
Problem Set 3

Name: Thomas Ren

Collaborators: None

Problem 3-1.

- (a) For this part, I have written the code to implement this hash table. And the results shows as follows:

```
1 [' (36)<->(92)', '', '', '', ' (56)', ' (47)<->(61)<->(33)', ' (52)']
```

The specific code implemented are as follows:

```
1 # Problem Session 3 -> Problem 3-1 Hash It Out
2 def Hash_Out(A, a, b, k):
3     '''
4     h(key) = (a*key+b) mod k
5     :param A: a list of integer keys
6     :param a: a in hash function h(key)
7     :param b: b in hash function h(key)
8     :param k: k in hash function h(key)
9     :return: a hash table chaining with sequence i.e. list in Python
10    '''
11    HASH = [[] for i in range(k)]
12    for item in A:
13        new_item = (a*item+b)%k
14        HASH[new_item].append(item)
15
16    return HASH
17
18
19 # Problem Set 3 -> Problem 3-1
20 def Hash_Chaining_linked_list(A,a,b,n):
21    '''
22    :param A: a list of integer keys
23    :param a: a in hash function h(key)
24    :param b: b in hash function h(key)
25    :param n: n in hash function h(key)
26    :return: a hash table chaining with doubly linked list
27    '''
28
29    Hash_table = Hash_Out(A,a,b,n)
```

```

30     chaining_hash_table = []
31     for slot in Hash_table:
32         # if slot is not None:
33         chaining = Doubly_Linked_List_Seq()
34         chaining.build(slot)
35         chaining_hash_table.append(chaining.__str__())
36
37     return chaining_hash_table

```

- (b) If we consider hashing by hash function $h(k) = (10k + 4) \bmod c$, we would get c slots. Since the list A has **length 7**, we cannot fit all the items without collision for any $c < 7$.

Thereby, we will try to enumerate c starting from 7.

Here, we give the implementation as follows:

```

1  def Hash_Out_double_modulo(A, a, b, c, n):
2      '''
3      h(key) = (a*key+b) mod n
4      :param A: a list of integer keys
5      :param a: a in hash function h(key)
6      :param b: b in hash function h(key)
7      :param c: c in hash function h(key)
8      :param n: n in hash function h(key)
9      :return: a hash table chaining with sequence i.e. list in Python
10     '''
11     HASH = [[] for i in range(n)]
12     for item in A:
13         new_item = (a*item+b)%c%n
14         HASH[new_item].append(item)
15
16     return HASH
17
18
19 def collision(A,a,b,c,n):
20     '''
21     h(key) = (a*key+b) mod c mod n
22     :param A: a list of integer keys
23     :param a: a in hash function h(key)
24     :param b: b in hash function h(key)
25     :param c: c in hash function h(key)
26     :param n: n in hash function h(key)
27     :return: a boolean value tells whether a collision exists in our hash table
28     '''
29     Hash_table = Hash_Out_double_modulo(A,a,b,c,n)
30     for slot in Hash_table:
31         if len(slot)>1:
32             print('collision exists for slot',slot)
33             return False
34

```

```
35     return True
36
37
38 def find_min_c(A,a,b,n):
39     '''
40     h(key) = (a*key+b) mod c mod n
41     :param A: a list of integer keys
42     :param a: a in hash function h(key)
43     :param b: b in hash function h(key)
44     :param n: n in hash function h(key)
45     :return: c_min: the minimum value satisfying the non-collision situation
46     '''
47     c_min = n
48     while True:
49         if collision(A,a,b,c_min,n):
50             print('The minimum c without causing collision is:',c_min)
51             return c_min
52         else:
53             print('collision exists for c =',c_min)
54             c_min += 1
```

The results for the case given in the problem is as follows:

```
1 collision exists for slot [36, 92]
2 collision exists for c = 7
3 collision exists for slot [36, 52, 56, 92]
4 collision exists for c = 8
5 collision exists for slot [61, 52]
6 collision exists for c = 9
7 collision exists for slot [47, 61, 36, 52, 56, 33, 92]
8 collision exists for c = 10
9 collision exists for slot [52, 92]
10 collision exists for c = 11
11 collision exists for slot [56, 92]
12 collision exists for c = 12
13 The minimum c without causing collision is: 13
14
15 Process finished with exit code 0
```

Problem 3-2. For simplicity, we assign r_1 to be the room ID of k_1 and r_2 to be the room ID of k_2 .

- (a) This could work if $k_1 \equiv k_2 \pmod n$.
 Then, we would have $r_1 = (ak_1 + b) \pmod n$.
 Since $k_1 \equiv k_2 \pmod n$, we can write $k_2 = sn + k_1$, for some constant s .
 Then, we could write $r_2 = ak_2 + b \pmod n = a(sn + k_1) + b \pmod n = asn + ak_1 + b \pmod n = ak_1 + b \pmod n = r_1$
- (b) Since u is a really large number, so if k changes for a small value relative to u , the results won't change.
 For example, if $k_2 = k_1 + 1$, then $r_1 = r_2$, except for the threshold.
 So, if we choose k_1 to be 0 and k_2 to be 1. Then $r_1 = r_2 = 0$.
 Thereby, they will be assigned to the same room.
- (c) This hash function is known to be the **universal hash family**. According to the lecture, the probability of any two different values being assigned to the same hash slot is at most $\frac{1}{n}$.
 Thereby, their highest probability of being roommates is $\frac{1}{n}$.

Problem 3-3.

- (a) For each ice core, we get a **unique** identifier.

The first thought comes to our mind should be using **Direct Access Array Sort**, since there is no collision exists. However, **Direct Access Array Sort** takes $O(u)$ time to run, here, $u = 128^{16\lceil\log_4(\sqrt{n})\rceil}$ for which we assign a unique value for an ASCII character. Then, we would deduce that $O(u) = O(n^c)$, for some constant c . This is clearly worse than comparison model e.g. merge sort which only takes $O(n \log n)$ time.

Since we get that $u = O(n^c)$, for some constant c , we would know immediately that we can achieve linear complexity by using **Radix Sort**. The running time is $\Theta(n + n \log_n u) = \Theta(n + cn) = \Theta(n)$

- (b) Since there is a bound for u that $u = O(800,000)$, we don't need **Radix Sort** to achieve linear running time.

Here, we will use **Counting Sort** because there is no promise for no collision. Still, we will get a sorted array in $\Theta(800,000 + n) = \Theta(n)$.

- (c) Multiply every key of thickness by n^3 to get all keys to be integers i.e. m which ranging from 0 to $4n^3$.

Then use **Radix Sort** to sort integer keys m takes $\Theta(n + n \log_n 4n^3) = \Theta(n)$.

- (d) This part of problem limit our algorithm to **Comparison Models**, which gives best case $\Theta(n \log n)$ running time.

To be more specific, we will use **Merge Sort**, which has "two-finger algorithm" for *merge* operation.

Problem 3-4.

- (a) Since this part requires **expected** $O(n)$ -time, we build a **Hash Table** in **linear time** so that we can achieve certain item by its key in **constant time**.

The algorithm is as follows:

```

1  for b_i in B:
2      b_j = find(r-b_i)
3      if find(r-b_i) is None:
4          return False
5      else:
6          if |j-i| <= n/10:
7              return True
8          else:
9              pass

```

- (b) Since $r < n^2 \Rightarrow r = O(n^2)$, we can use **Radix Sort** to first sort the B in **linear time**. Then we will use the two-finger algorithm to approach our best pair sum t to target r , after that, determine whether they are **close pair**, if so, return the pair, else, continue moving our two fingers.

The algorithm is as follows:

```

1  # delete all the items that are larger than r
2  for i in range(len(S)):
3      if S[i] >= h:
4          del S[i:]
5          break
6  # two-finger algorithm
7  i, j = 0, len(B)-1
8  target_i, target_j = i, j
9  while j - i > 1:
10     pair_sum = S[i] + S[j]
11     if pair_sum > h:
12         j -= 1
13     elif pair_sum < h:
14         i += 1
15     else:
16         if abs(s[i].key-s[j].key)<len(B)/10:
17             return s[i].key,s[j].key
18         else:
19             pass

```

Problem 3-5.

- (a) Since we need to return the anagram substring count of B in A in $O(k)$ -time, we cannot traverse the whole set of A , which is $O(n)$ -time, we should, instead, give each substring a key, which is the same for anagrams, so that we can determine whether B exists in A and how many or not.

We solve this by giving the *key* of each substring a frequency table of alphabets, which is implemented in Python using **tuple** data structure. This frequency table should be hashable as well, by which we mean immutable.

Then we collect all the frequency table to a hash table, where the key is the frequency table and the item is the number of this frequency table, we implemented this hash table using **dictionary** in Python.

Then, we take $O(|A|)$ -time to build the **dictionary** of substring table by shift one space and synchronize the frequency table in $O(1)$ -time. For each B , we build the frequency table of B in $O(k)$ -time, then we can find the number of anagram of B in $O(1)$ -time in the hash table.

The Python code is implemented as follows:

```

1 def single_count_anagram_substring(A,B):
2     '''
3     :param A: a string
4     :param B: a substring with length less than A
5     :return: the number of anagram substring B in A
6     '''
7     k = len(B)
8     substring_table = {}
9     for i in range(len(A)-k+1): # O(|A|+k) = O(|A|)
10         if i == 0:
11             substring = A[:k]
12         else:
13             substring = substring[1:]+A[i+k-1]
14             frequency_table = build_frequency_table(substring)
15             if frequency_table in substring_table:
16                 substring_table[frequency_table] += 1
17             else:
18                 substring_table[frequency_table] = 1
19
20     B_frequency_table = build_frequency_table(B)
21     if B_frequency_table not in substring_table:
22         print('There is no anagram substring of B in A')
23         return 0
24     else:
25         return substring_table[B_frequency_table]
26
27
28
29
30 def build_frequency_table(str):

```

```

31     '''
32     :param str: a string
33     :return: a tuple of frequency table of alphabets
34     '''
35     alphabet = [0 for i in range(26)]
36     for letter in str:
37         index = ord(letter) - 97
38         alphabet[index] += 1
39
40     return tuple(alphabet)

```

- (b) We can first build the substring hash table in $O(|T|)$ -time and for each string in S , we do the aforementioned step to find the number of anagram of this string, which takes $O(k)$ -time. For traversing the whole S , takes $O(n)$ -time, thus the whole process takes $O(|T| + nk)$ -time.

The Python implementation code is as follows:

```

1  def count_anagram_substrings(T, S):
2      '''
3      Input:  T | String
4              S | Tuple of strings S_i of equal length k < |T|
5      Output: A | Tuple of integers a_i:
6              | the anagram substring count of S_i in T
7      '''
8      A = []
9      #####
10     # YOUR CODE HERE #
11     #####
12     k = len(S[0])
13     n = len(S)
14     A = [0 for i in range(n)]
15     substring_table = {}
16     for i in range(len(T) - k + 1): # O(|A|+k) = O(|A|)
17         if i == 0:
18             substring = T[:k]
19         else:
20             substring = substring[1:] + T[i + k - 1]
21             frequency_table = build_frequency_table(substring)
22             if frequency_table in substring_table:
23                 substring_table[frequency_table] += 1
24             else:
25                 substring_table[frequency_table] = 1
26
27     for index, s in enumerate(S):
28         s_frequency_table = build_frequency_table(s)
29         if s_frequency_table in substring_table:
30             A[index] = substring_table[s_frequency_table]
31         else:
32             A[index] = 0

```



```
33  
34  
35     return tuple(A)
```

(c) Submit your implementation to `alg.mit.edu`.