

Problem Set 3

All parts are due Thursday, March 20 at 11:59PM. Please download the .zip archive for this problem set, and refer to the `README.txt` file for instructions on preparing your solutions. Remember, your goal is to communicate. Full credit will be given only to a correct solution which is described clearly. Convoluting and obtuse descriptions might receive low marks, even when they are correct. Also, aim for concise solutions, as it will save you time spent on write-ups, and also help you conceptualize the key idea of the problem.

Your Name: Eric Klinkhammer

Collaborators: Name1, Name2

Part A

Problem 3-1.

(a) Data Structure Description

A hash table has to be used if the user wants constant time access to a random element. A hash table is a data structure that maps every object it stores (and can store) into an integer, and ultimately into an index in an array.

In this way, the hash table knows where to put every object (neglected collisions for the moment) before it begins. This is what saves time on the tasks.

With regards to collisions, or when different objects are mapped to the same index, this can be handled in different ways. The two most obvious are through either an additional data structure (a list, for example, if insert is the desired fast operation) in a process known as chaining, or by probing, where the object is put into another slot in the hash table through some process (for example just putting it in the next open slot).

(b) Algorithm Description

The creation of a hash table from an unsorted array involves just iterating through the array once and mapping each element into the hash table. To reduce collisions, I would create a table with size on the order of but slightly bigger than n , the number of elements in the array. I would then hash them all individually. Each hash has an amortized cost of $O(1)$, so the entire array can be converted in $O(n)$ time.

To get an unsorted array from the hash table, you would iterate through the array of the table and return the results. This would take $O(m)$ time, where m is the size of the table, but since m was defined to be order n (and it should be, to achieve the right mix of collisions and space saving), the run time of this is also $O(n)$.

(c) Set Membership

To test for set membership, I would take the key of the object I am searching for and look up the location associated with the key. If the location was empty, I would be done and return false. If it had the object we were looking for, I would be equally done. The other case, if it has a different object, would require looking wherever or however the hash table chooses to resolve its collisions.

This test for membership on average takes $O(1)$, but in the worst case, assuming all of the data hashes to the same spot (and then the collision resolution works identically for all of them), it will be $O(n)$ time.

(d) Insertion and Delete

Insert and delete work in a similar fashion as the test for set membership, except, in the case of insert, when the empty slot is found the object is put into that slot, and, in the case of delete, when the object is found, it is deleted. It is important that the delete makes the slot blank (so it won't confuse future collisions). Accordingly, they have a similar run time: $O(1)$ on average, but $O(n)$ in the worse case.

(e) Set Intersection

The intersection of two sets is all of their common elements. Therefore, to build an intersection set, one only has to traverse the members of one set, check the membership of that element in the other set, and, if found, insert it in the resulting set.

The unsorted array of elements can be taken from the first hash table in $O(n)$ time, and then there will be $O(n)$ lookups and at most $O(n)$ inserts. This means that making the new set will also be an $O(n)$.

(f) Set Union

The union of two sets is all of their elements, without duplicates. Since the insert method in a set already handles duplicates, constructing the union of two sets involves only inserting every element from both into a new set.

It would be $O(n)$ to get the elements in each table, as before, and then $O(n)$ again to insert them into the new set. Again, this entire operation is $O(n)$.

If the sets were of different sizes, say n_1 and n_2 , then the total time would be $O(n_1 +$

n_2).

(g) Set Difference

The set difference between sets A and B (A - B) would be all of the elements in A not also found in B.

This can be created by finding the intersection of A and B, and then iterating through all of the intersection points, deleting them from set A. In practice, you could also copy set A and delete them from the copy, so as to preserve the original sets. Either way, the intersection would be $O(n)$ as above, and the iterating and deleting would be $O(n)$ operations each taking $O(1)$ on average. If you were creating a copy, that could be done in $O(n)$ time, so, regardless, the set difference operation can be done in $O(n)$ time.

Problem 3-2. Probabilistic Analysis of Hashing With Chaining

(a) Expected Time of Search

Given a random key, it is equally likely to be map to any value. For all but one of those values, it will immediately return false. The expected value is calculated as follows:

$$\begin{aligned}
 E_t &= \sum_{h \in m} t(h) * p_h(h) \\
 E_t &= \sum_{i=0}^n t(i) * p_h(i) \\
 E_t &= \sum_{i=0}^{n-1} t(i) * p_h(i) + t(n) * p_h(n) \\
 E_t &= \sum_{i=0}^{n-1} c * \frac{1}{m} + cn * \frac{1}{m} \\
 E_t &= c \frac{n-1}{m} + c \frac{n}{m} \\
 E_t &= \frac{2n-1}{m} \\
 E_t &= O(1)
 \end{aligned}$$

The last line assumes that m and n are the same order (which they should be in a hash table).

(b) Expected Time for Successful Search

In this case, the key is not random, but is one already present in the hash table. Given

that the hash table is essentially a linked list at this point, the expected time for a successful search is roughly $n/2$, and is $O(n)$.

As before, the expected time to complete the search is the sum over all possible inputs of the time it would take to find that input and the probability of that time (from the definition of expected value).

$$E_t = \sum_{i=1}^n p_i t_i$$

All times have a probability of $\frac{1}{n}$

$$= \frac{1}{n} \sum_{i=1}^n i$$

The expected time of each element is the position of the element in the list.

$$\begin{aligned} &= \frac{(n)(n+1)}{2n} \\ &= \frac{n+1}{2} \end{aligned}$$

That number is the middle of n , assuming n starts at 1.

(c) Expected Time for Search with normal inserts

To determine the expected time for search, we first have to determine the average number of collisions for a given key value.

We define a new random variable CH_i as the number of keys hashed to position i . It has a probability distribution as below. It is important to recognize that such a random variable is a Binomial distribution (with success in the Bernoulli trial being defined as a collision).

$$CH_i = \binom{n}{i} \left(\frac{1}{m}\right)^i \left(1 - \left(\frac{1}{m}\right)\right)^{n-i}$$

Because it is a binomial, we know that the expected value is $\frac{n}{m}$. Search is a constant time operation unless it runs into a collision. Since the expected number of collisions

is $\frac{n}{m}$, the expected run time (not the worse case) of a successful search is the time per collision (constant) times the number of expected collisions (constant in this case) - $O(\frac{n}{m}) = O(1)$.

Problem 3-3. Meandering Hashes

(a) Naive Brute Force

Assuming only one orientation is acceptable, you would have to perform up to k^2 checks at each of the potential starting positions of the smaller array, of which there are $(n - k)^2$. The running time is then $O((n - k)^2 k^2)$, which, since n is much larger than k , simplifies to $O((nk)^2)$.

If additional orientations desired, you can rotate the matrices accordingly and recompute. That won't change the algorithmic complexity.

(b) Algorithm Description

A rolling hash could be used here in a similar way as the 1-D case. The entire target matrix could be hashed and the value stored. Then, a hash could be computed for the initial k by k block. Then, the hash would roll across the row, each time computing the value of the window that contains k different elements (a new column).

You could store the hash values already computed (in a $n-k \times n-k$ matrix) and this way always roll the hash "down," but it should be equally efficient to simply go back and forth across the matrix, going down one row each time.

I could see storing as you go being easier to run in parallel, but they're the same complexity. If desired, this would need only $O(1)$ memory on the side.

The algorithm's code below assumes that the hash function works in such a way that elements can be rolled off the top and onto the bottom and off the left side and onto the right side. I cannot immediately think of a hash function that would allow for this in a simple manner. If such a method of rolling hashes does not exist, then calculate the entire row and go down all possible columns.

That would take more time (nk^2) for the first row, but it would be the same asymptotically.

(c) Pseudocode

```
int 2D compare( S, T )
```

```

hash = hash(T);
rolling_hash = hash(T sized matrix with left corner at S(0,0)
for rows in range(n-k-1):
    rolling_hash.remove(top_row) # Roll off top k values
    rolling_hash.add(new_row) # Roll on k values
    temp <- rolling_hash
    for cols in range(n-k-1):
        rolling_hash = rolling_hash.remove(old_column)
        rolling_hash = rolling_hash.add(new_column)
        check rolling_hash with T:
            if match, compare element by element
    rolling_hash = temp

```

As described above, this computes a hash function for every square sub matrix of order k in the larger matrix. The only reason this is more efficient than the brute force method is that the hash can be computed in $O(k)$ steps instead of $O(k^2)$ steps.

The temporary variable is used to start again on the left side of the matrix. It is equally possible to traverse back and forth, but that is harder to code in a simple loop.

(d) Time Complexity

This algorithm will run in $O(kn^2 + k)$. It must do the initial calculation and the final one (hopefully just one) that both take $O(k)$. The algorithm then does $(n - k)^2$ rolls, but each roll involves k steps. Therefore, the hash rolling and comparison will take $O(kn^2)$. Finally, the constant time to check if the equality is a collision or the real value will bring up the running time to $O(kn^2 + k)$.

This is asymptotically faster than $O(kn^2 + k)$

(e) Space Complexity

This algorithm should use $O(1)$ space, seeing as it only has to store a few integers (hash(T), the value of the one on the left most column (to go down), and the current rolling hash value.