

①

↳ Pseudo Code for add()

add(HashTable T, tuple T)

O(1)

- Calculate hashValue of the key of the tuple

O(1)

- Go to the array list pointed by the $T[\text{hashValue}]$

O(nK)

- for each element in the array list

O(K) ° if the array list contains a tuple with the same key & value, only increment the count of the tuple in the array by 1.

° if we reach the end of the ArrayList and the tuple was not found, then we add the tuple to the array list

O(n)

- calculate load Factor of the hashTable T

O(1)

- If ($\text{loadFactor} > 0.7$)

° create a new hashTable T_1 of size that is the prime which is immediately greater than twice the size of hashTable T.

° create a new Hash Function

O(nK)

- for each element in the HashTable T

° find the hash value of the key & add the tuple to the array list at that location.

i. $\text{add}(\text{hashTable } T, \text{tuple } t) \in O(nK)$

- However, on the average case we consider the add operation to take $O(K)$ time.

Pseudo Code for search (hashTable T, tuple t)

search (hashTable T, tuple t)

- O(1) - calculate hash value of the key of the tuple.
- go to the array list pointed by the $T[\text{hashValue}]$
- O(K) - for each tuples in the ArrayList
 - O(K) - if s_0 equals (t)
 - return $t \cdot \text{getCount}()$
 - else return 0

$\therefore \text{search}(T, t) \in O(K)$

	Similarity	Runtime	(2)
4)	Brute Force	0.39867774	62.748 s
	Hash String	0.39867774	24.448 s
	Hash Code	0.39868575	4.842 s

5) Hash Code Similarity produces the shortest runtime.

Brute Force Similarity produces the longest runtime.

- Although, it's not true in my implementation, but the runtime for Hash String & Hash Code should be near to each other.
- Brute Force is obviously going to take the maximum time. The number of distinct elements in S1 is 47535 (with length 8) & in S2 is 50271. B
- Brute force would have to compare each of the 8-length shingles in S1 to every other shingle in S2, thus producing the highest runtime.
- Hash String compares the strings individually as well, but only the strings which have been hashed in the same bucket. This is a mere fraction of the total number of strings.
- In this particular example, the maxload() returned a value of 6. That means even in the very worst case, the Hash String was comparing each 8-length string with only 6 strings. Contrast this with the brute force where each string was being compared with 50,000+ strings & it becomes evident that Hash String would run much faster.
- Hash Code on the other hand does not compare strings at all & deals only with integers. Thus is a constant time operation & therefore takes absolutely minimum amount of time.

6>

The Brute Force & HashString return the same value for similarity. However, the HashCode class returns a different value for similarity.

- This because in HashCode class we are composing the hash values obtained by rollover hash function and not the strings themselves.
- Therefore, we might run into a case where the rolling hash value obtained from 2 different strings ends up being the same. In such a case, although the strings are different, we would still end up calculating them as similar cases & thus lead to a higher value.

37 Brute Force Similarity:

a) Data Structures:

- String `str1` & `str2` are split into equal length strings of length "`length`" each. Each of these strings of length `length` are stored in an `ArrayList` of `Strings`.
- In my code, this corresponds to `array1 S1-array` and `array2 S2-array`. Both of these contain repetitions.
- The unique elements from both of the arrays are merged to form a single `ArrayList` `arrayU` in my code. This contains distinct elements and is the same as `set U` given in the problem.
- We also define a new class CountingTuples that stores - the string Value & the number of times that particular string appears in `S1-array` & `S2-array`.
- We define an ArrayList of class CountingTuples called `myCount`.
- The basic idea behind this is that when we want to find out the count of a particular string, instead of querying the entire `S1-array` / `S2-array`, we instead query the `ArrayList` `myCount` which has already stored the string & count for that string in each of the arrays `S1-array` & `S2-array` for us.

b) Pseudo Code:

$O(n-k)$ → split string s_1, s_2 into length of k length each and store them in an array list s_1 -array & s_2 -array respectively.

2) calculate length of s_1 :

$O(1)$ → 2.1) initialize 2 array lists $\text{marked } s_1$ & $\text{marked } s_2$

$O(n)$ → 2.2) for each element in s_1 -array {

$O(nk)$ → 1b) (element is not present in $\text{marked } s_1$) {
add element to $\text{marked } s_1$

$O(n)$ → for (each element in s_1 -array starting from the current)
{ { - count the number of times the element occurs
in s_1 -array } }

$O((m+n)k)$ → - check if there is a tuple of value string in myCount array list

$O(1)$ → - create a new tuple of class CountingTuples with fields Value (string itself), count-In-S1 = 0,
count-In-S2 = 0 if no element found in prev step

$O(1)$ → - set count-In-S1 to the count obtained in the previous step

$O(1)$ → - sum += count * count // calculate a running
sum to give vector length

vectorSum1 = Math.sqrt(sum)

$O(km(m+n))$ 3) similarly calculate length of $s_2()$

4) Similarity :

$O(n+m)$ 4.1) for each tuple i in myCount arrayList

sum += count-In-S1 * count-In-S2.

similarity = sum / (length of s_1 * length of s_2)

∴ Runtime $\in O(mnK)$

Hash Code Similarity :

a) Data Structure :

- String str1 and str2 are split into equal length strings of length length each. Each of these strings of length length are stored in an arraylist of strings.
- This corresponds to the array s1-array & s2-array. Both of these contain repetitions.
- The unique elements from both of these arrays are merged to form a single arraylist arrayU. This set contains distinct elements and is equivalent to set U given in the problem.
- We create 2 hash tables hashTable_s1 and hashTable_s2 that contain unique elements from the multisets s1-array & s2-array respectively.
- These hash tables are built using the HashTable class and therefore store elements as Tuples (with value of string and count)
- Basic Idea is that for every string value stored in the Array U, we make 2 searches. The first one is in the hashTable - S1 that returns the count of the string in s1-array. The second one is in hashTable - S2 that returns the count of that string in the s2-array.
- Once we have these counts, we can calculate length of S1(), length of S2 & similarity()

b) Pseudo Code:

$O(m-k)$ 1) Split string s_1, s_2 into lengths of k length each and store them in an arraylist $s1_array$ & $s2_array$ respectively.

2) Create hashTable- $s1$ from the strings stored in $s1_array$

2.1) Initialize an instance of the hashTable class with hash table size of 8 \Rightarrow HSS-Table-1

$O(k)$ 2.2) Compute & store an exponent array that contains powers of the alpha value upto $(\text{alpha})^{\text{length}} \Rightarrow \text{exponent}[]$

$O(m)$ 2.3) Create a new array which stores each character of the string $s1 \Rightarrow$ array1[]

2.4) Calculate the first hashValue by doing the following:
 - for ($i=0$; $i < \text{length}$; $i++$)
 {

`add()` takes constant

time

since we are

Comparing primitives

```
char c = array1[i].charAt(0) // take the first
// single length string from the array1 and
// convert it into a char
asciiValue = (int) c; // calculate the ascii value
// of the character.
```

O(1) 2.5) Create a new Tuple and add it into the hashTable_s1 with values (h1, s1-array[0])

$O(m-k)$ 2.6) calculate remaining values by using Rollover Hashing.

$O(1)$ - Create a new tuple and add it to the hashTable - $s1$ with values ($h\text{-next}$, $s1\text{-array}[i]$)

return HashTable - S1

(4)

$O(n) \ 3 \triangleright$ Use the same process outlined in step 2 to create the hashTable-S2 from the contents of the array-S2.

$4 \triangleright$ Create Array U that contains distinct elements from the union of hashTable-S1 and hashTable-S2.

$O(m) \ 4.1 \triangleright$ Copy all the tuple.Key() of hashTable-S1 into Array U.

$O(n) \ 4.2 \triangleright$ For each tuple in the hashTable-S2 {

- get the Key of the tuple

Average $O(1)$ - Search for the key in hashTable-S1

- If key does not exist in hashTable-S1, add the key to the Array U

}

5) Calculate length of S1()

$O(m+n) \ 5.1 \triangleright$ for each key in the Array U {

Average: $O(1)$ - search for the tuple with the same key in hashTable-S1

- return the count of the tuple

- calculate running sum

$$\text{sum} += (\text{count} * \text{count})$$

}

vector sum = Math.sqrt(sum);

$O(m+n) \ 6 \triangleright$ Calculate length of S2() similarly.

7) Calculate Similarity()

$O(m+n) \ 7.1 \triangleright$ for each key in the Array U {

$O(1)$ - Search for the tuple with the same key in hashTable-S1

- return the count of the tuple = count_1

$O(1)$ - Search for the tuple with the same key in hashTable-S2

- return the count of the tuple = count_2.

$$\text{sum} += \text{count}_1 * \text{count}_2.$$

}

Runtime $\in O(m+n)$

Hash String Similarity

- Data Structures:

- Strings `str1` and `str2` are split into equal length strings of length `sLength` each. Each of these strings are stored in an array list of strings.
- This corresponds to the array `S1-array` & `S2-array`. Both of these contain repetitions.
- We create 2 hash tables `hashTable-S1` and `hashTable-S2` that contain unique elements from the multi-set `S1-array` & `S2-array` respectively.
- These hash tables are built using the HashTable class & therefore store elements as tuples (with value of string & count)
- We create another hash table `hashTable-U`. The unique elements from both the hash tables `hashTable-S1` and `hashTable-S2` are merged to form a distinct set of elements. This is the same as the set `U` in the problem.
- The basic idea is that we are going to iterate through the `hashTable-U` and for each element in `table-U` we are going to find out the count for that element by searching in `hashTable-S1` and `hashTable-S2`.

Pseudo Code for Hash String Similarity

$O(m-k)$ 1) Split the strings s_1, s_2 into lengths of $s.length$ each & store them in an array list $s1_array \& s2_array$ respectively.
 $O(n-k)$ 2) Create hashTable - $s1$ from the strings stored in the $s1_array$
 $O(k)$ 2.1) Initialize an instance of the `HashTable` class with `HashTable size of 8` \Rightarrow `HSS-Table-1`.
 $O(k)$ 2.2) Compute & store an exponent array that contains powers of the alpha values upto $(\text{alpha})^{\text{length}}$ \Rightarrow `exponent []`
 $O(m)$ 2.3) Create a new array which stores each character of the string `char1 \Rightarrow array1 []`
 2.4) Calculate the first hash Value by doing the following:
 $O(k)$ - for ($i=0$; $i < \text{length}$; $i++$)
 $\quad \quad \quad \{$

```
char c = array1[i].charAt(0) //take the first  
//single length string from the array1 and convert it  
//into a char.
```

```
asciiValue = (int)c ; // calculate the ascii value of  
// the character -
```

since we will be comparing the value() and the key() with each tuple.

$O(K)$ 2.5 > Create a new Tuple with values (h_1 , $s1$ -array [0]) and add it to the hashTable- $S1$.

2.6) Calculate remaining values using rollover hashing.

$O(m \cdot k)$ — for (int $i = 1$; $i < \text{array1.length} - \text{length} + 1$; $i++$) {
 h-next = ($h^i_1 - \text{subtract}$) * alpha + $\text{array1}[\text{value} - \text{add}]$.
 ↓ ↓ | ↓

next hash rollover prev value ascii value of the ascii value of the
 value of hash character to be character to be
 rollover removed added -

$O(k)$ - create a new Tuple t ($h\text{-next}$, $\underline{S1}$ -array[i])
and add it to the hashTable- $\underline{S1}$

3

rotisserie hashTable S1.

(G)

3> Use the same process outlined in step 2 to create the hashTable-S2 from the contents of the array - S2.

4> Create hashTable-U

4.1> Initialize an instance of the hashTable class with a hashtable of size 8 \Rightarrow HSS-Table-U

$O(m)$ 4.2> for each tuple in the hashTable-S1

average $O(1)$ - add the tuple to hashTable-U
time.

$O(nk)$ 4.3> for each tuple in the hashTable-S2. {

$O(k)$ - search the tuple.get Value() in the hashTable-U

$O(1)$ - if hashTable-U does not contain the tuple then
add the tuple to the hashTable-U.

{}

5> calculate lengthOfS1()

5.1> for each tuple in the hashTable-U {

$O(k)$ - search the tuple.get Value() in the hashTable-S1.

- return the tuple.get Count() \Rightarrow count

- calculate the running sum

sum += (count * count)

{}

5.2> vector Length = Math.sqrt(sum)

$O((m+n)k)$ 6> calculate length of S2() similarly.

7> calculate Similarity()

$O((m+n)k)$ 7.1> for each tuple in hashTable-U {

$O(k)$ - search the tuple.get Value() in the hashTable-S1

- return the count of the tuple = count-1

$O(k)$ - search the tuple.get Value() in the hashTable-S2

- return the count of the tuple = count-2.

sum += count-1 * count-2.

{}

Run Time $\in O((m+n)k)$