

String Matching Algorithms

Naive and KMP algorithms implementation and analysis

Elia Innocenti

November 2023

Index

1	Introduction	2
1.1	Hardware and software specifications	2
2	The string matching problem	2
2.1	Definition	2
2.2	Example	3
2.3	Practical application	3
3	The naive string-matching algorithm	4
3.1	Pseudocode	4
4	String matching with finite automata	5
5	Knuth-Morris-Pratt algorithm	5
5.1	Pseudocode	5
6	Tests and results	6
6.1	Expected performances	6
7	Conclusions	6

1 Introduction

In this report we would like to compare two algorithms for the **string-matching problem**, namely the **naive string-matching algorithm** and the **Knuth-Morris-Pratt algorithm**.

1.1 Hardware and software specifications

Below are the specifications of the machine used to conduct the tests:

- **CPU:** Apple M1 SoC (System on Chip) with an 8-core CPU (4 high-performance and 4 high-efficiency cores), 7-core GPU, and Neural Engine
- **RAM:** 8GB unified memory
- **OS:** Sonoma 14.0
- **Python version:** 3.11

2 The string matching problem

2.1 Definition

The **string-matching problem** is the problem of finding all occurrences of a pattern P in a text T .

We formalize the string-matching problem as follows.

We assume that the text is an array $T[1..n]$ of length n and that the pattern is an array $P[1..m]$ of length $m \leq n$. We further assume that the elements of T and P are characters drawn from a finite alphabet Σ . For example, we may have $\Sigma = \{0, 1\}$ or $\Sigma = \{a, b, c, \dots, z\}$, the set of lowercase English letters. The character arrays T and P are often called strings of character.

In fact, the string-matching problem can also be seen as the problem of finding all valid shifts with which a given pattern P occurs in a text T .

2.2 Example

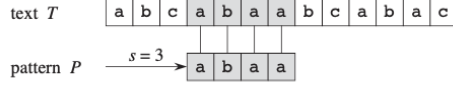


Fig. 1: An example of the string-matching problem. Given the text $T = abcabaabcbabac$ and the pattern $P = abaa$, we want to find all the occurrences of P in T . The pattern occurs only once in T , starting at shift $s = 3$.

$$\begin{cases} 0 \leq s \leq n - m \\ T[s+1..s+m] = P[1..m] \\ T[s+j] = P[j] \text{ for } (1 \leq j \leq m) \end{cases}$$

2.3 Practical application

String matching algorithms have many practical applications in different fields:

- Text search: Find the occurrences of words or phrases within documents or long texts.
- Natural language processing: Analysing text to recognise language patterns, lexical analysis and parsing of texts.
- Computational biology: Analysing DNA, RNA and protein sequences to identify biological patterns and correlations.
- Log and structured text analysis: Recognise patterns in large datasets, such as system logs, sensor records, etc.
- Data filtering and analysis: They allow information to be extracted from large volumes of structured and unstructured data.
- Search engines: Working behind the scenes to identify the most relevant matches between user search queries and content.
- Data compression: Used in compression algorithms to find and replace repeated patterns.
- Computer security: Identify signatures and patterns of attacks by detecting byte sequences or malicious behaviour in data.

3 The naive string-matching algorithm

The naive algorithm finds all valid shifts using a loop that checks the condition $P[1..m] = T[s + 1..s + m]$ for each of the $n - m + 1$ values of s .

3.1 Pseudocode

Data: Text T of size n and pattern P of size m

```
1 NAIVE-STRING-MATCHER(T,P)
2  $n \leftarrow T.length$ 
3  $m \leftarrow P.length$ 
4 for  $s \leftarrow 0$  to  $n - m$  do
5   | if  $P[1..m] = T[s + 1..s + m]$  then
6   |   | print "Pattern occurs with shift"  $s$ ;
7   | end
8 end
```

4 String matching with finite automata

5 Knuth-Morris-Pratt algorithm

5.1 Pseudocode

Data: Text T of size n and pattern P of size m

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

```

1 COMPUTE-PREFIX-FUNCTION(P)
2  $m \leftarrow P.length$ 
3 let  $\pi[1..m]$  be a new array;
4  $\pi[1] \leftarrow 0$ ;
5  $k \leftarrow 0$ ;
6 for  $q \leftarrow 2$  to  $m$  do
7   while  $k > 0$  and  $P[k + 1] \neq P[q]$  do
8      $k \leftarrow \pi[k]$ ;
9   end
10  if  $P[k + 1] = P[q]$  then
11     $k \leftarrow k + 1$ ;
12  end
13   $\pi[q] \leftarrow k$ ;
14 end
15 return  $\pi$ ;

```

6 Tests and results

6.1 Expected performances

7 Conclusions