

# Z algorithm

Ulises Méndez Martínez

Algorist Weekly Talks

*ulisesmdzmtz@gmail.com*

August 7, 2015

# Overview

- 1 Motivation
  - Exact string matching
  - Prefix, suffix problem
- 2 Two Pointers
  - With two list
  - In a single list
- 3 Z algorithm
  - Definition
  - Algorithm
  - Sample
  - Implementation

At the end of the day we will be able to solve...

## Exact string matching

Given a (long) string **S** of length **n** and a shorter string pattern **P** of length **m** find all occurrences of **P** in **S**.

Occurrences of **P** are allowed to overlap.

## CF Beta Round 93. Problem B. Password

Given a string **S** find the longest substring which is both prefix and suffix of **S** and also appears inside **S**.

# Two Pointer Technique

The 2 pointer technique is mostly applicable in sorted arrays where we try to perform a search in  $O(N)$ .

# Two Pointer Technique

The 2 pointer technique is mostly applicable in sorted arrays where we try to perform a search in  $O(N)$ .

## With two list

Given two arrays (**A** and **B**) sorted in ascending order and an integer **X**, we need to find **i** and **j**, such that  $a[i] + b[j]$  is equal to **X**.

# Two Pointer Technique

The 2 pointer technique is mostly applicable in sorted arrays where we try to perform a search in  $O(N)$ .

## With two list

Given two arrays (**A** and **B**) sorted in ascending order and an integer **X**, we need to find **i** and **j**, such that **a[i] + b[j]** is equal to **X**.

## Solution

```
i = 0; j = b.size() - 1;
while( i < a.size() )
{
    while(a[i]+b[j]>X && j>0) j--;
    if(a[i]+b[j]==X) processAnswer(i,j);
    ++i;
}
```

## In a single list

Given a list of **N** integers, your task is to select **K** integers from the list such that its unfairness is minimized.

If  $(x_1, x_2, x_3, \dots, x_k)$  are **K** numbers selected from the list **N**, the unfairness is defined as  $\max(x_1, x_2, \dots, x_k) - \min(x_1, x_2, \dots, x_k)$  where max denotes the largest integer among the elements of **K**, and min denotes the smallest integer among the elements of **K**.

## In a single list

Given a list of **N** integers, your task is to select **K** integers from the list such that its unfairness is minimized.

If  $(x_1, x_2, x_3, \dots, x_k)$  are **K** numbers selected from the list **N**, the unfairness is defined as  $\max(x_1, x_2, \dots, x_k) - \min(x_1, x_2, \dots, x_k)$  where max denotes the largest integer among the elements of **K**, and min denotes the smallest integer among the elements of **K**.

## Solution

```
mn = INT_MAX;
sort(a, a+n);
for(i=0; i<=n-k; i++)
{
    mn = min(mn , a[i+k-1]-a[i]);
}
cout<<mn<<endl;
```



# Z function

## Definition

Given a string **S** of length **n**, the **Z Algorithm** produces an array **Z** where **Z[i]** is the length of the longest substring starting from **S[i]** which is also a **prefix** of **S**, i.e. the maximum **k** such that **S[j] = S[i + j]** for all  $0 \leq j < k$ . Note that **Z[i] = 0** means that **S[0] ≠ S[i]**.

# Strings with their Z values

aaaaa	aaabaab	abacaba
$z[0] = 0$	$z[0] = 0$	$z[0] = 0$
$z[1] = 4$	$z[1] = 2$	$z[1] = 0$
$z[2] = 3$	$z[2] = 1$	$z[2] = 1$
$z[3] = 2$	$z[3] = 0$	$z[3] = 0$
$z[4] = 1$	$z[4] = 2$	$z[4] = 3$
	$z[5] = 1$	$z[5] = 0$
	$z[6] = 0$	$z[6] = 1$

Table : Example of Z function

## Z-Box

The algorithm relies on a single, crucial invariant. As we iterate over the letters in the string (index  $i$  from  $1$  to  $n - 1$ ), we maintain an interval  $[L, R]$  which is the interval with **maximum**  $R$  such that  $1 \leq L \leq i \leq R$  and  $S[L...R]$  is a **prefix-substring**.

# Procedure

Now suppose we have the correct interval  $[L, R]$  for  $i - 1$  and all of the  $Z$  values up to  $i - 1$ . We will compute  $Z[i]$  and the new  $[L, R]$  by the following steps:

# Procedure

Now suppose we have the correct interval  $[L, R]$  for  $i - 1$  and all of the  $Z$  values up to  $i - 1$ . We will compute  $Z[i]$  and the new  $[L, R]$  by the following steps:

## If $i > R$

Then there does not exist a prefix-substring of  $S$  that starts before  $i$  and ends at or after  $i$ . If such a substring existed,  $[L, R]$  would have been the interval for that substring rather than its current value. Thus we “reset” and compute a new  $[L, R]$  by comparing  $S[0\dots]$  to  $S[i\dots]$  and get  $Z[i]$  at the same time ( $Z[i] = R - L + 1$ ).

Otherwise,  $i \leq R$ , so the current  $[L, R]$  extends at least to  $i$ . Let  $k = i - L$ . We **know that**  $Z[i] \geq \min(Z[k], R - i + 1)$  because  $S[i...]$  matches  $S[k...]$  for at least  $R - i + 1$  characters (they are in the  $[L, R]$  interval which we know to be a prefix-substring).

Now we have a few more cases to consider.

Otherwise,  $i \leq R$ , so the current  $[L, R]$  extends at least to  $i$ . Let  $k = i - L$ . We **know that**  $Z[i] \geq \min(Z[k], R - i + 1)$  because  $S[i...]$  matches  $S[k...]$  for at least  $R - i + 1$  characters (they are in the  $[L, R]$  interval which we know to be a prefix-substring).

Now we have a few more cases to consider.

**If  $Z[k] < R - i + 1$**

Then there is no longer prefix-substring starting at  $S[i]$  (or else  $Z[k]$  would be larger), meaning  $Z[i] = Z[k]$  and  $[L, R]$  stays the same. The latter is **true** because  $[L, R]$  only changes if there is a prefix-substring starting at  $S[i]$  that extends beyond  $R$ , which we know is not the case here.

$$Z[k] \geq R - i + 1$$

Then it is possible for  $S[i\dots]$  to match  $S[0\dots]$  for more than  $R - i + 1$  characters (i.e. past position  $R$ ). Thus we need to update  $[L, R]$  by setting  $L = i$  and matching from  $S[R + 1]$  forward to obtain the new  $R$ . Again, we get  $Z[i]$  during this.

- The process computes all of the  $Z$  values in a single pass over the string, so we are done.
- Correctness is inherent in the algorithm and is pretty intuitively clear.



Example **S** = **aabcaabxaaaz**

$v \setminus i$	0	1	2	3	4	5	6	7	8	9	10	11
$S_i$	a	a	b	c	a	a	b	x	a	a	a	z
$L$	0	1	2	3	4	4	4	7	8	9	9	11
$R$	0	1	1	2	6	6	6	6	9	10	10	11
$k$	-	-	-	-	-	1	2	-	-	-	1	-
$Z[k]$	-	-	-	-	-	1	0	-	-	-	1	-
$Z[i]$	0	1	0	0	3	1	0	0	2	2	1	0

Table : Example of Z algorithm

See algorithm running step by step in the following [▶ Link](#)

# Algorithm Code

## Example ( C++ Implementation )

```
int L = 0, R = 0;
for (int i = 1; i < n; i++) {
    if (i > R) {
        L = R = i;
        while (R < n && s[R-L] == s[R]) R++;
        z[i] = R-L; R--;
    } else {
        int k = i-L;
        if (z[k] < R-i+1) z[i] = z[k];
        else {
            L = i;
            while (R < n && s[R-L] == s[R]) R++;
            z[i] = R-L; R--;
        }
    }
}
```

# Solution to our original problems

## String matching

We can do this in  $O(n + m)$  time by using the **Z** Algorithm on the string **P\$S** (that is, concatenating **P**, **\$**, and **S**) where **\$** is a character that matches nothing. The indices **i** with  $Z[i] = m$  correspond to matches of **P** in **S**.

## Problem B

We simply compute **Z** for the given string **S**, then iterate from **i** to **n - 1**. If  $Z[i] = n - i$  then we know the suffix from **S[i]** is a prefix, and if the largest **Z** value we've seen so far is at least **n - i**, then we know some string inside also matches that prefix.

```
for (int i = 1; i < n; i++){
    if (z[i] == n-i && maxx >= n-i) { res = n-i; break; }
    maxx = max(maxx, z[i]);
}
```

# Q & A

## References

- <http://codeforces.com/blog/entry/3107>
- <https://www.hackerrank.com/challenges/pairs/topics/two-pointer-technique>
- <http://e-maxx-eng.github.io/string/z-function.html>