

Suffix Arrays

While we've seen that suffix trees require $O(m)$ space, the constants in that bound can be very large. A related data structure, suffix arrays, allow a much more compact representation at the cost of a small penalty at search time.

Simple idea: sort the suffixes of string S , store the starting indices of sorted suffixes in an array.

Ex:

$S = \text{BANANA}\$$

1. BANANA\$
2. ANANA\$
3. NANA\$
4. ANA\$
5. NA\$
6. A\$
7. \$

sort
————— S

7. \$
6. A\$
4. ANA\$
2. ANANA\$
1. BANANA\$
5. NA\$
3. NANA\$

store
indices

7
6
4
2
1
5
3

Another example:

$S = \text{cat\#cat}\$$

Searching for a query in a suffix array

Use binary search: main idea, suffixes that start with the same prefix will be stored in a contiguous range in the suffix array.

$r = |S|$

$l = 1$

For ~~each~~ c in pattern: $i = 1, \dots, |P|$
set l to smallest index within l, \dots, r that matches $P[i]$ in position i of suffix, if any

set r to largest index within l, \dots, r that matches $P[i]$ in position i of suffix, if any

if $l \leq r$ return suffixes in range l, \dots, r .

Constructing suffix arrays

① Naïve algorithm: sort the suffixes, $O(u \log u)$ comparisons, each comparison is $O(u)$. Total: $O(u^2 \log u)$.

② Linear algorithm: Construct sorted suffix tree, obtain indices ~~and~~ by traversing leaves in order.

$S = \text{BANANA}\$$

