



String matching algorithms

Naive, Knuth-Morris-Pratt, Rabin-Karp



Типы алгоритмов

Алгоритмы поиска подстроки делятся на 2 типа алгоритмов/задач:

1. Точные алгоритмы поиска подстроки (Exact String Matching Algorithms)
2. Поиск похожих строк / Нечеткий поиск (Approximate String Matching Algorithms)

Поговорим о точных алгоритмах

Exact String Matching Algorithms

Задача: найти одно, несколько, или все совпадения определенной строки (шаблона) $\text{pat}[m]$ в длинной строке (текст и др.) $\text{txt}[n]$ ($n > m$) так, что совпадение будет точным.

Алфавит шаблона должен совпадать с алфавитом строк.

Naive algorithm

Идея: “накладываем” шаблон на текст и проверяем символы один за другим. Если первый символ шаблона и текста совпадут, то проверяем второй символ и т.д. Если не совпадут, сдвигаем шаблон на 1.

```
1 def NaiveSearch(txt, pat):
2     n = len(txt)
3     m = len(pat)
4     for i in range(n-m+1):
5         j = 0
6         while j < m:
7             if txt[i+j] != pat[j]:
8                 break
9             j += 1
10        if j == m:
11            return i
12    return -1
```

н	а		д	в	о	р	е		т	р	а	в	а	,		н	а		т	р	а	в	е		д	р	о	в	а
т	р	а	в	е																									
	т	р	а	в	е																								
		т	р	а	в	е																							
			т	р	а	в	е																						
				т	р	а	в	е																					
					т	р	а	в	е																				
						т	р	а	в	е																			
							т	р	а	в	е																		
								т	р	а	в	е																	
									т	р	а	в	е																
										т	р	а	в	е															
											т	р	а	в	е														
												т	р	а	в	е													
													т	р	а	в	е												
														т	р	а	в	е											
															т	р	а	в	е										
																т	р	а	в	е									
																	т	р	а	в	е								
																		т	р	а	в	е							
																			т	р	а	в	е						
																				т	р	а	в	е					

																		т	р	а	в	е								
																			т	р	а	в	е							
																			т	р	а	в	е							
																			т	р	а	в	е							
																			т	р	а	в	е							

Naive algorithm

Лучший случай: когда первый символ шаблона не присутствует в тексте

```
txt[] = "AABCCAADDEE";
```

```
pat[] = "FAA";
```

Время работы алгоритма в лучшем случае: $O(n)$.

Naive algorithm

Худший случай: 2 варианта:

1. Когда все символы в тексте и в шаблоне одинаковые

```
txt[] = "AAAAAAAAAAAAAAAAAAAA";
```

```
pat[] = "AAAAA";
```

2. Когда отличается только последний символ

```
txt[] = "AAAAAAAAAAAAAAAAAAAB";
```

```
pat[] = "AAAAB";
```

Время работы алгоритма в худшем случае: $O(m \cdot (n - m + 1))$

Naive algorithm

Однако если m достаточно мало по сравнению с n , то тогда асимптотика получается близкой к $O(n)$, поэтому этот алгоритм достаточно широко применяется на практике.

Naive algorithm

Преимущества

- Требуется $O(1)$ памяти.
- Приемлемое время работы на практике. Благодаря этому алгоритм применяется, например, в браузерах и текстовых редакторах (при использовании `Ctrl + F`), потому что обычно паттерн, который нужно найти, очень короткий по сравнению с самим текстом. Также примитивный алгоритм используется в стандартных библиотеках языков высокого уровня (C++, Java), потому что он не требует дополнительной памяти.
- Простая и понятная реализация.

Недостатки

- Требуется $O(m \cdot (n - m + 1))$ операций, вследствие чего алгоритм работает медленно в случае, когда длина паттерна достаточно велика

KMP (Knuth-Morris-Pratt algorithm)

Время работы алгоритма в худшем случае: $O(n)$.

Основная идея алгоритма KMP заключается в следующем: всякий раз, когда мы обнаруживаем несоответствие (после нескольких совпадений), мы уже знаем некоторые символы в тексте следующего окна. Мы используем эту информацию, чтобы избежать сопоставления символов, которые, как мы знаем, в любом случае будут совпадать. Давайте рассмотрим пример ниже, чтобы понять это

KMP (Knuth-Morris-Pratt algorithm)

```
txt = "AAAAABAAABA"
```

```
pat = "AAAA"
```

Мы сравниваем первое окно **txt** с **pat**

```
txt = "AAAAABAAABA"
```

```
pat = "AAAA"
```

Нашли совпадение. Здесь всё аналогично Naive string matching

KMP (Knuth-Morris-Pratt algorithm)

На следующем шаге мы сравниваем следующее окно **txt** с **pat**.

```
txt = "AAAAABAAABA"
```

```
pat = "AAAA" [pat сдвинулся на одну позицию]
```

Здесь KMP оптимизирует примитивный алгоритм

Во втором окне мы сравниваем только четвёртую А шаблона с четвёртым символом окна текста, чтобы решить, совпадает окно или нет.

Так как мы знаем, что первые три символа совпадут, мы пропустили их сравнение

Префикс-функция

Прежде чем перейти к описанию алгоритма, необходимо рассмотреть понятие префикс-функции.

Дана строка $s[0..n-1]$. Требуется вычислить для неё префикс-функцию, т.е. массив чисел $\pi[0..n-1]$, где $\pi[i]$ вычисляется следующим образом: это такая наибольшая длина собственного суффикса (окончания) подстроки $s[0..i]$, совпадающего с её префиксом (начало подстроки). В частности, значение $\pi[0]$ полагается равным 0.

Математическое определение префикс-функции можно записать следующим образом:

$$\pi[i] = \max_{k=0 \dots i} \{ k : s[0 \dots k-1] = s[i-k+1 \dots i] \}.$$

Например, для строки "abcabcd" префикс-функция равна: [0, 0, 0, 1, 2, 3, 0]

- у строки "a" нет нетривиального префикса, совпадающего с суффиксом;
- у строки "ab" нет нетривиального префикса, совпадающего с суффиксом;
- у строки "abc" нет нетривиального префикса, совпадающего с суффиксом;
- у строки "abca" префикс длины 1 совпадает с суффиксом;
- у строки "abcab" префикс длины 2 совпадает с суффиксом;
- у строки "abcabc" префикс длины 3 совпадает с суффиксом;
- у строки "abcabcd" нет нетривиального префикса, совпадающего с суффиксом

Другой пример — для строки "aabaab" она равна [0, 1, 0, 1, 2, 3]


```
1 def prefix(s):
2     v = [0]*len(s)
3     for i in range(1,len(s)):
4         k = v[i-1]
5         while k > 0 and s[k] != s[i]:
6             k = v[k-1]
7         if s[k] == s[i]:
8             k = k + 1
9         v[i] = k
10    return v
11
12 def kmp(txt, pat):
13     index = -1
14     f = prefix(pat)
15     print('prefix:', f)
16     k = 0
17     for i in range(len(txt)):
18         while k > 0 and pat[k] != txt[i]:
19             k = f[k-1]
20         if pat[k] == txt[i]:
21             k = k + 1
22         if k == len(pat):
23             index = i - len(pat) + 1
24             break
25     return index
```

KMP (Knuth-Morris-Pratt algorithm)

Особенности КМП-поиска:

1. требуется порядка $O(m+n)$ сравнений символов для получения результата;
2. схема КМП-поиска дает подлинный выигрыш только тогда, когда неудаче предшествовало некоторое число совпадений. Лишь в этом случае образ сдвигается более чем на единицу. К несчастью совпадения встречаются значительно реже чем несовпадения. Поэтому выигрыш от КМП-поиска в большинстве случаев текстов весьма незначителен.
3. Затраты памяти $O(m)$

Rabin-Karp algorithm

Алгоритм Рабина — Карпа — это алгоритм поиска подстроки с использованием хеширования.

Хеш-функция

Хеш-функция (англ. *hash function* от *hash* — «превращать в фарш», «мешанина»), или **функция свёртки** — функция, осуществляющая преобразование массива входных данных произвольной длины в выходную битовую строку установленной длины, выполняемое определённым алгоритмом. Преобразование, производимое хеш-функцией, называется **хешированием**.

Rabin-Karp algorithm

Как и примитивный алгоритм, алгоритм Рабина-Карпа также сдвигает шаблон на одно значение. Но в отличие от примитивного алгоритма, алгоритм Рабина Карпа сопоставляет хеш-значение шаблона с хеш-значением текущей подстроки текста, и если хеш-значения совпадают, то только тогда он начинает сопоставлять отдельные символы.

Rabin-Karp algorithm

Таким образом, алгоритм Рабина Карпа должен вычислять хеш-значения для следующих строк:

- 1) Сам паттерн
- 2) Все подстроки текста длины m .

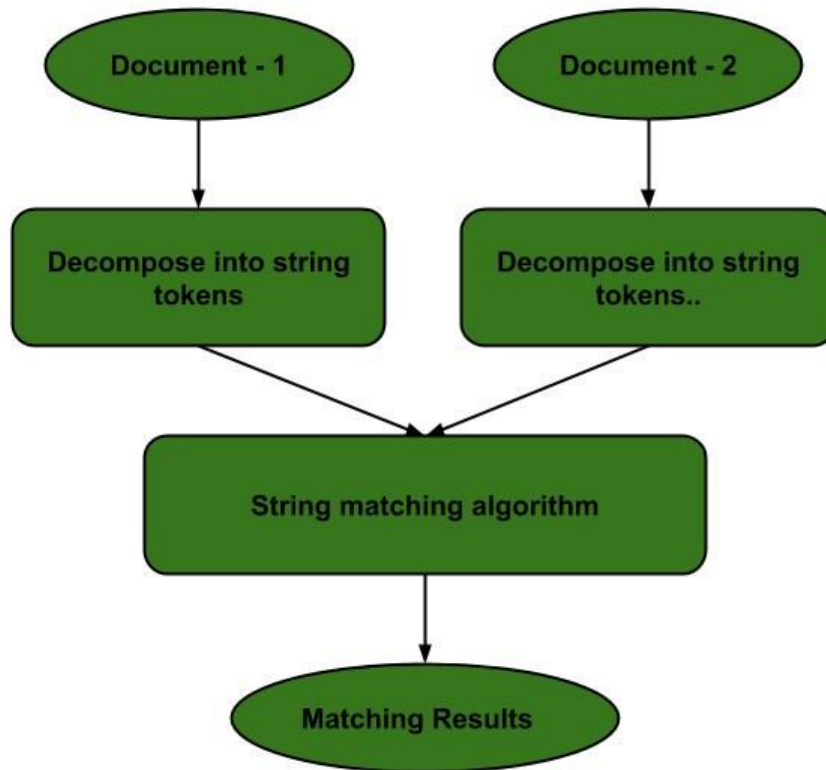
Rabin-Karp algorithm

Алгоритм редко используется для поиска одиночного шаблона, но имеет значительную теоретическую важность и очень эффективен в поиске совпадений множественных шаблонов одинаковой длины. Для текста длины n и шаблона длины m его среднее и лучшее время исполнения равно $O(n)$ при правильном выборе хеш-функции, но в худшем случае он имеет эффективность $O(nm)$, что является одной из причин того, почему он не слишком широко используется.

```
1 def RabinKarp(string, pattern):
2     n, m = len(string), len(pattern)
3     hpattern = hash(pattern);
4     for i in range(n-m+1):
5         hs = hash(string[i:i+m])
6         if hs == hpattern:
7             if string[i:i+m] == pattern:
8                 return i
9     return -1
```

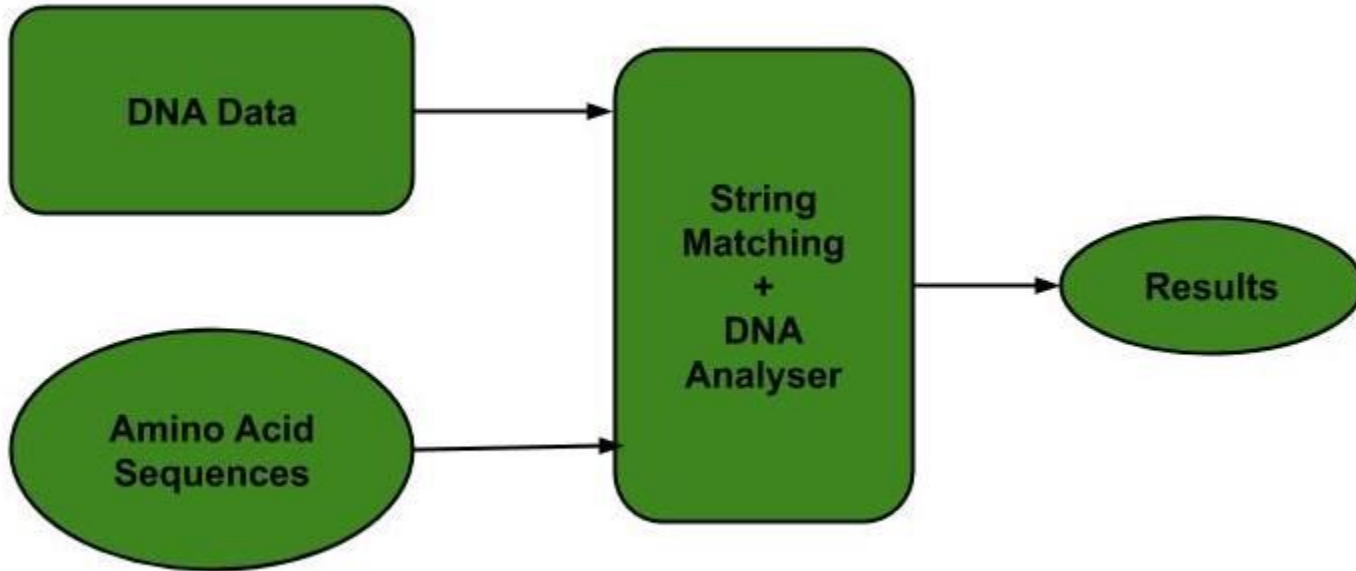

При

Антипла
сравнива
алгоритм
является

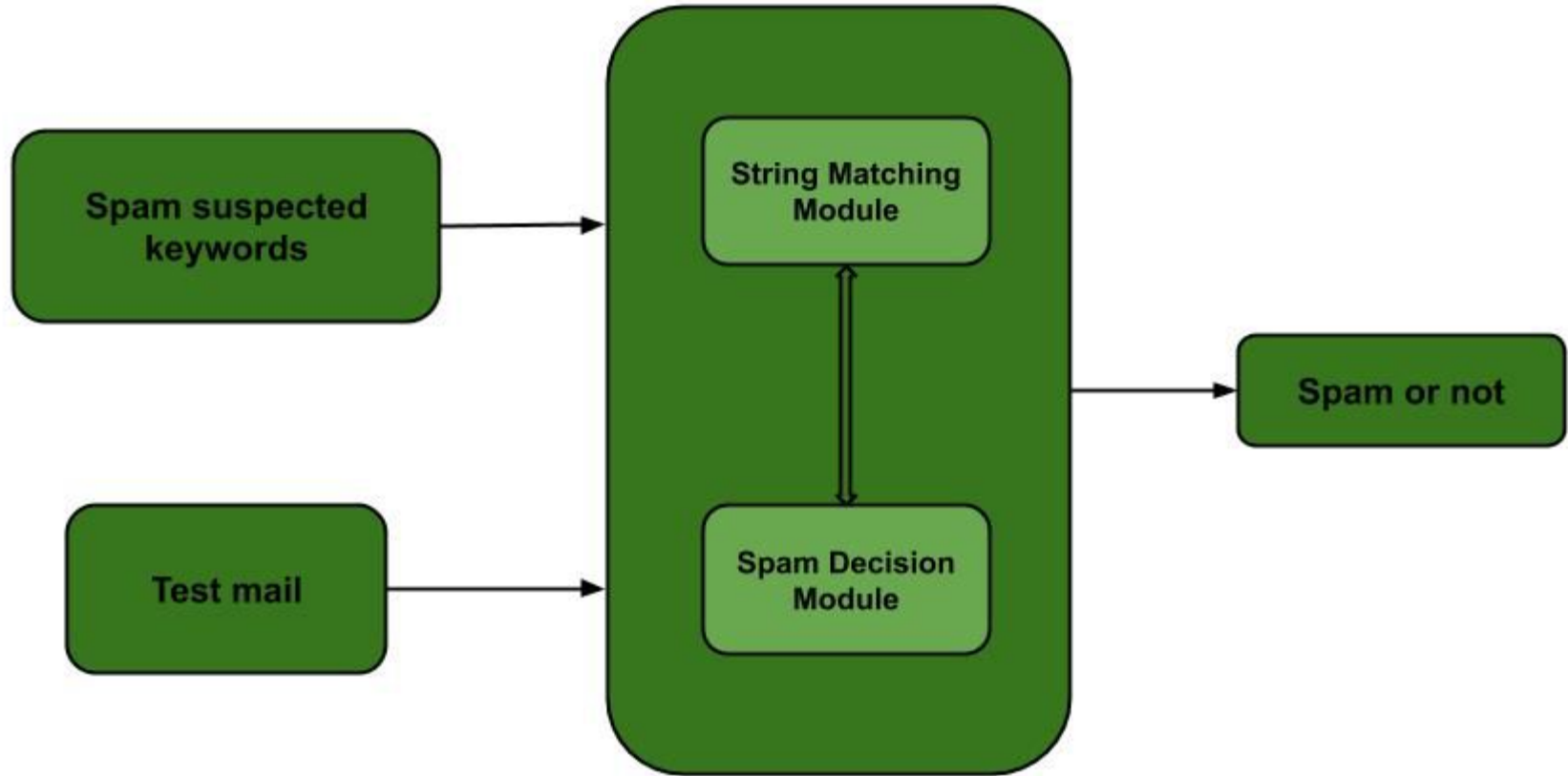


«И

ЭТИ
ТОГО,



- **Цифровая криминалистика:** Алгоритмы поиска подстроки используются для поиска конкретных текстовых строк, представляющих интерес в цифровом криминалистическом тексте, которые нужны в расследовании.
- **Проверка орфографии**



Поисковые системы или поиск контента в больших базах данных: для эффективной категоризации и организации данных используются алгоритмы поиска подстроки. Классификация осуществляется на основе ключевых слов поиска. Таким образом, алгоритмы сопоставления строк облегчают поиск информации

Система обнаружения вторжений: пакеты данных, содержащие ключевые слова, связанные со вторжением, обнаруживаются с помощью алгоритмов поиска подстроки. Весь вредоносный код хранится в базе данных, и все поступающие данные сравниваются с сохраненными данными. Если совпадение найдено, генерируется сигнал тревоги.

Выбор алгоритма

На сегодняшний день существует огромное разнообразие алгоритмов поиска подстроки. Программисту приходится выбирать подходящий в зависимости от таких факторов.

1. Нужна ли вообще оптимизация, или хватает примитивного алгоритма? Как правило, именно его реализуют стандартные библиотеки языков программирования.
2. «Враждебность» пользователя. Другими словами: будет ли пользователь намеренно задавать данные, на которых алгоритм будет медленно работать?
3. Грамматика языка может быть недружественной к тем или иным эвристикам, которые ускоряют поиск «в среднем».
4. Размер алфавита. Многие алгоритмы (особенно основанные на сравнении с конца) имеют эвристики, связанные с несовпавшим символом. На больших алфавитах таблица символов будет занимать много памяти, на малых — соответствующая эвристика будет неэффективной.
5. Возможность проиндексировать *строку*. Если таковая есть, поиск серьезно ускорится.
6. Требуется ли одновременный поиск нескольких строк? Приблизительный поиск? Побочные свойства некоторых алгоритмов (Ахо-Корасик, двоичного алгоритма) позволяют такое.