

# Fast sequence clustering using a suffix array algorithm

Stepachev Maxim

March 9, 2013

## Основные проблемы суффиксного дерева:

1. Проблема дополнительной памяти.
2. Аллокация памяти.
3. Достаточно высокая сложность реализации.

Пусть у нас есть некоторый образец  $P$  и текст  $S$ , требуется найти все вхождения образца в тексте.

$S$ : ATGTGTGC

$P$ : TGT

Суффиксный массив — это массив лексикографически отсортированных суффиксов.

Ранг суффикса — то место, которое займет суффикс при сортировке, то есть его позиция в суффиксном массиве.

S - Текст в котором происходит поиск, P - шаблон который мы ищем.

Lcp — longest common prefix

	rank	suffix array	
rosalind	7	4	alind
osalind	6	8	d
salind	8	6	ind
alind	1	5	lind
lind	4	7	nd
ind	3	2	osalind
nd	5	1	rosalind
d	2	3	salind

Сколько памяти требует суффиксный массив?

Простейший способ узнать, встречается ли образец в тексте, используя суффиксный массив, — взять первый символ образца и бинарным поиском по суффиксному массиву найти диапазон с суффиксами, начинающимися на такую же букву.

<u>s</u> ip	<u>s</u> ip	<u>s</u> ip
l	i	i
lppi	ippi	ippi
lssippi	issippi	issippi
lssissippi	ississippi	ississippi
Mississippi	mississippi	mississippi
Pi	pi	pi
Ppi	ppi	ppi
<u>S</u> ippi	<u>s</u> ippi	<u>s</u> ippi
<u>S</u> issippi	<u>s</u> issippi	sissippi
<u>S</u> sippi	ssippi	ssippi
<u>S</u> sissippi	ssissippi	ssissippi

Какова сложность простейшего поиска?

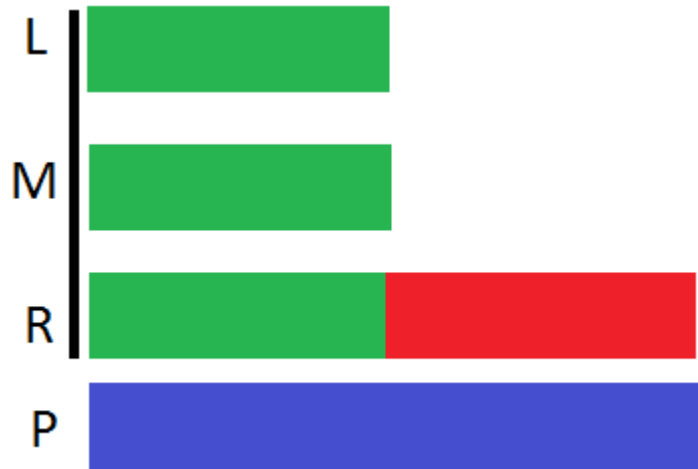
Какова сложность простейшего поиска?

Бинарный поиск работает за время равное  $O(\log|S|)$ , а сравнение суффикса с образцом не может превышать длины образца.

Таким образом время работы алгоритмы  $O(|P|\log|S|)$ , где  $S$ — текст,  $P$ — образец.



Основная идея заключается в вычислении общего префикса на краях диапазона.



Обозначения:

$L = \text{lcp}(\text{suffix}[l], P)$

$R = \text{lcp}(\text{suffix}[r], P)$

Пример:

L: ACCTCCCC

R: ACCT<sup>....</sup>GGGG

P: ACCTGGAA

Утверждение первое: Префикс для любой строки внутри диапазона  $\text{lcp}$  не меньше, чем минимум этих двух чисел.

Утверждение второе: Если общий префикс образца и любой строки внутри диапазона не меньше  $m = \min(L, R)$ , то  $m$  символов можно пропускать сразу.

LCP-массив — это массив значений  $lcp$  между лексикографически-соседними суффиксами.

Построить массив  $lcp$  достаточно просто. Решение в лоб, это идти по суффиксному массиву и считать  $lcp$  соседних — это  $O(N^2)$ .

Но если идти по массиву в правильном порядке, то можно использовать информацию с предыдущего шага для расчетов.

А именно, если перебирать суффиксы в том порядке, в котором они встречаются в исходной строке. Мы получим  $lcp$ -массив за линейное время  $O(N)$ .

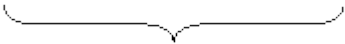
#	суффикс	шаг	$lcp$
1	i	10	1
2	ippi	7	1
3	issippi	4	4
4	ississippi	2	0
5	mississippi	1	0
6	pi	9	1
7	ppi	8	0
8	sippi	6	2
9	sissippi	3	1
10	ssippi	5	3
11	ssissippi	skip	

Range Minimum (Maximum) Query – запрос минимума (максимума) на отрезке в массиве.

### Постановка задачи:

Пусть дан массив  $A[1..n]$ . Нам необходимо уметь отвечать на запрос вида «найти минимум на отрезке с  $i$ -ого элемента по  $j$ -ый».

A[i]	3	8	6	4	2	5	9	0	7	1
i	1	2	3	4	5	6	7	8	9	10

  
 $RMQ(2, 7) = 2$

Более-менее простые алгоритмы позволяют отвечать на запросы за  $\langle O(n \log n), O(1) \rangle$ . Нас может устроить и такое решение, если мы строили массив за  $O(n \log n)$ , но можно быстрее...

Алгоритм применяется для решения за  $\langle O(n), O(1) \rangle$  времени специального случая задачи RMQ (поиск минимума на отрезке), в котором соседние элементы входной последовательности различаются на  $\pm 1$ .

### Применение в нашей задаче:

Построим так называемый lcp-массив — это массив значений lcp между лексикографически-соседними суффиксами.

Условные обозначения:

$L_r$  и  $R_r$  - левая и правая границы диапазона ответов в суффиксном массиве.

$L$  - левая граница диапазона поиска (изначально равна 0).

$R$  - правая граница диапазона поиска (изначально равна  $|S|-1$ ).

$M = (R+L)/2$

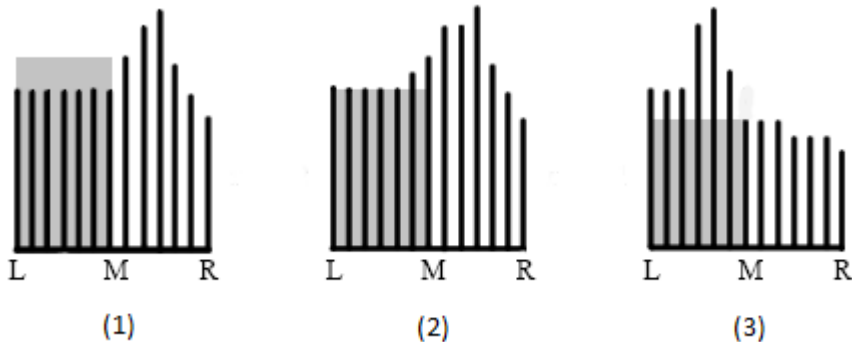
$l = \text{lcp}(\text{suffix}[L], p)$

$r = \text{lcp}(\text{suffix}[R], p)$

$ml = \text{lcp}(\text{suffix}[L], \text{suffix}[M])$

$mr = \text{lcp}(\text{suffix}[M], \text{suffix}[R])$

Подсчет  $l$  можно производить за  $O(1)$ , если применять алгоритм Фарака-Колтона и Бендера. Любая пара суффиксов  $\text{suffix}[]$  из диапазона  $[L, M]$  имеет хотя бы  $ml$  совпадений в префиксах. Аналогично любая пара суффиксов  $\text{suffix}[]$  из диапазона  $[M, R]$  имеет хотя бы  $mr$  совпадений в префиксах.



1.  $ml > l, L = M, l$  не меняем.
2.  $ml = l, (R = M, r = l + k + 1)$  или  $(L = M, l = l + k + 1)$
3.  $ml < l, R = M, r = ml$

Операция	Suffix tree	Suffix array
Построение:	$O( s )$	$O( s )$
Поиск:	$O( p )$	$O(\log s + p )$
Память:	$O( s )$	$O( s )$

<http://neerc.ifmo.ru>

<http://habrahabr.ru/post/115346/>

<http://software.intel.com/ru-ru/blogs/2012/05/22/accelerate-2012-3>

[http://en.wikipedia.org/wiki/Suffix\\_array](http://en.wikipedia.org/wiki/Suffix_array)

<http://en.wikipedia.org/wiki/Substring>

<http://e-maxx.ru/algo/>

Simple Linear Work Suffix Array Construction (<http://www.cs.helsinki.fi/u/tpkarkka/publications/icalp03.pdf>)

Suffix arrays: A new method for on-line string searches (<http://webglimpse.net/pubs/suffix.pdf> )

Fast sequence clustering using a suffix array algorithm

Handmade ;)