

Scientific Computing (M3SC)

Peter J. Schmid

January 22, 2017

1 PATTERN SEARCHING (RABIN-KARP ALGORITHM)

Pattern searching is an important discipline of computer science and applied mathematics. It is concerned with the extraction of a given pattern p of size m from a (commonly very large) text t of length n . The pattern and text can consist of numbers, letters or other distinct structures or tokens. Typical applications of pattern searching are encountered in plagiarism detection software or in the sequencing of genetic material.

Here we will consider sequential pattern searching, where a large string of text is searched for occurrences of shorter patterns. For example, in a long sequence of nucleotides within a DNA molecule we may be interested in searching for locations and frequency of a short pattern made up of the DNA alphabet $\{A, C, G, T\}$ which denotes the four building blocks, i.e., the nucleic acids Adenine, Cytosine, Guanine and Thymine.

1.1 NAIVE SEARCH

A naive approach to pattern searching is the sequential comparison of the pattern to text snippets of equal length, then sliding along the entire text sequence. While this naive pattern searching algorithm often works better than expected, it exhibits its worst-case performance (and run times) when looking for rather long patterns in very long search strings. With a text of length n and a pattern of length m , we have to make a comparison of m entries in the pattern and substring for each of the positions in the text t . The run time of the algorithm thus scales as $\mathcal{O}(nm)$.

An implementation of a naive search algorithm is given below, written as a python function.

```
1 def naive_search (string, pattern):
```

```

2     # naive string matching algorithm
3     n = len(string)
4     m = len(pattern)
5     ic = 0
6     # element-by-element matching
7     for i in range(n-m+1):
8         # call the (logical) match function
9         if match(string, pattern, i):
10             ic += 1
11             print i
12     print "number of matches: ", ic

```

This function uses another python function `match` that we will use later in the more advanced Rabin-Karp algorithm (see below). It is given by

```

1 def match(string, pattern, i):
2     # element-by-element matching of a pattern
3     # in a string; i: position in string
4     for k in range(len(pattern)):
5         if pattern[k] == string[i+k]:
6             pass
7         else:
8             return False
9     return True

```

1.2 SPEEDING UP THE NAIVE-SEARCH ALGORITHM BY HASHING

One possibility of speeding up the naive-search algorithm is by applying more sophisticated manners of traversing and skipping through the text, rather than the slow sliding-window approach. Another and alternative possibility is the speed-up of the matching procedure. This latter technique is what the Rabin-Karp algorithm proposes. The Rabin-Karp algorithm, named after Michael Rabin and Richard Karp, speeds up the matching section of the naive-search algorithm by introducing **hash functions**.

In general, a hash function is a mapping of a string or number sequence onto a numerical value, the **hash value**. Based on this hash function, the search pattern and a string segment match when their respective hash values match. This hash value is akin to the familiar checksum, when comparing two data files. The Rabin-Karp algorithm thus replaces a direct comparison of a given pattern and an extracted string segment by a comparison of their associated hash values. In other words, in the text sequence we search for substrings that have the same hash value as the pattern's hash value.

When introducing hash functions into the naive-search algorithm, two issues arise.

First, the mapping of a pattern onto a hash value may produce false positives, i.e., the hash values of the pattern and the text substring coincide, even though the actual substrings do not match. This is also referred to as **hash collision**. It is a consequence of the fact that

different text-substrings can map to the same hash value. While this non-uniqueness does not pose a problem, we simply need to compare the pattern and the real text segment after we receive a flag from the matching hash values. In other words, hash value matching does a fast preprocessing step, before direct comparison confirms or rejects the potential match. Since a direct comparison is relatively costly, we have to choose a proper hash function that does not produce too many false positives but rather keeps direct comparisons a rare event.

Second, the computing of the hash value has to be efficient, otherwise there will not be any savings over a direct comparison. Ideally, the comparison should not scale with the pattern size, but instead be constant in time. In this latter case, the run time $\mathcal{O}(nm)$ of the naive search would reduce to $\mathcal{O}(n)$, which is as good as one can expect from a sequential search algorithm.

With the hash function introduced, the Rabin-Karp string searching algorithm calculates a hash value for the pattern, and for each m -character subsequence of the text. If the hash values are unequal, the algorithm will calculate the hash value for the next m -character sequence. If the hash values are equal, the algorithm will compare the pattern and the m -character sequence directly. In this way, there is only one hash comparison per text subsequence, and character matching is only needed when hash values match.

1.3 A BIT OF MATHEMATICS

Consider an m -character sequence as an m -digit number expressed in base b , where b is at least the number of letters in the alphabet used in our text t . The text subsequence $t[i, \dots, i + m - 1]$ can then be mapped to the number x_i according to

$$x_i = t[i] \cdot b^{m-1} + t[i+1] \cdot b^{m-2} + \dots + t[i+m-1]. \quad (1.1)$$

This expression already acts as a hash function, assigning a number x_i to a string $t[i, \dots, i + m - 1]$. If the text to be processed consists of letters (such as in plagiarism detection, or in DNA pattern searching), an additional mapping from the letters to the $t[i]$'s has to be introduced (see the example below). The above way of expressing a text sequence as a number using a polynomial in b is known as the **Rabin fingerprint**. It has the considerable advantage that the subsequent hash value x_{i+1} of the right-shifted m -digit subsequence $t[i+1, \dots, i+m]$ can be computed in constant time. This can be easily shown by

$$x_{i+1} = t[i+1] \cdot b^{m-1} + t[i+2] \cdot b^{m-2} + \dots + t[i+m], \quad (1.2a)$$

$$= x_i \cdot b \quad (\text{shift left one digit}) \quad (1.2b)$$

$$- t[i] \cdot b^m \quad (\text{subtract leftmost digit}) \quad (1.2c)$$

$$+ t[i+m]. \quad (\text{add new rightmost digit}) \quad (1.2d)$$

In this manner, we avoid the explicit computation of a new hash value; rather, we simply adjust the current hash value as we shift our substring in the text by one character to the right. Hash functions that allow this type of shortcut are referred to as **rolling hash** functions; there are many others besides the Rabin fingerprint.

In the above expression, there still remains one problem. If the size m of the pattern gets large and the alphabet of the text is rather rich (and thus requires a large base b), the value of b^m that appears in the expression above will become exceedingly high. To counteract this trend, we limit the size of the hash value by introducing the modulo function, i.e., we clip the hash values to a finite number range. More specifically, we compute our hash values by taking the modulo (with an appropriate number q) of each component in the above expression.

The modulo function `mod` is particularly convenient since it satisfies the mathematical expressions

$$[(x \bmod q) + (y \bmod q)] \bmod q = (x + y) \bmod q, \quad (1.3a)$$

$$(x \bmod q) \bmod q = x \bmod q. \quad (1.3b)$$

Applied to our rolling hash expression, and introducing the notation $h_i = x_i \bmod q$, we obtain

$$\begin{aligned} h_{i+1} &= (t[i+1] \cdot b^{m-1} \bmod q + t[i+2] \cdot b^{m-2} \bmod q + \dots + t[i+m] \bmod q) \bmod q \\ &= (h_i \cdot b \bmod q \quad \text{(shift left one digit)} \end{aligned} \quad (1.4a)$$

$$- t[i] \cdot b^m \bmod q \quad \text{(subtract leftmost digit)} \quad (1.4b)$$

$$+ t[i+m] \bmod q) \bmod q \quad \text{(add new rightmost digit)} \quad (1.4c)$$

This latter expression is at the heart of the Rabin-Karp algorithm.

It leaves the final question of what value to choose for q . A small number will increase the frequency of hash collisions, the necessity to perform an excessive number of direct comparisons and, consequently, a serious degradation in performance. For example, when choosing $q = 3$, we only produce hash values of $h_i = 0, 1, 2$, which may be far too limited to properly represent all the different text substrings. Instead, a large prime number is commonly chosen for q .

1.4 ALGORITHM

With this we can state the Rabin-Karp algorithm.

The Rabin-Karp pattern search algorithm computes a hash value for a search pattern p of length m and for each m -character substring of a text t . The two hash values are compared, rather than the actual strings. Only after the two values match does the algorithm compare the two patterns (search pattern and substring) to confirm or reject a match, thus resolving a hash collision. If no match is found, the algorithm shifts to the next substring of t . When a rolling hash is used, this latter shift can be computed very efficiently.

We have the following python functions, implementing the Rabin-Karp string search algorithm.

```
1 import sys
2 import time
```

```

3
4 def RabinKarp(str, pat, base=256, modu=16647133):
5     # Rabin-Karp string matching algorithm
6     #     str: text to be searched
7     #     pat: pattern to match
8     #     base: base of the Rabin-fingerprint
9     #     modu: modulus of the hash-function
10    ic = 0
11    n = len(str)
12    m = len(pat)
13    # compute the hash values for pattern
14    # and initial string (Horner's scheme)
15    hp = 0
16    for i in pat:
17        hp = (base*hp + ord(i))%modu
18    hs = 0
19    for i in str[:m]:
20        hs = (base*hs + ord(i))%modu
21    # check for initial match
22    if hs == hp and match(str[:m], pat):
23        print 0
24        ic = 1
25    # compute b^m
26    bm = 1
27    for i in range(m-1):
28        bm = (bm*base)%modu
29    # main loop
30    for i in range(1, n - m + 1):
31        hs = (base*(hs - bm*ord(str[i-1])) + \
32              ord(str[i+m-1]))%modu
33        if (hs == hp):
34            if match(str[i:i+m], pat):
35                print i,
36                ic += 1
37    print "\n number of matches found: ", ic

```

Finally, the main code contains material that applies the Rabin-Karp algorithm and the naive search algorithm to two applications: a data-base of genetic material from the DNA of the fruit fly (*Drosophila melanogaster*) which contains more than 23 million nucleotide, and the first 100000 digits of π .

```

1 if __name__ == '__main__':
2
3     # read 'Genetic Material' file
4     f = open('GeneticMaterialLarge','r')
5     strGM = f.read()
6     f.close()
7     strGM = strGM.replace("\n","")
8
9     # patterns to search for
10    p1 = 'CACAATATATGATCGC'
11    p2 = 'GTGCCAACATATTGTGCTCTCTATATAATGACTGCCTCT'
12    p3 = 'GCGAATATGGAAGGAGCAGACACACACTTGTGATTTCTTTATGCG'
13    pa = 'ATTTCTACTGTAAAAATCTTTCTTTGTCCAAAAATCTTTGATAGTGTTT'
14    pb = 'CCCTATGTTTGTGTGCACGACACGTGCTAAAGCCAAATATTTTGCTAGT'
15    pc = 'CAATATATAACAAAAATATTTTCGGGCTTGCGTACTGTAAAGATGAAT'
16    p4 = pa + pb + pc
17    patGM = p1
18
19    # read 'Digits of Pi' file
20    f = open('DigitsOfPi','r')
21    strPi = f.read()
22    f.close()
23
24    # pattern to search for
25    p1 = '8888'
26    p2 = '2479114016957900338356'
27    patPi = p1
28
29    # scan the Genetic Material file
30    t0 = time.time()
31    RabinKarp(strGM,patGM)
32    t1 = time.time()
33    total = t1-t0
34    print 'Genetic Material (Rabin-Karp algorithm): ',total
35
36    t0 = time.time()
37    naiveSearch(strGM,patGM)
38    t1 = time.time()
39    total = t1-t0
40    print 'Genetic Material (naive search): ',total
41
42    # scan the Digits of Pi file
43    t0 = time.time()

```

```

44  RabinKarp(strPi,patPi)
45  t1      = time.time()
46  total  = t1-t0
47  print 'Digits of Pi (Rabin-Karp algorithm): ',total
48
49  t0      = time.time()
50  naiveSearch(strPi,patPi)
51  t1      = time.time()
52  total  = t1-t0
53  print 'Digits of Pi (naive search): ',total

```

The main code includes a time comparison between the naive search and the Rabin-Karp algorithm. Even though the difference may not appear this drastic, one has to keep in mind that the algorithm is commonly applied to far larger data-bases and to larger patterns. In this case, the difference in run-time between naive searching and Rabin-Karp will be more pronounced.

1.5 IMPROVEMENTS AND APPLICATIONS

The Rabin-Karp algorithm has been presented in its most basic form. However, even in this rudimentary form it is commonly applied in bioinformatics (when searching for patterns in nucleotides), text processing (such as plagiarism detection), music prototyping (for search engines) or data compression.

Nonetheless, there are improvements to the algorithms that increase its efficiency and applicability to a wide range of data. One of the most common is the use of prefix and suffix trees. Prefix and suffix strings are smaller substrings attached in front or at the end of a string that allow the early dismissal of mismatching strings and the skipping ahead by more than one record. Taking advantage of prefix and suffix strings can produce substantial speedup in the overall search for patterns.

The Rabin-Karp algorithm can quite easily be extended to multiple-pattern searches. In this case, we are looking for occurrences in the form of a vector-valued shift (s_1, s_2) (in the case of two patterns). Facial recognition and medical imaging are areas where two-dimensional pattern searches are commonly encountered.

Another extension of the Rabin-Karp algorithm involves the occurrence of wild cards, which are taken as matching *any* pattern. This situation is often encountered when dealing with DNA transcription factors. These factors need to be ignored (labelled as wild cards) when searching for valid patterns.

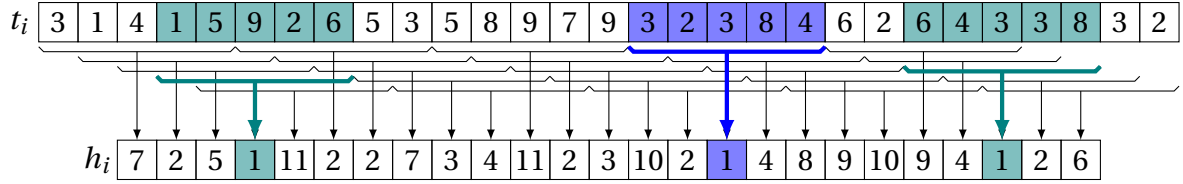


Figure 1.1: Demonstration of the Rabin-Karp algorithm: search for the string '32384' in the digits of π . We employ the algorithm with radix 10 (i.e., a representation of the digits in base 10) and a modulo 13 function. The hash value of the pattern is 1. We recognize the string in position 16 in the text (highlighted in blue), but we also encounter two hash collisions (in green) in positions 4 and 23, where the hash value matches, but not the actual substring.

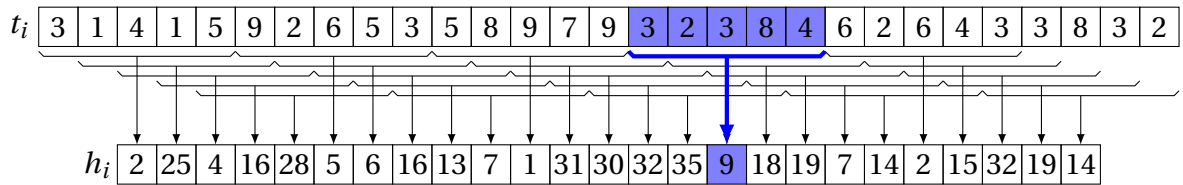


Figure 1.2: Demonstration of the Rabin-Karp algorithm: search for the same string '32384' in the digits of π . This time, we choose a modulo 37 function. The hash value of the pattern is 9, and we do not encounter hash collisions.