

Strings

Algoritmo de Rabin-Karp

Prof. Edson Alves - UnB/FGA

2019

1. Algoritmo de Rabin-Karp
2. Variantes do algoritmo de Rabin-Karp

Algoritmo de Rabin-Karp

Definição

- O algoritmo de Rabin-Karp é um algoritmo que contabiliza o número de ocorrências da string P , de tamanho m , na string S , de tamanho n
- Ele foi proposto por Michael O. Rabin e Richard M. Karp em 1987
- A ideia principal do algoritmo é computar o *hash* $h_P = h(P)$ e compará-lo com todas as substrings $h_{ij} = S[i..j]$ de S de tamanho m
- Caso $h_P \neq h_{ij}$, segue que $P \neq S[i..j]$ e o algoritmo pode prosseguir
- Se $h_P = h_{ij}$, as strings ou são iguais ou houve uma colisão
- Esta dúvida pode ser sanada através da comparação direta, enquanto strings, entre $S[i..j]$ e P
- O algoritmo tem complexidade $O(mn)$ no pior caso, por conta do custo do cálculo dos *hashes* e das possíveis comparações diretas entre as strings

Pseudocódigo do algoritmo de Rabin-Karp

Algoritmo 1 Algoritmo de Rabin-Karp – Naive

Input: Duas strings P e S

Output: O número de ocorrências occ de P em S

```
1: function RABINKARP( $P, S$ )
2:    $m \leftarrow |P|$ 
3:    $n \leftarrow |S|$ 
4:    $occ \leftarrow 0$ 
5:    $h_P \leftarrow h(P)$ 
6:   for  $i \leftarrow 1$  to  $n - m + 1$  do
7:      $h_S \leftarrow h(S[i..(i + m - 1)])$ 
8:     if  $h_S = h_P$  then
9:       if  $S[i..(i + m - 1)] = P$  then
10:         $occ \leftarrow occ + 1$ 
11:   return  $occ$ 
```

Implementação do algoritmo de Rabin-Karp em Haskell

```
1 import Data.Char
2
3 f :: Char -> Int
4 f c = (ord c) - (ord 'a') + 1
5
6 h :: String -> Int
7 h s = sum (zipWith (*) fs ps) `mod` m where
8     p = 31
9     m = 10^9 + 7
10    fs = map f s
11    ps = map (\x -> p ^ x) $ take (length s) [0..]
12
13 rabin_karp :: String -> String -> Int
14 rabin_karp s p = sum rs where
15     n = length s
16     m = length p
17     hp = h p
18     xss = [take m (drop i s) | i <- [0..(n - m)]]
19     rs = [fromEnum (h xs == hp && xs == p) | xs <- xss]
20
21 main = print $ rabin_karp "abababababab" "aba"
```

Implementação do algoritmo de Rabin-Karp em C++

```
1 #include <bits/stdc++.h>
2
3 int f(char c)
4 {
5     return c - 'a' + 1;
6 }
7
8 int h(const std::string& s)
9 {
10     long long ans = 0, p = 31, m = 1000000007;
11
12     for (auto it = s.rbegin(); it != s.rend(); ++it)
13     {
14         ans = (ans * p) % m;
15         ans = (ans + f(*it)) % m;
16     }
17
18     return ans;
19 }
20
```

Implementação do algoritmo de Rabin-Karp em C++

```
21 int rabin_karp(const std::string& s, const std::string& p)
22 {
23     int n = s.size(), m = p.size(), occ = 0, hp = h(p);
24
25     for (int i = 0; i < n - m; i++)
26     {
27         auto b = s.substr(i, m);
28         occ += (h(b) == hp && b == p) ? 1 : 0;
29     }
30
31     return occ;
32 }
33
34 int main()
35 {
36     auto s = "abababababab", p = "aba";
37
38     std::cout << rabin_karp(s, p) << '\n';
39
40     return 0;
41 }
```


Pseudocódigo do algoritmo de Rabin-Karp

Algoritmo 2 Algoritmo de Rabin-Karp – Naive

Input: Duas strings P e S

Output: O número de ocorrências occ de P em S

```
1: function RABINKARP( $P, S$ )
2:    $m \leftarrow |P|$ 
3:    $n \leftarrow |S|$ 
4:    $occ \leftarrow 0$ 
5:    $h_P \leftarrow h(P)$ 
6:   for  $i \leftarrow 1$  to  $n - m + 1$  do
7:      $h_S \leftarrow h(S[i..(i + m - 1)])$ 
8:     if  $h_S = h_P$  then
9:       if  $S[i..(i + m - 1)] = P$  then
10:         $occ \leftarrow occ + 1$ 
11:   return  $occ$ 
```

Variantes do algoritmo de Rabin-Karp

Diminuição da complexidade para o cálculo dos *hashes*

- Da maneira como foi apresentada, o algoritmo de Rabin-Karp tem complexidade $O(mn)$ no pior caso, o mesmo da busca completa, e com *runtime* maior, por conta do cálculo dos *hashes*
- Uma primeira melhoria que pode ser feita é usar o *rolling hash*, e computar $h(S[(i+1)..(i+m)])$ a partir de $h(S[i..(i+m-1)])$ com custo $O(1)$
- Isto é possível, pois

$$\begin{aligned}h(S[(i+1)..(i+m)]) &= S_{i+1} + S_{i+2}p + \dots + S_{i+m}p^{m-1} \\&= \frac{S_i + S_{i+1}p + \dots + S_{i+m-1}p^{n-1} + S_{i+m}p^m - S_i}{p} \\&= \frac{S_i + S_{i+1}p + \dots + S_{i+m-1}p^{n-1} - S_i}{p} + S_{i+m}p^{m-1} \\&= \frac{h(S[i..(i+m-1)]) - S_i}{p} + S_{i+m}p^{m-1}\end{aligned}$$

Diminuição da complexidade para o cálculo dos *hashes*

- Como $S_i < p$, para todo i , então

$$h(S[(i+1)..(i+m)]) = \left\lfloor \frac{h(S[i..(i+m-1)])}{p} \right\rfloor + S_{i+m}p^{m-1}$$

- Se a constante $k \equiv p^{m-1} \pmod{m}$ for precomputado, cada atualização do *hash* tem custo $O(1)$
- Observe que a divisão deve ser feita por meio da multiplicação pelo inverso multiplicativo de p módulo m
- Este inverso também pode ser precomputado, como no caso da constante k
- O pior caso ainda tem complexidade $O(nm)$, mas o caso médio passa a ter complexidade $O(n+m)$

Pseudocódigo do algoritmo de Rabin-Karp

Algoritmo 3 Algoritmo de Rabin-Karp – Naive

Input: Duas strings P e S

Output: O número de ocorrências occ de P em S

```
1: function RABINKARP( $P, S$ )
2:    $m \leftarrow |P|$ 
3:    $n \leftarrow |S|$ 
4:    $occ \leftarrow 0$ 
5:    $h_P \leftarrow h(P)$ 
6:   for  $i \leftarrow 1$  to  $n - m + 1$  do
7:      $h_S \leftarrow h(S[i..(i + m - 1)])$ 
8:     if  $h_S = h_P$  then
9:       if  $S[i..(i + m - 1)] = P$  then
10:         $occ \leftarrow occ + 1$ 
11:   return  $occ$ 
```

1. CP-Algorithms. [String Hashing](#), acesso em 06/08/2019.
2. **CROCHEMORE**, Maxime; **RYTTER**, Wojciech. *Jewels of Stringology: Text Algorithms*, WSPC, 2002.
3. **HALIM**, Steve; **HALIM**, Felix. *Competitive Programming 3*, Lulu, 2013.
4. Wikipédia. [Rabin-Karp algorithm](#), acesso em 08/08/2019.