

# MATH578A Assignment 1

Raktim Mitra

USC ID: 1487079265 email: raktimmi@usc.edu

February 23, 2020

**Q1.**

---

## Q2. (Chapter 1, exercise 11)

**Q:** Let  $T$  be a text string of length  $m$  and let  $S$  be a multiset of  $n$  characters. The problem is to find all substrings in  $T$  of length  $n$  that are formed by the characters of  $S$ . For example, let  $S = \{a, a, b, c\}$  and  $T = \text{'abahgcabab'}$ . Then 'caba' is a substring of  $T$  formed from the characters of  $S$ . Give a solution to this problem that runs in  $O(m)$  time. The method should also be able to state, for each position  $i$ , the length of the longest substring in  $T$  starting at  $i$  that can be formed from  $S$ .

**Ans:** Let,  $A[ ]$  contains unique elements of  $S$ .  $\text{Count}[ ]$  array contains number of occurrences of each element of  $A$  in  $S$ .

Now, we shall construct the array  $D[ ]$  where  $D[i]$  contains length of the longest substring of  $T$  that ends at position  $i$  which can be formed by elements from  $S$ .

$i$  ranges from 0 to  $m-1$ .

We also need a  $\text{RunningCount}[ ]$  array similar to  $\text{Count}[ ]$  which we will use to keep track of counts of elements of  $S$  faced for  $D[i]$ .

$A$ ,  $\text{Count}$ ,  $\text{RunningCount}$  all have the length  $\text{count}(\text{unique}(S))$ . Note:  $A.\text{index}[c]$  gives position of char  $c$  in  $A$ , constant time since  $A$  is of constant length assuming alphabet size is constant. Therefore, we can write the following recurrence for  $D[i]$ .

$$D[i] = \begin{cases} 0 & \text{if } T[i] \text{ is not in } A \text{ (case 1)} \\ \text{if } T[i] \text{ is in } A \begin{cases} D[i-1] + 1 & \text{if } \text{RunningCount}[A.\text{index}[T[i]]] < \text{Count}[A.\text{index}[T[i]]] \\ & (\uparrow \text{case 2a}) \\ i - j & \text{if } \text{RunningCount}[A.\text{index}[T[i]]] == \text{Count}[A.\text{index}[T[i]]] \\ & \text{where } j \text{ is the first occurrence of } T[i] \text{ after } i-1-D[i-1] \\ & (\uparrow \text{case 2b}) \end{cases} \end{cases}$$

For the above recurrence to work properly, we need to update  $\text{RunningCount}$  in the following manner:

$$\begin{cases} \text{(case 1) } \text{RunningCount}[i] = 0 \text{ for all } i \\ \left\{ \begin{array}{l} \text{(case 2a) } \text{RunningCount}[A.\text{index}[T[i]]] += 1 \\ \text{(case 2b) } \text{RunningCount}[A.\text{index}[T[k]]] -= 1, \text{ for } k \text{ in interval } (i, j) \end{array} \right. \end{cases}$$

After we fill up the array D we can report all starting occurrence position as:

$$\{i - \text{size}(S) \text{ for } i \text{ if } D[i] == \text{size}(S) \}$$

We can also construct length of longest substrings starting at i and formed by elements of S from D. Let, this array be B[ ]. Shown in full pseudocode below:

Algorithm 1: Multiset Matching

```

Input: Text T and Multiset S
1: A ← Unique(S)
2: initialize Count[0...length(A)] to all 0
3: initialize RunningCount[0...length(A)] to all 0
4: for element e in S do
5:   Count[A.index[e]]++
6: end for
7: initialize D[0...m-1]
8: D[-1] ← 0 //For convenience, In actual implementation we would
   //hardcode D[0] and start the following loop at 1.
9: for {i=0; i < m; i++} do
10:  if T[i] not in A then //Case 1
11:    D[i] ← 0
12:    {RunningCount[j] ← 0 for j in interval [0,length(a))}
13:  else if RunningCount[A.index[T[i]] < Count[A.index[T[i]]] //Case 2a
   then
14:    D[i] ← D[i-1] + 1
15:    RunningCount[A.index[T[i]]]++
16:  else //Case 2b
17:    j ← i - D[i-1] + 1 //Case 2b
18:    while j ≤ i and T[j] != T[i] do
19:      RunningCount[A.index[T[j]]]--
20:    end while
21:    D[i] ← i - j
22:  end if
23: end for
24: Initialize B[0..m-1] //length of substrings longest substrings formed
   //by elements of S from starting positions.
25: i ← 0
26: Initialize R //List of starting positions of occurrences
27: while {i < m} do //O(m)
28:   B[i - D[i]] ← D[i]
29:   if D[i] == size(S) then
30:     R.add(i - D[i])
31:   end if
32: end while
33: return R, B

```

**Time Complexity:** Clearly all the preprocessing can be done in  $O(m)$  time and so is generating R and B from D.

That leaves constructing D. The for loop (line 9-23) runs  $m$  times. Case 1 and Case 2a are clearly  $O(1)$  [Assuming alphabet size is  $O(1)$ ]. However, the while loop (Case 2b, line 18-20) can be  $O(\text{size}(S))$  in worst case. But, we can use amortized analysis to show that amortized cost of each iteration of the for loop is  $O(1)$ . The reason is, case2b can decrement RunningCount (line 19)  $\text{size}(S)$  times only when Case(2a) has been executed  $\text{size}(S)$  times. (Because only Case 2a increments RunningCount (line 15).)

Therefore, Case 2b can be  $O(\text{size}(S))$  only after  $O(\text{size}(s))$  iterations of case 2b which are all  $O(1)$  iterations. Therefore, on average each iteration takes  $O(1)$ , making the whole for loop  $O(m)$ .

Hence, time complexity of Multiset Matching is  $O(m)$ .

---

### Q3. (Chapter 6, exercise 1.)

**Q:** Construct an infinite family of strings over a fixed alphabet, where the total length of the edge-labels on their suffix trees grows faster than  $O(m)$  ( $m$  is the length of the string). That is, show that linear-time suffix tree algorithms would be impossible if edge-labels were written explicitly on the edges.

**Ans:**

---

### Q4. (Chapter 7, exercise 1.)

**Q:** Given a set  $S$  of  $k$  strings, we want to find every string in  $S$  that is a substring of some other string in  $S$ . Assuming that the total length of all the strings is  $n$ , give an  $O(n)$ -time algorithm to solve this problem.

**Ans:**

---

### Q5. (Chapter 7, exercise 2.)

**Q:** For a string  $S$  of length  $n$ , show how to compute the  $N(i)$ ,  $L(i)$ ,  $L'(i)$  and  $sp_i$  values (discussed in Sections 2.2.4 and 2.3.2) in  $O(n)$  time directly from a suffix tree for  $S$ .

**Ans:**

---