

Comparing and Analyzing the Effectiveness of Various String Matching Algorithms Across Different Domains

Research Question - How do different string matching algorithms perform across various application domains, and what factors contribute to their varying effectiveness?

Words: not enough

Contents

1 Introduction

The impacts and applications of string matching algorithms extend far beyond their immediate textual applications; its quotidian use in domains such as bioinformatics, web crawling and security systems make it a pertinent area of study within computer science. There is also constant research being carried out to identify more efficient algorithms, and find ways in which to combine them with other disciplines to carry out certain tasks.

For instance, multiple studies have been recently published performing various comparative analyses of new string matching algorithms [4], and their applicability to a wide range of domains such as biological sequences [5], quantum computing [6], and intrusion detection systems [7]. These findings indicate the modern-day relevance and applicability of string matching algorithms and demonstrate its worthiness as a subject of investigation.

Similarly, this paper seeks to adopt an investigative style by examining **how different string matching algorithms perform across various application domains, and the factors that contribute to their varying effectiveness.**

To formulate a methodology that will directly address this aim, the research question can be broken down into three parts: measuring the performance of the mentioned string algorithms, analyzing the results by making links back to the theory in order to understand the reasons behind their varying efficiency, and linking these to different domains.

To examine this subject, three popular string matching algorithms were programmed and tested on their ability to carry out three tasks from three different domains. They are then evaluated on a set of criteria to determine their efficiency and applicability. To maintain focus on the research question, the reasoning and rationale behind all the decisions will be explained to demonstrate how they directly address the question.

The choice of this research question is both relevant and justifiable, reflecting the persistent demand for evaluating and optimizing string matching algorithms across a range of applications.

2 Theoretical Background

2.1 Topic Overview

String matching algorithms (often used interchangeably with string or pattern searching algorithms) are a subset of string algorithms. The term is composed of three parts: “string”, “matching”, and “algorithm”. Strings refer to an abstract data type consisting of a sequence of characters. Matching relates to its function to identify specific patterns within a larger dataset that match a criteria. As for algorithms, these can be defined as a set of instructions designed to carry out a function. These algorithms take in two inputs: a pattern (P), which is a sequence of characters being searched for within a larger sequence, a text (T). Within string matching problems, there exists two general variants [2]:

1. **Find occurrences of a pre-defined pattern in a previously unseen dataset**

Carried out using finite automata models [2] (an idealized machine used to recognize patterns in a given text; it accepts or rejects the input depending on whether the pre-defined pattern occurs in the text) or through the combinatorial properties of strings.

2. **Find occurrences of any identifiable patterns in a given text**

Carried out using finite automata and binary trees.

Linking into the mentioned problem variants, these algorithms can be further classified into two categories [1]: exact, and approximate string matching algorithms. As its name suggests, the former refers to finding occurrences of a pattern that match it to the character, whereas the latter allows for some variation or deviation. For the purpose of this investigation, the first problem variant, as well as an exact-string matching approach was chosen for all the algorithms, both for simplicity and control. Three exact string matching algorithms were chosen for comparison based on their popularity: Naive, Knuth-Morris-Pratt, and Rabin-Karp. Prior to their explanations, it is important to note that the following notation will be used:

Table 1: String Matching Algorithm Notation

m	=	Length of P
n	=	Length of T
T_j	:	For an input text T where $0 \leq j \leq n - 1$
P_i	:	For an input pattern P where $0 \leq i \leq m - 1$

2.2 Naive Algorithm

[8] Also known as the “brute-force” approach, this algorithm compares the first indexes of both P and T . If P_i equals T_j then i and j are incremented and further compared. If P_i reaches $m - 1$ then a variable containing the number of occurrences of P is incremented and P_i is reset to $P_{i=0}$. On the other hand, if they do not match, then P_i is reset and j incremented. This repeats until $T_{j=m-1}$ is reached. Since there is no pre-processing phase, the only time complexity taken into account is during its matching phase, which is expected to be $O(n \cdot m)$ [2]. The algorithm is better expressed in pseudocode:

Algorithm 1: Naive

```

1  $n \leftarrow T.length()$  ;
2  $m \leftarrow P.length()$  ;

3 for  $i \leftarrow 0$  to  $n - m$  do
4   while  $j < m$  and  $P[j] == T[i + j]$  do
5      $j \leftarrow j + 1$ 
6   end
7   if  $j == m$  then
8     output "Pattern found at index" +  $i$ 
9   end
10 end
```

The best case scenario of this process is when $P_{i=0}$ is not present in T , yielding an $O(n)$ number of comparisons, making it particularly advantageous for solving problems with small n and m values. When larger sets are used however, the algorithm’s efficiency significantly deteriorates. The worst case scenario for instance would be when all the characters for P and T are the same ($P_{i=0} \dots P_{m-1} = T_{j=0} \dots T_{n-1}$) or when all the characters for both P and T are the same except for their last characters ($(P_{i=0} \dots P_{m-2} =$

$T_{j=0} \dots T_{n-2}) \wedge (P_{m-1} = T_{n-1}))$. This would yield an $O(m \cdot (n - m + 1))$ number of comparisons.

2.3 Knuth-Morris-Pratt Algorithm

[2] [9] This algorithm differs from the Naive approach by implementing an LPS table as part of the pre-processing phase. Its purpose is to keep track of the comparisons made; after a mismatch, it is used to calculate where to begin the next match without needing to reset P . The table is constructed by determining the largest prefix of P that matches its largest suffix – the main idea is to identify how many characters can be skipped, by not matching characters that have already been calculated to match (eliminating redundancy). Although appearing complicated, this can be better explained through pseudocode and a practical example represented through a truth table. Upon initialization, the LPS table is an array of 0s of length m . For demonstration purposes, P will equal "ABABA":

Index (i)	Pattern Character (P_i)	Length (of Suffix-Prefix match)	LPS Value ($LPS[i]$)
0	A	0	0
1	B	0	0
2	A	0	1
3	B	1	2
4	A	2	3

Table 2: Truth Table for Pattern "ABABA"

Algorithm 2: CalculateLPSArray

```
1  $length \leftarrow 0$  ;
2  $LPS \leftarrow m * [0]$  ;

3 while  $i < m$  do
4   if  $P[i] == P[m]$  then
5      $length \leftarrow length + 1$ ;
6      $LPS[i] \leftarrow 0$ ;
7      $i \leftarrow i + 1$ ;
8   end
9   else if  $length \neq 0$  then
10     $length \leftarrow LPS[length - 1]$ ;
11  end
12  else
13     $LPS[i] \leftarrow 0$ ;
14     $i \leftarrow i + 1$ ;
15  end
16 end
17 return  $LPS$ 
```

Afterwards, the main phase involves applying the calculated LPS table on T . The indexes i and j both begin at 0, and are incremented while P_i matches T_j . When a mismatch occurs, it becomes apparent that the characters in P match with the characters of T up to the mismatch ($P_{0...(i-1)} = T_{j-i...j-1}$). Moreover, from the LPS table, it is known that $LPS[i - 1]$ represents the length of the longest part of the prefix and suffix of P . With these pieces of information, it can be concluded that the characters in $P_{0...i-1}$ do not need to be checked with $T_{j-i...j-1}$ since it is already known that they match. Thus, these can be skipped in both P and T . The exact process is presented in pseudocode here below:

Algorithm 3: KMP Match

```
1  $n \leftarrow T.length()$  ;
2  $m \leftarrow P.length()$  ;
3  $LPS \leftarrow CalculateLPSArray(P)$  ;
4  $i \leftarrow 0$  ;
5  $j \leftarrow 0$  ;

6 while  $j < n$  do
7   if  $P[i] == T[j]$  then
8      $i \leftarrow i + 1$ ;
9      $j \leftarrow j + 1$ ;
10  end
11  if  $i == m$  then
12    output "Pattern found at index"  $+(j - i)$ ;
13     $i \leftarrow LPS[i - 1]$ ;
14  end
15  else if  $j < n$  and  $P[i] \neq T[j]$  then
16    if  $i \neq 0$  then
17       $i \leftarrow LPS[i - 1]$ ;
18    end
19    else
20       $j \leftarrow j + 1$ ;
21    end
22  end
23 end
```

By skipping the unnecessary character comparisons, the algorithm proves to be much more efficient than the Naive approach. This is seen in its worst-case time complexity, which is $O(n + m)$ (best-case time complexity remains the same). As for the time complexities in the pre-processing and matching phases, they are $O(m)$ and $O(n)$ respectively

2.4 Rabin-Karp Algorithm

[10] The most similar out of the two to the Naive approach, the Rabin-Karp algorithm has a similar process of checking every substring for a match. Unlike the Naive approach however, this algorithm compares “hash values”. These are calculated using rolling hash functions, which allow for the removal and introduction of characters to the calculation, eliminating the need to recalculate the value entirely from scratch. To perform the hashing process on a substring, the following steps are carried out:

Rolling Hash Function

1. The base (b) and modulus (p) values are chosen (usually prime numbers).
2. A hash value is initialised to 0 and the initial hash value of P calculated (for each character in P , its contribution is added to the hash value as $c \cdot (b^{m-i-1}) \% p$).
3. The hash value for the first substring of length m in T is calculated.
4. Everytime j is incremented, the contribution of the leftmost character is removed, the contribution of the rightmost character added.
5. Compare the text and pattern matches each time; if they match then it is a **potential match** (rolling function outputs can be the same for different inputs). To verify the match, a character match is performed.

The time complexity for the pre-processing phase (calculating P 's hash value) is $O(m)$. As for the time complexities in the matching phase, they are $O(m \cdot n)$ for the worst case scenario and $O(n + m)$ for the expected case. Although this algorithm maintains the simplicity of the Naive algorithm all the while improving its efficiency, it has a disadvantage that, depending on the specific T case, may impact its efficiency – a spurious hit. This refers to the case where a hash value match is made between P and T but the characters do not match, increasing the time complexity. This can however be minimized by implementing a robust hash function.

The pseudocode for the algorithm is presented here below:

Algorithm 4: Rabin Karp

```
1  $n \leftarrow T.length()$  ;
2  $m \leftarrow P.length()$  ;
3  $h \leftarrow 1$  ;
4  $hashT \leftarrow 0$  ;
5  $hashP \leftarrow 0$  ;

6 for  $i \leftarrow 1$  to  $m - 1$  do
7    $h \leftarrow (h \cdot b) \% p$  ;
8 end

9 for  $i \leftarrow 0$  to  $m - 1$  do
10   $hashP \leftarrow (b \cdot hashP + P[i]) \% p$  ;
11   $hashT \leftarrow (b \cdot hashT + T[i]) \% p$  ;
12 end

13 for  $i \leftarrow 0$  to  $n - m$  do
14   $match \leftarrow TRUE$  ;
15  for  $j \leftarrow 0$  to  $m - 1$  do
16    if  $P[j] \neq T[i + j]$  then
17       $match \leftarrow FALSE$  ;
18    end
19  end
20  if  $match == TRUE$  then
21    end
22    output "Pattern found at index " +  $i$ ;
23  if  $i < n - m$  then
24     $hashT \leftarrow (b \cdot hashT + T[i]) + T[i + m] \% p$  ;
25    if  $hashT < 0$  then
26       $hashT \leftarrow hashT + p$  ;
27    end
28  end
29 end
```

3 Experimental Methodology

3.1 Measuring Performance

To measure the detailed algorithms' performance, their mentioned pseudocode algorithms were programmed in *Python* 3.12. Python was chosen due to it's wide array of third-party modules, two of which are used as part of the implementation: *Time* and *Matplotlib*. *Time* is used to measure the algorithms' runtimes for each algorithm, and *Matplotlib* is

used to generate the graphical plots. Using their implementations, the following variables will be measured:

1. Dependent Variables

Time (ms) - The amount of time taken to search a pattern P in a text T .

Number of Comparisons (pattern) - The number of comparisons made between the pattern characters.

Number of Comparisons (text) - The number of references made to the text (until all instances of the pattern are found).

2. Control Variables

Language - Depending on whether the language employs a compiler or an interpreter, runtime speeds can differ. Moreover, the type of programming language (procedural, functional, etc.) can have an impact on how the code is run, further impacting runtime speed.

Pattern Length - By controlling this factor, it is ensured that variations in number of comparisons and runtime are due to the algorithms themselves. This

The algorithms are then evaluated on their performance based on a set of constructed criteria:

Criteria	Testing
Time complexity	A
Space complexity	B
Simplicity	A
Applicability	B

Table 3: Evaluation Criteria

3.2 Application to Domains

Plagiarism detection, DNA sequencing, Search Engine

4 Experimental Results Analysis

4.1 Advantages and Disadvantages

5 Conclusion

reference style

References