is the longest suffix to i
108     also prefix of the pattern
109     """
110     m = len(pattern)
111     prefix = [-1] * m
112     k = -1
113     for i in range(1, m):
114         # next character of the pattern does not match the current
115         while k > -1 and pattern[k + 1] != pattern[i]:
116             k = prefix[k]
117         if pattern[k + 1] == pattern[i]:
118             k = k + 1
119         prefix[i] = k
120     return prefix
121
122 def good_suffix(pattern):
123     """
124     Create a suffix table where you can find the next left motif
125     similar to the suffix of the pattern
126     """
127     m = len(pattern)
128     pi = prefix_function(pattern)
129     pattern_prime = pattern[::-1]

```

```
129     pi_prime = prefix_function(pattern_prime)
130
131     gamma = []
132
133     for j in range(m):
134         gamma.append(m - 1 - pi[m - 1])
135
136     for l in range(m):
137         j = m - 1 - pi_prime[l] - 1
138
139         if gamma[j] > 1 - pi_prime[l]:
140             gamma[j] = 1 - pi_prime[l]
141
142     return gamma
143
144
145 def last_occurrence(pattern, alphabet):
146     """
147     Create a dictionary with the last occurrence (left most) of the
148     character in the pattern or -1 if the character does not occur
149     in the pattern
150     """
151     m = len(pattern)
152     delta = {}
153
154     for letter in alphabet:
155         delta[letter] = -1
156
157     for j in range(m):
158         delta[pattern[j]] = j
159
160     return delta
161
162
163 def create_random_dico(alphabet):
164     """
165     Create a dictionary of different length pattern with 10 words
166     each
167     """
168     dico = {}
169     for i in range(10, 510, 10):
170         dico[i] = []
171         while len(dico[i]) < 10:
172             word = np.random.choice(alphabet, i)
```

```
170         word = ''.join(word)
171         if word not in dico[i]:
172             dico[i].append(word)
173     return dico
174
175
176 def create_random_text(alphabet):
177     """
178     Create a list of random texts of different lengths
179     """
180     random_list = []
181     for i in range(1000, 11000, 1000):
182         text = np.random.choice(alphabet, i)
183         text = ''.join(text)
184         random_list.append(text)
185
186     return random_list
187
188
189 def main(text, dico, alphabet):
190     """
191     Run the algorithm to calculate the average time taken for the
192     matching
193     """
194     #create the prime list
195     list_primes = []
196     with open("primes.txt", "r") as primes:
197         for line in primes:
198             line = line.strip().split(' ')
199             for ele in line:
200                 if ele.isdigit():
201                     list_primes.append(int(ele))
202
203     current_algo = []
204     for size in sorted(dico):
205         current_size = []
206
207         for word in dico[size]:
208             t = time.process_time()
209             naive(text, word)
210             elapsed = time.process_time() - t
211             current_size.append(elapsed)
```

```
212         current_algo.append(np.longfloat(np.mean(np.array(
213             current_size))))
214     time_algo = np.array([current_algo])
215
216     current_algo = []
217     for size in sorted(dico):
218         current_size = []
219         for word in dico[size]:
220             q = np.random.choice(list_primes)
221             while q < len(word):
222                 q = np.random.choice(list_primes)
223             t = time.process_time()
224             Rabin_Karp(text, word, len(alphabet), q)
225             elapsed = time.process_time() - t
226             current_size.append(elapsed)
227         current_algo.append(np.longfloat(np.mean(np.array(
228             current_size))))
229
230     time_algo = np.append(time_algo, [current_algo], axis=0)
231
232     current_algo = []
233     for size in sorted(dico):
234         current_size = []
235
236         for word in dico[size]:
237             t = time.process_time()
238             Knutt_Morris_Pratt(text, word)
239             elapsed = time.process_time() - t
240             current_size.append(elapsed)
241         current_algo.append(np.longfloat(np.mean(np.array(
242             current_size))))
243
244     time_algo = np.append(time_algo, [current_algo], axis=0)
245
246     current_algo = []
247     for size in sorted(dico):
248         current_size = []
249         for word in dico[size]:
250             t = time.process_time()
251             Boyer_Moore(text, word, alphabet)
252             elapsed = time.process_time() - t
253             current_size.append(elapsed)
```



```
252         current_algo.append(np.longfloat(np.mean(np.array(
253             current_size))))
254     time_algo = np.append(time_algo, [current_algo], axis=0)
255
256     return time_algo
257
258
259 def plot_all(average_runs):
260     """
261     Plot the different average times obtained for the algorithms
262     for different length patterns
263     """
264     names = ('Brute Force', 'Rabin-Karp', 'Knutt-Morris-Pratt', '
265             Boyer-Moore')
266
267     fig = plt.figure()
268     ax = plt.axes()
269     colors = ['r', 'g', 'b', 'y']
270     list_names = list(names)
271     x = np.arange(10, 510, 10)
272
273     for algo in average_runs:
274         ax.plot(x, algo, color=colors.pop(), linewidth=2, label=
275             list_names.pop(0))
276     fig.suptitle('Plot of the average time for different patterns')
277     ax.set_xlim([0, 500])
278     ax.set_xlabel('Length of pattern')
279     ax.set_ylabel('Average time taken')
280
281     box = ax.get_position()
282     ax.set_position([box.x0, box.y0, box.width * 0.8, box.height])
283
284     # Put a legend to the right of the current axis
285     ax.legend(loc='center left', bbox_to_anchor=(1, 0.5))
286
287     plt.show()
288
289 def plot_comparison(average_list, upperbound_list):
290     """
291     Plot the average of comparisons confronted to the upper bound
292     limit
293     """
```

```

291 fig = plt.figure()
292 ax = plt.axes()
293 colors = ['r', 'b']
294 names = ['Results of the run', 'Upper Bound']
295 x = np.arange(1000, 11000, 1000)
296
297 ax.plot(x, average_list, marker='.', markersize = 10, color=colors.
298         pop(), linewidth=2, label=names.pop(0))
299 ax.plot(x, upperbound_list, marker='.', markersize = 10, color=
300         colors.pop(), linewidth=2, label=names.pop(0))
301 fig.suptitle('Plot of the number of comparisons compared to the
302             upper bound')
303 ax.set_xlim([0, 10000])
304 ax.set_xlabel('Length of the text')
305 ax.set_ylabel('Average number of comparisons taken')
306
307 box = ax.get_position()
308 ax.set_position([box.x0, box.y0, box.width * 0.8, box.height])
309
310 # Put a legend to the right of the current axis
311 ax.legend(loc='center left', bbox_to_anchor=(1, 0.5))
312
313 plt.show()
314
315 def average_time():
316     """
317     Count the average time taken by the different algorithms, average
318     on 10 runs and plot them all together
319     """
320     alphabet = ['0', '1']
321
322     # 100 000 characters
323     text = ''
324     for i in range(100000):
325         text += np.random.choice(alphabet)
326
327     average_runs = []
328
329     for i in range(10):
330         dico = create_random_dico(alphabet)
331         average_runs.append(main(text, dico, alphabet))
332     average_runs = np.mean(np.array(average_runs), axis=0)

```

```
331     plot_all(average_runs)
332
333
334 def number_comparisons():
335     """
336     Count the average number of comparisons for different algorithms
337     and plot it against their upper bound limit
338     """
339     alphabet=['A','G','C','T']
340
341     text_list = create_random_text(alphabet)
342
343     dico = []
344     while len(dico) < 10:
345         word = np.random.choice(alphabet, 100)
346         word = ''.join(word)
347         if word not in dico:
348             dico.append(word)
349
350     average_list = []
351     upper_bound_list = []
352     for text in text_list:
353         sum = 0
354         for word in dico:
355             sum += naive(text, word)
356             average_list.append(sum/len(dico))
357             upper_bound_list.append(2*(len(text)-100 + 1))
358     plot_comparison(np.array(average_list), np.array(upper_bound_list))
359
360     average_list = []
361     upper_bound_list = []
362     for text in text_list:
363         sum = 0
364         for word in dico:
365             sum += Knutt_Morris_Pratt(text, word)
366             average_list.append(sum/len(dico))
367             upper_bound_list.append(2*len(text))
368     plot_comparison(np.array(average_list), np.array(upper_bound_list))
369
370     average_list = []
371     upper_bound_list = []
372     for text in text_list:
```

```
372     sum = 0
373     for word in dico:
374         sum += Boyer_Moore(text, word, alphabet)
375         average_list.append(sum/len(dico))
376         upper_bound_list.append(3*len(text))
377     plot_comparison(np.array(average_list), np.array(upper_bound_list)
378 )
379
380 # worst case scenario
381 worst_text = "aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa"
382 worst_pattern="ab"
383 print("Worst case - simulation :", naive(worst_text, worst_pattern)
384       , "Bound :", len(worst_pattern)*(len(worst_text)-len(
385       worst_pattern)))
386 print("Worst case - simulation :", Knutt_Morris_Pratt(worst_text,
387       worst_pattern), "Bound :", 2*len(worst_text))
388 print("Worst case - simulation :", Boyer_Moore(worst_text,
389       worst_pattern, alphabet), "Bound :", 3*len(worst_text))
390
391 if __name__ == "__main__":
392     number_comparisons()
393     average_time()
```