

BIF CM1

Exact pattern matching without index

Téo Imane, Claire Lemaitre, Catherine Belleannée

10 september 2020

Structure

1. Exact pattern matching
 - 1.1 Introduction
 - 1.2 Definitions
 - 1.3 Problem
 - 1.4 Pattern matching without index
 - 1.4.1 Naive algorithm
 - 1.4.2 Rabin-Karp algorithm
 - 1.4.3 Multiple pattern
 - 1.5 Pattern matching with index
 - 1.5.1 Suffix trees
 - 1.5.2 Suffix array
 - 1.5.3 Burrows-Wheeler Transform
 - 1.5.4 FM-index
2. Grammar (Catherine Belleannée)
3. Sequence comparison: dynamic programming, mapping, assembly (Claire Lemaître)

Pattern Matching

PM - What

- ▶ Are you here?
- ▶ Where are you?

PM - Why

- ▶ CTRL-F
- ▶ Search engines (50 billion web pages)
- ▶ DB requests
- ▶ Music (youtube, shazam)
- ▶ Biology (Zettabytes (ZB, 10^{21}) from 2020)

Sub problems

Distinct PM

- ▶ **Exact** : Search *ACTG* in *ACGCTAACGGACGCA*
- ▶ **Approximate**: Search *ACTG* in *ACGCTAACGGACGCA* with at most x substitutions, insertions and deletions

PM - two approaches

- ▶ **On the fly**: reference non indexed
- ▶ **Indexed**: reference pre-treated

Definitions (I)

- ▶ An alphabet Σ is a finite non-empty set of symbols. We denote $|\Sigma|$ the alphabet cardinality.
- ▶ Examples:
 - ▶ Unary alphabet: contains only one symbol
 - ▶ Binary alphabet: two symbols
 - ▶ In our case $\Sigma = \{A, C, G, T\}$, $|\Sigma| = 4$
- ▶ A string over an alphabet Σ is a finite ordered sequence of symbols from Σ .
- ▶ A set of all strings over a given alphabet Σ is denoted Σ^* . For example $AACCC \in \{A, C\}^*$ or $\{A, C, G, T\}^*$.
- ▶ A set of all strings of size n over a given alphabet Σ is denoted Σ^n .
- ▶ A length of a string **S** is denoted $|\mathbf{S}|$.
- ▶ The empty string is denoted ϵ and can be over any alphabet.
 $|\epsilon| = 0$

Definitions (II)

- ▶ For any string S :
 - ▶ $S[i]$ denotes the i^{th} symbol.
 - ▶ $S[i, j]$, with $i, j \in \mathbb{N} \cup \{0\}$, is the contiguous substring of S that starts at position i and ends at position j .
 - ▶ $S[i, j]$ is the empty string if $i > j$.
 - ▶ $S[0, i]$ and $S[i, |S| - 1]$ are respectively a **prefix** and a **suffix** of string S .

Definitions (III)

- ▶ Let $S = aagcgccgaa$ and $a, c, g \in \Sigma$.
- ▶ Let \mathbf{p} and \mathbf{s} a prefix and a suffix of S .
- ▶ If $|\mathbf{p}| < |S|$, then \mathbf{p} is a **proper prefix** of S .
- ▶ If $|\mathbf{s}| < |S|$, then \mathbf{s} is a **proper suffix** of S .
- ▶ $aagc$ is a proper prefix of S . $aagcgccgaa$ is a prefix of S but not a proper prefix.
- ▶ A **border** of S is a substring which is both a proper prefix and a proper suffix of S .
- ▶ **The** border of S is the longest border of S . It's denoted $Border(S)$.
- ▶ a is a border of S , aa is **the** border of S .

Exact pattern matching problem

Problem

Let T a string. $T \in \Sigma^*$.

Let P a searched pattern. $P \in \Sigma^*$.

Let $|P| < |T|$.

Let O the set of the starting positions of P in T .

Find $O = \{i_1, \dots, i_n\}$ such as $\forall i \in O, T[i, (i + |P| - 1)] = P$.

Solutions

- ▶ Naive algorithm.
- ▶ Rabin-Karp algorithm.
- ▶ (Knuth-Morris-Pratt algorithm).
- ▶ (Boyer-Moore algorithm).

Naive algorithm

```

Data:  $P \in \Sigma^*, T \in \Sigma^*$ 
1  $m = \text{length}(P);$ 
2  $n = \text{length}(T);$ 
3  $\text{list } match = \text{emptyList};$ 
4 for  $i = 0; i < n - m + 1; i = i + 1$  do
5      $j = 0;$ 
6     while  $j < m$  do
7         if  $T[i + j] \neq P[j]$  then
8              $\text{break}$ 
9         end
10         $j = j + 1;$ 
11    end
12    if  $j == m$  then
13         $\text{print } i;$ 
14         $match.\text{insert}(i);$ 
15    end
16 end
17 return  $match$ 

```

Naive algorithm complexity

```

Data:  $P \in \Sigma^*, T \in \Sigma^*$ 
1  $m = \text{length}(P);$ 
2  $n = \text{length}(T);$ 
3  $\text{list } match = \text{emptyList};$ 
4 for  $i = 0; i < n - m + 1; i = i + 1$  do
5      $j = 0;$ 
6     while  $j < m$  do
7         if  $T[i + j] \neq P[j]$  then
8              $\text{break}$ 
9         end
10         $j = j + 1;$ 
11    end
12    if  $j == m$  then
13         $\text{print } i;$ 
14         $match.\text{insert}(i);$ 
15    end
16 end
17 return  $match$ 

```

- To perform all the necessary comparisons we need to visit $n - m$ positions of T (the **for** loop).
- For each position, we compare $T[i + j]$ and $P[j]$ (the **while** loop). This corresponds to m operations ($|P|$).
- Time complexity : $O(m(n - m))$ that can be simplified such as $O(mn)$.
- Space complexity : $O(1)$.

Rabin-Karp algorithm

- ▶ Introduced by Richard M. Karp and Michael O. Rabin :
R. M. Karp and M. O. Rabin, "Efficient randomized pattern-matching algorithms," in IBM Journal of Research and Development, vol. 31, no. 2, pp. 249-260, March 1987, doi: 10.1147/rd.312.0249.
- ▶ We no longer compare each symbol of a sub-string to pattern symbols but sub-string hashes to a pattern hash.
- ▶ Avoids a large number of symbol comparisons.

Rabin-Karp algorithm

Data: $P \in \Sigma^*$, $T \in \Sigma^*$

```

1   $m = \text{length}(P)$ ;
2   $n = \text{length}(T)$ ;
3   $\text{list match} = \text{emptyList}$ ;
4   $\text{hash\_pattern} = \text{hash}(P)$ ;
5  for  $i = 0$ ;  $i < n - m + 1$ ;  $i = i + 1$  do
6       $\text{hash\_text} = \text{hash}(T[i, i + m - 1])$ ;
7      if  $\text{hash\_pattern} == \text{hash\_text}$  then
8           $j = 0$ ;
9          while  $j < m$  do
10             if  $T[i + j] \neq P[j]$  then
11                  $\text{break}$ ;
12             end
13              $j = j + 1$ ;
14         end
15         if  $j == m$  then
16              $\text{print } i$ ;
17              $\text{match.insert}(i)$ ;
18         end
19     end
20 end
21 return match

```

Rabin-Karp complexity

Data: $P \in \Sigma^*, T \in \Sigma^*$

```

1   $m = \text{length}(P);$ 
2   $n = \text{length}(T);$ 
3   $\text{list } match = \text{emptyList};$ 
4   $\text{hash\_pattern} = \text{hash}(P);$ 
5  for  $i = 0; i < n - m + 1; i = i + 1$  do
6       $\text{hash\_text} = \text{hash}(T[i, i + m - 1]);$ 
7      if  $\text{hash\_pattern} == \text{hash\_text}$  then
8           $j = 0;$ 
9          while  $j < m$  do
10             if  $T[i + j] \neq P[j]$  then
11                  $\text{break};$ 
12             end
13              $j = j + 1;$ 
14         end
15         if  $j == m$  then
16              $\text{print } i;$ 
17              $\text{match.insert}(i);$ 
18         end
19     end
20 end
21 return  $match$ 
```

- The hashes comparison is $O(1)$ but is executed for each position of T so its impact is $O(n)$ (line 7).
- Lines 4, 6, 10 are $O(m)$. But line 4 is executed once and 10 is executed if the hashes match: rarely with a good hash function.

Rolling hash

A rolling hash function allows to compute, for a text T , $\text{hash}(T[i, i + m - 1])$ from $\text{hash}(T[i - 1, i + m - 2])$.

Example:

Let $\Sigma = \{A, C, G, T\}$

Let $h(A) = 0, h(C) = 1, h(G) = 2, h(T) = 3$

Let $T = ACGTAT$

Let $\text{hash}(T[i, i + m - 1]) = \sum_{n=i}^{n \leq i+m-1} h(T[n])$

For $|P| = m = 3$:

$$\blacktriangleright \text{hash}(T[0, 2]) = h(A) + h(C) + h(G)$$

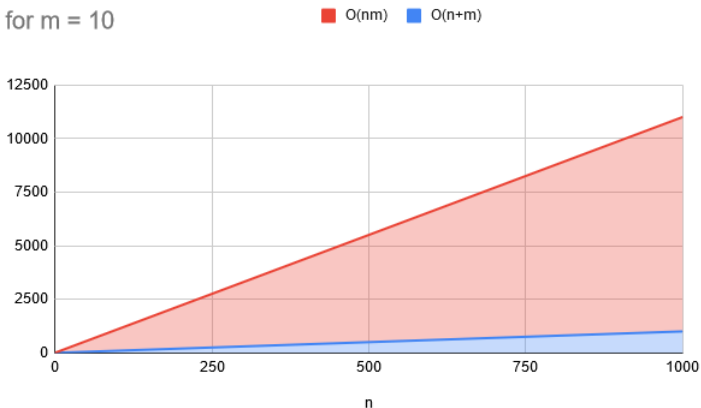
$$\blacktriangleright \text{hash}(T[1, 3]) = h(C) + h(G) + h(A) \Leftrightarrow \text{hash}(T[1, 3]) = \text{hash}(T[0, 2]) - h(T[0]) + h(T[3])$$

$$\text{hash}(T[i + 1, i + m - 1]) = \text{hash}(T[i, i + m - 2]) - h(T[i]) + h(T[i + m - 1])$$

Rabin-Karp and Rolling hash

With a rolling hash function, Rabin-Karp algorithm is average $O(n + m)$, worst $O(nm)$.

for $m = 10$



Rabin fingerprint

Compute fingerprint

For a string S and a prime p :

$$\text{hash}(S) = (S[0] \times p^0) + (S[1] \times p^1) + \dots + (S[|S| - 1] \times p^{|S|-1})$$

$$\text{hash}(S[i, i + m - 1]) = (S[i] \times p^0) + (S[i + 1] \times p^1) + \dots + (S[i + m - 1] \times p^{m-1})$$

Example using ASCII values as hash values

- ▶ Let $S = \text{ACGCCGCGG}$, $m = 3$, $p = 11$.
- ▶ $h(x) = \text{ascii}(x)$
- ▶ $H_0 = \text{hash}(S[0, 0 + 3 - 1]) = (h(A) \times 11^0) + (h(C) \times 11^1) + (h(G) \times 11^2)$
- ▶ $H_0 = \text{hash}(S[0, 0 + 3 - 1]) = (65 \times 11^0) + (67 \times 11^1) + (71 \times 11^2)$
- ▶ $H_1 = \text{hash}(S[1, 1 + 3 - 1]) = (67 \times 11^0) + (71 \times 11^1) + (67 \times 11^2)$
- $\Leftrightarrow H_1 = (65 \times 11^0) + \frac{(67 \times 11^1)}{11} + \frac{(71 \times 11^2)}{11} - 65 + (67 \times 11^2)$

Rabin fingerprint

Example using ASCII values as hash values

- ▶ Let $S = ACGCCGCGG$, $m = 3$, $p = 11$.
- ▶ $h(x) = \text{ascii}(x)$
- ▶ $H_0 = \text{hash}(S[0, 0 + 3 - 1]) = (h(A) \times 11^0) + (h(C) \times 11^1) + (h(G) \times 11^2)$
- ▶ $H_0 = \text{hash}(S[0, 0 + 3 - 1]) = (65 \times 11^0) + (67 \times 11^1) + (71 \times 11^2)$
- ▶ $H_1 = \text{hash}(S[1, 1 + 3 - 1]) = (67 \times 11^0) + (71 \times 11^1) + (67 \times 11^2)$
- $\Leftrightarrow H_1 = (65 \times 11^0) + \frac{(67 \times 11^1)}{11} + \frac{(71 \times 11^2)}{11} - 65 + (67 \times 11^2)$

$$\text{hash}(S[(i + 1), (i + m)]) = \text{hash}(S[i, i + m - 1]) - h(S[i]) + h(S[i + m]) \times p^{m-1}$$

Rabin-Karp and Rabin fingerprint

Data: $P \in \Sigma^*$, $T \in \Sigma^*$, $p \in \mathcal{P}$

```

1   $m = \text{length}(P)$ ;
2   $n = \text{length}(T)$ ;
3   $\text{list match} = \text{emptyList}$ ;
4   $\text{hash\_pattern} = \text{hash}(P, p)$ ;
5   $\text{hash\_text} = \text{hash}(T[0, 0 + m - 1], p)$ ;
6  for  $i = 0$ ;  $i < n - m + 1$ ;  $i = i + 1$  do
7      if  $\text{hash\_pattern} == \text{hash\_text}$  then
8           $j = 0$ ;
9          while  $j < m$  do
10             if  $T[i + j] \neq P[j]$  then
11                  $\text{break}$ ;
12             end
13              $j = j + 1$ ;
14         end
15         if  $j == m$  then
16              $\text{print } i$ ;
17              $\text{match.insert}(i)$ ;
18         end
19     end
20      $\text{hash\_text} = \text{next\_hash}(\text{hash\_text}, i, T[i, i + m - 1], p)$ ;
21 end
22 return match

```

Rabin fingerprint functions

```
1 Function hash (string, prime) :  
2    $h = 0$   
3   for  $m = 0$ ;  $m < \text{length}(\text{string})$ ;  $m = m + 1$  do  
4      $h = h + (\text{value}(\text{string}[m]) \times \text{prime}^m)$   
5   return  $h$   
6  
7 Function next_hash (previous, offset, string, prime) :  
8    $h = \text{previous} - \text{string}[\text{offset}]$   
9    $h = h / \text{prime}$   
10   $m = \text{length}(\text{string})$   
11   $h = h + (\text{string}[\text{offset} + m] \times \text{prime}^{m-1})$   
12  return  $h$   
13
```

Rabin-Karp for multiple pattern

For a single exact pattern matching, other algorithms which pre-process the pattern are in practice more efficient than Rabin-Karp, such as Boyer-Moore or Knuth-Morris-Pratt. But for multiple exact pattern matching, Rabin-Karp can be a very good choice.

Rabin-Karp for multiple pattern

```

Data:  $P = \{P_0, P_1, \dots, P_n\}, T \in \Sigma^*, p \in P$ 
1  $m = \text{length}(P);$ 
2  $n = \text{length}(T);$ 
3  $\text{list } match = \text{emptyList};$ 
4  $\text{set } hash\_patterns = \text{emptySet};$ 
5 foreach  $pattern \in P$  do
6    $hash\_patterns.insert(hash(pattern));$ 
7 end
8  $hash\_text = hash(T[0, 0 + m - 1], p);$ 
9 for  $i = 0; i < n - m + 1; i = i + 1$  do
10   if  $hash\_text \in hash\_patterns$  and  $T[i, i + m - 1] \in P$  then
11      $\text{print } i;$ 
12      $match.insert(i);$ 
13   end
14    $hash\_text = \text{next\_hash}(hash\_text, i, T[i, i + m - 1], p)$ 
15 end
16 return  $match$ 

```
