

Suffix Array | Set 1 (Introduction)

We strongly recommend to read following post on suffix trees as a pre-requisite for this post.

Pattern Searching | Set 8 (Suffix Tree Introduction)

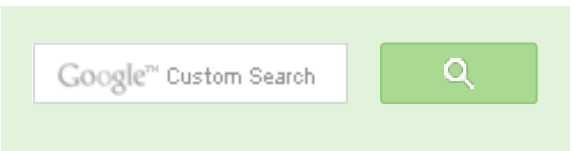
A suffix array is a sorted array of all suffixes of a given string. The definition is similar to **Suffix Tree** which is compressed trie of all suffixes of the given text. Any suffix tree based algorithm can be replaced with an algorithm that uses a suffix array enhanced with additional information and solves the same problem in the same time complexity (Source [Wiki](#)). A suffix array can be constructed from Suffix tree by doing a DFS traversal of the suffix tree. In fact Suffix array and suffix tree both can be constructed from each other in linear time. Advantages of suffix arrays over suffix trees include improved space requirements, simpler linear time construction algorithms (e.g., compared to Ukkonen's algorithm) and improved cache locality (Source: [Wiki](#))

Example:

Let the given string be "banana".

0	banana		5	a
1	anana	Sort the Suffixes	3	ana
2	nana	----->	1	anana
3	ana	alphabetically	0	banana
4	na		4	na
5	a		2	nana

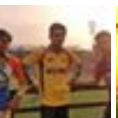
So the suffix array for "banana" is {5, 3, 1, 0, 4, 2}



GeeksforGeeks



53,519 people like [GeeksforGeeks](#).



[Facebook](#)

[Interview Experiences](#)

[Advanced Data Structures](#)

[Dynamic Programming](#)

[Greedy Algorithms](#)

[Backtracking](#)

[Pattern Searching](#)

[Divide & Conquer](#)

[Mathematical Algorithms](#)

[Recursion](#)

[Geometric Algorithms](#)

Naive method to build Suffix Array

A simple method to construct suffix array is to make an array of all suffixes and then sort the array. Following is implementation of simple method.

```
// Naive algorithm for building suffix array of a given text
#include <iostream>
#include <cstring>
#include <algorithm>
using namespace std;

// Structure to store information of a suffix
struct suffix
{
    int index;
    char *suff;
};

// A comparison function used by sort() to compare two suffixes
int cmp(struct suffix a, struct suffix b)
{
    return strcmp(a.suff, b.suff) < 0? 1 : 0;
}

// This is the main function that takes a string 'txt' of size n as an
// argument, builds and return the suffix array for the given string
int *buildSuffixArray(char *txt, int n)
{
    // A structure to store suffixes and their indexes
    struct suffix suffixes[n];

    // Store suffixes and their indexes in an array of structures.
    // The structure is needed to sort the suffixes alphabetically
    // and maintain their old indexes while sorting
    for (int i = 0; i < n; i++)
    {
        suffixes[i].index = i;
        suffixes[i].suff = (txt+i);
    }

    // Sort the suffixes using the comparison function
    // defined above.
    sort(suffixes, suffixes+n, cmp);

    // Store indexes of all sorted suffixes in the suffix array
    int *suffixArr = new int[n];
    for (int i = 0; i < n; i++)
```



Popular Posts

All permutations of a given string

Memory Layout of C Programs

Understanding "extern" keyword in C

Median of two sorted arrays

Tree traversal without recursion and without stack!

Structure Member Alignment, Padding and Data Packing

Intersection point of two Linked Lists

Lowest Common Ancestor in a BST.

Check if a binary tree is BST or not

Sorted Linked List to Balanced BST

```

        suffixArr[i] = suffixes[i].index;

    // Return the suffix array
    return suffixArr;
}

// A utility function to print an array of given size
void printArr(int arr[], int n)
{
    for(int i = 0; i < n; i++)
        cout << arr[i] << " ";
    cout << endl;
}

// Driver program to test above functions
int main()
{
    char txt[] = "banana";
    int n = strlen(txt);
    int *suffixArr = buildSuffixArray(txt, n);
    cout << "Following is suffix array for " << txt << endl;
    printArr(suffixArr, n);
    return 0;
}

```

Output:

```

Following is suffix array for banana
5 3 1 0 4 2

```

The time complexity of above method to build suffix array is $O(n^2 \log n)$ if we consider a $O(n \log n)$ algorithm used for sorting. The sorting step itself takes $O(n^2 \log n)$ time as every comparison is a comparison of two strings and the comparison takes $O(n)$ time.

There are many efficient algorithms to build suffix array. We will soon be covering them as separate posts.

Search a pattern using the built Suffix Array

To search a pattern in a text, we preprocess the text and build a suffix array of the text. Since we have a sorted array of all suffixes, **Binary Search** can be used to search. Following is the search function. Note that the function doesn't report all occurrences of pattern, it only report one of them.

Deploy Early. Deploy Often.

DevOps from
Rackspace:

Automation

[FIND OUT HOW ►](#)



```
// This code only contains search() and main. To make it a complete ru...
// above code or see http://ideone.com/1Io9eN
```

```
// A suffix array based search function to search a given pattern
// 'pat' in given text 'txt' using suffix array suffArr[]
void search(char *pat, char *txt, int *suffArr, int n)
{
    int m = strlen(pat); // get length of pattern, needed for strncmp

    // Do simple binary search for the pat in txt using the
    // built suffix array
    int l = 0, r = n-1; // Initilize left and right indexes
    while (l <= r)
    {
        // See if 'pat' is prefix of middle suffix in suffix array
        int mid = l + (r - l)/2;
        int res = strncmp(pat, txt+suffArr[mid], m);

        // If match found at the middle, print it and return
        if (res == 0)
        {
            cout << "Pattern found at index " << suffArr[mid];
            return;
        }

        // Move to left half if pattern is alphabtically less than
        // the mid suffix
        if (res < 0) r = mid - 1;

        // Otherwise move to right half
        else l = mid + 1;
    }

    // We reach here if return statement in loop is not executed
    cout << "Pattern not found";
}
```

```
// Driver program to test above function
int main()
{
    char txt[] = "banana"; // text
    char pat[] = "nan"; // pattern to be searched in text

    // Build suffix array
    int n = strlen(txt);
    int *suffArr = buildSuffixArray(txt, n);
```

Recent Comments

kzs please provide solution for the problem...

Backtracking | Set 2 (Rat in a Maze) · 3 minutes ago

Sanjay Agarwal bool

tree::Root_to_leaf_path_given_sum(tree...

Root to leaf path sum equal to a given number · 28

minutes ago

GOPI GOPINATH @admin Highlight this

sentence "We can easily...

Count trailing zeroes in factorial of a number · 30

minutes ago

newCoder3006 If the array contains negative
numbers also. We...

Find subarray with given sum · 54 minutes ago

newCoder3006 Code without using while
loop. We can do it..

Find subarray with given sum · 1 hour ago

Sanjay Agarwal You can also use the this
method:...

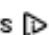
Count trailing zeroes in factorial of a number · 1
hour ago

AdChoices 

[▶ C++ Code](#)

[▶ Java Array](#)

[▶ Java Source Code](#)

AdChoices 

```
// search pat in txt using the built suffix array
search(pat, txt, suffArr, n);

return 0;
}
```

Output:

Pattern found at index 2

The time complexity of the above search function is $O(m \log n)$. There are more efficient algorithms to search pattern once the suffix array is built. In fact there is a $O(m)$ suffix array based algorithm to search a pattern. We will soon be discussing efficient algorithm for search.

Applications of Suffix Array

Suffix array is an extremely useful data structure, it can be used for a wide range of problems. Following are some famous problems where Suffix array can be used.

- 1) Pattern Searching
- 2) Finding the longest repeated substring
- 3) Finding the longest common substring
- 4) Finding the longest palindrome in a string

See [this](#) for more problems where Suffix arrays can be used.

This post is a simple introduction. There is a lot to cover in Suffix arrays. We have discussed a $O(n \log n)$ algorithm for Suffix Array construction [here](#). We will soon be discussing more efficient suffix array algorithms.

References:

<http://www.stanford.edu/class/cs97si/suffix-array.pdf>
http://en.wikipedia.org/wiki/Suffix_array

Please write comments if you find anything incorrect, or you want to share more information about the topic discussed above

► [C++ Array](#)

► [An Array](#)

► [Suffix Prefix](#)

AdChoices ►

► [Array in String](#)

► [Array of Arrays](#)

► [Array Function](#)



Related Topics:

- Remove minimum elements from either side such that $2 \times \text{min}$ becomes more than max
- Divide and Conquer | Set 6 (Search in a Row-wise and Column-wise Sorted 2D Array)
- Bucket Sort
- Kth smallest element in a row-wise and column-wise sorted 2D array | Set 1
- Find the number of zeroes
- Find if there is a subarray with 0 sum
- Divide and Conquer | Set 5 (Strassen's Matrix Multiplication)
- Count all possible groups of size 2 or 3 that have sum as multiple of 3



13



Tweet

3



1

Writing code in comment? Please use ideone.com and share the link here.

14 Comments

GeeksforGeeks

Sort by Newest ▼





from the discussion...



ram • 8 days ago

This uses hell of memory. isnt it ?

^ | v • Reply • Share ›



kaushik Lele • 23 days ago

I have followed "bucket" or "map" methodology to developed code.

e.g. for string "banab"

Starting character is 'b' so bucket of 'b' with value "anab" ('b' is node value)

Then next substring is "anab". So bucket of 'a' with value "nab" ('a' is node value)

Similarly another bucket 'n' for value "anab"

Now next string is "ab". We already have bucket for 'a' with value "anab".

So we will create subbuckets

Bucket 'a' -> sub-bucket 'n' with value "ab" (starting "a" and "n" are from bucket)

Bucket 'a' -> sub-bucket 'b' with value null (starting "a" from bucket)

If there was substring "anan..". Then we would have repeated steps

Bucket 'a' -> sub-bucket 'n' -> sub-sub-bucket 'a' and so on

Refer my code for word "banana". It prints value as expected.

5 a

3 ana

[see more](#)

^ | v • Reply • Share ›



sambhavsharma • 3 months ago

What if we search for "bana" or "ban"? I doesn't work for these patterns?

• Reply • Share ›



anonymous → sambhavsharma • a month ago

we are using `strncmp(text, pattern, pattern_length)` that why it matches

^ | v • Reply • Share ›



Kartik → sambhavsharma • 3 months ago

Please take a closer look. The program works fine for these cases also

^ | v • Reply • Share ›



sambhavsharma → Kartik • 3 months ago

I think $mid = l + (r-1)/2$ might not be the right condition,
It should be $(l+r)/2$

^ | v • Reply • Share ›



Kartik → sambhavsharma • 3 months ago

Please see <http://googleresearch.blogspot...>

^ | v • Reply • Share ›



sambhavsharma → Kartik • 3 months ago

This code doesn't work for the current condition.
for value "bana".
and the document says one possible fix could be:

```
int mid = low + ((high - low) / 2);
```

^ | v • Reply • Share ›



Saurabh Verma • 3 months ago

how to write code in comment box??

^ | v • Reply • Share ›



Kartik → Saurabh Verma • 3 months ago

Looks like code in Disqus comment system doesn't work very well. Use

link seems to be a good idea.

^ | v • Reply • Share ›



CREATORS93 • 3 months ago

this is great

^ | v • Reply • Share ›



geek • 3 months ago

There is a problem in the implementation of pattern search function search. strcmp should not be used, because it search for exact match. It is possible that pattern element, so instead of strcmp I would suggest using of Find() function. For example your implementation fails and returns pattern not found although it is present, I

1 ^ | v • Reply • Share ›



GeeksforGeeks → geek • 3 months ago

Geeks, thanks for pointing this out. We have updated the code to use <http://www.cplusplus.com/refer...>

^ | v • Reply • Share ›



Thewar • 4 months ago

Most awaited post :)

2 ^ | v • Reply • Share ›



Subscribe



Add Disqus to your site

