

Week 11 Studio Sheet

(Solutions)

Useful advice: The following solutions pertain to the theoretical problems given in the studio classes. You are strongly advised to attempt the problems thoroughly before looking at these solutions. Simply reading the solutions without thinking about the problems will rob you of the practice required to be able to solve complicated problems on your own. You will perform poorly on the exam if you simply attempt to memorise solutions to the studio problems. Thinking about a problem, even if you do not solve it will greatly increase your understanding of the underlying concepts. Solutions are typically not provided for Python implementation questions. In some cases, pseudocode may be provided where it illustrates a particular useful concept.

Studio Problems

Problem 1. A hashtable can be used to store different types of elements such as integers, strings, tuples, or even arrays. Assume we want to use a hashtable to store arrays of integers (i.e., each element is an array of integers). Consider the following potential hash functions for hashing the arrays of positive integers. Rank them in terms of quality and give a brief explanation of the problems with each of them. Assume that they are all taken mod m , where m is the table size.

- Return the first number in the array (e.g., if $\text{array} = [10, 5, 7, 2]$, hash index will be $10 \% m$)
- Return a random number in the array (e.g., if $\text{array} = [10, 5, 7, 2]$, hash index will be a randomly chosen element from the array mod m)
- Return the sum of the numbers in the array (e.g., if $\text{array} = [10, 5, 7, 2]$, hash index will be $24 \% m$)

Give a better hash function than these three and explain why it is an improvement.

Solution

From best to worst:

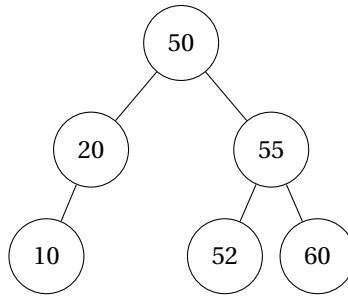
1. **Return the sum of the numbers in the array.** This is a decent hash function, as it accounts for all of the contents of the array. It is not amazing since it will produce identical hashes for permutations of the same elements (or arrays having the same sum), e.g., $\text{array1} = [10, 5, 7, 2]$, $\text{array2} = [2, 5, 7, 10]$ and $\text{array3} = [9, 6, 7, 2]$ will all collide.
2. **Return the first number in the array.** This is not very good since it only accounts for one element of the array, so any two arrays that begin with the same element will collide.
3. **Return a random number in the array.** This is not even a valid hash function since it might return a different value for the same array when used multiple times. Hash functions must be consistent, that is they must always produce the same value for the same key.

A better hash function would be

$$h(A) = A[0] + A[1]x + A[2]x^2 + \dots + A[n]x^n,$$

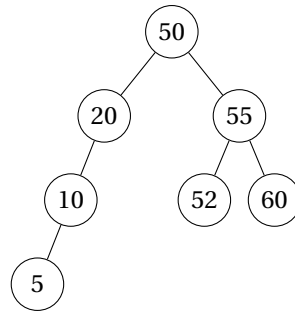
for some value of $x > 1$. This is better since if we use this hash function, all of the array is taken into account, and permuting the order of the elements is likely to change the value of the hash. We can make it even better by selecting $x > \text{any value we could see in our input}$ (for example, if we were hashing just a-z, giving values of 1-26, then we could choose $x=27$) and taking the sum mod p for a suitable prime p , which yields the standard polynomial hash function.

Problem 2. Insert 5 into the following AVL tree and show the rebalancing procedure step by step.

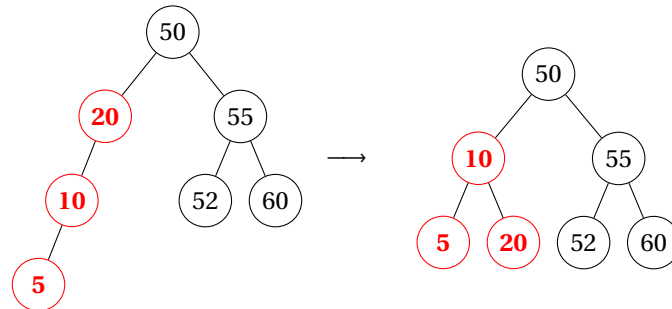


Solution

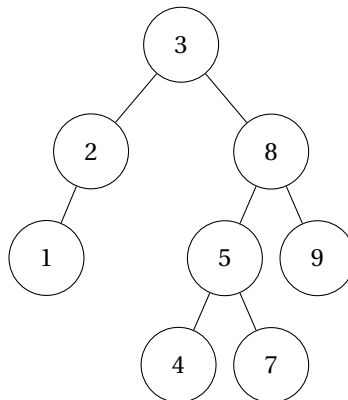
We insert 5 using the ordinary BST insertion procedure and we obtain



This tree is imbalanced, as the node containing 20 has a height 2 left subtree, and a height 0 right subtree, resulting in a balance factor of 2. To rectify this, we rotate the node 10:

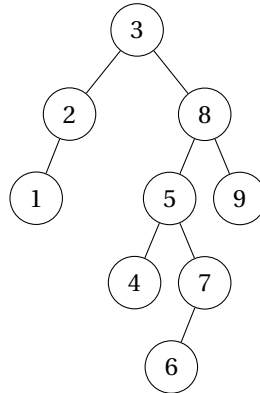


Problem 3. Insert 6 into the following AVL tree and show the rebalancing procedure step by step.



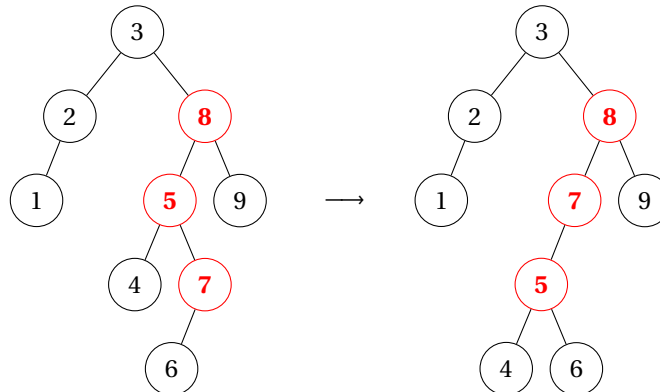
Solution

We insert 6 using the ordinary BST insertion procedure and we obtain the following tree.

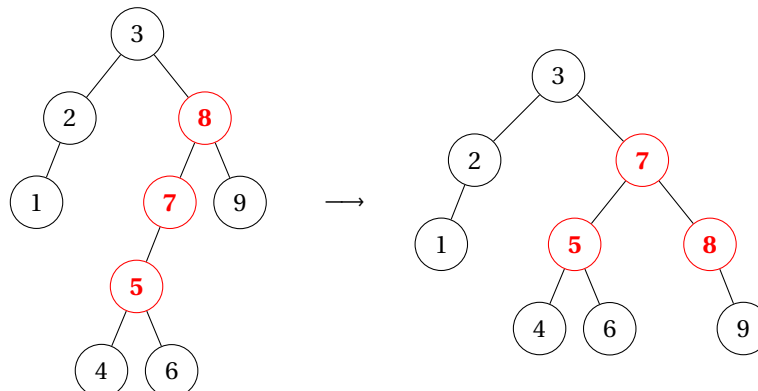


This tree is imbalanced at node 8 since it has a left subtree of height 3 and a right subtree of height 1, leading to a balance factor of 2. Note that the root node 3 is also imbalanced, but we always perform rotations on the lower levels first since this may fix the imbalances at higher levels.

Since node 8 has a taller left subtree and node 5 has a taller right subtree, this is a left-right imbalance so we must perform a double rotation on node 7 to correct it. We perform the first rotation to obtain:



We then perform the second rotation and are left with:



This is a balanced tree, so we are done.

Problem 4. Recall the algorithm for deleting an element from a binary search tree. If that element has no children, we do nothing. If it has one child, we can simply remove it and move its child upwards. Otherwise, if

the node has two children, we swap it with its *successor* and then delete it. This algorithm works because of the following fact: The successor of a node with two children has no left child. Prove this fact.

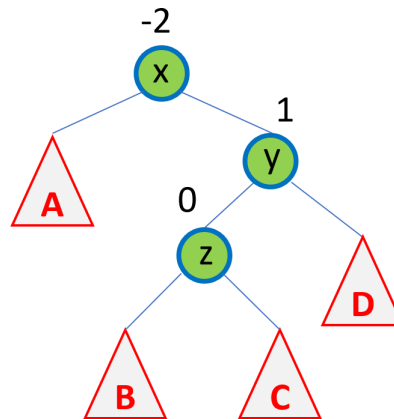
- First prove that the successor of a node x with two children must be contained in the right subtree of x .
- Use this fact to prove that the successor of x cannot have a left child.

Solution

Recall that the successor of a node is the minimum value node that is greater than x . Consider a node x with two children x_l and x_r . First, we will show that the successor of x lies in the right subtree r_x of x . All elements in the left subtree of x have values $< x$ and hence cannot be the successor of x . Suppose that x has a parent p_x . If $p_x < x$ then x is a right child and hence neither p_x nor its left subtree can contain the successor of x since every node v in the left subtree of p_x satisfies $v < p_x < x$. Otherwise x is a left child of p_x and hence $x_r < p_x$ and in particular $x_r < v$ for every v in the right subtree of p_x since $x_r < p_x < v$, therefore the successor of x is not a right descendant of p_x . This leaves the only possibility that the successor is in the subtree of r_x .

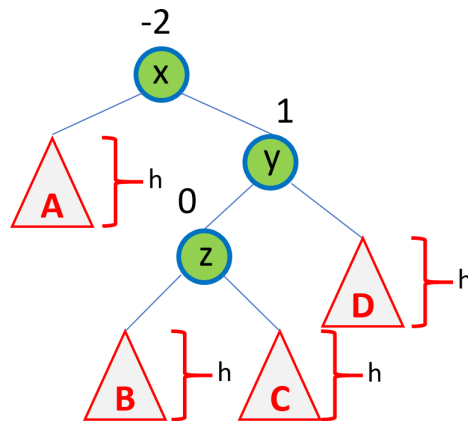
Suppose that the successor y of x had a left child y_l . By the BST property we have $y_l < y$, but since y_l is contained in the right subtree of x , we also have $x < y_l$. This means that $x < y_l < y$ which implies that y is not in fact the successor of x . This is a contradiction and hence y must have no left child.

Problem 5. Consider the binary search tree shown below, with balance factors as indicated. Assume that the subtrees A, B, C and D are AVL trees. Show that performing the rotation algorithm for the right-left case results in an AVL tree.

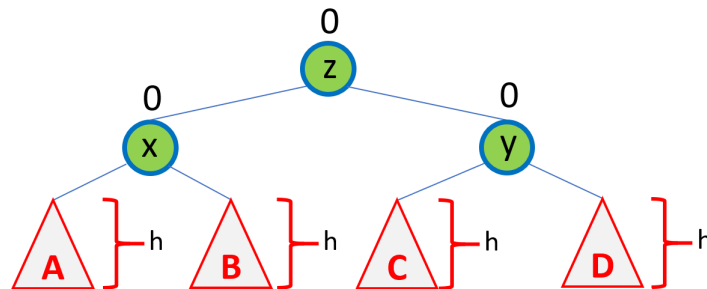


Solution

Let the height of subtree B be h . We can deduce the heights of subtrees C, D and A are also h (based on the balance factors).



After applying the two phases of the right-left case rotations, we obtain the following situation:



Based on the heights of A, B, C and D, we can deduce the balance factors of x, y and z as all being 0. Since A, B, C and D are AVL trees, the balance factors of all nodes in the tree are in $\{-1, 0, 1\}$. So the tree is now an AVL tree.

Problem 6. Consider the following potential hash functions for hashing integers. Rank them in terms of quality and give a brief explanation. Assume that they are all taken mod m , where m is the table size.

- Return $x \bmod 2^p$ for some prime p .
- Return $(ax + b)$ for some positive, and randomly chosen, a, b .
- Return a random integer.

Solution

Ranked from best to worst:

- Return $(ax + b)$ for some positive a, b . This hash is ok. It could be improved by taking the result modulo a prime number.
- Return $x \bmod 2^p$ for some prime p . Note that, if p is small, taking modulo powers of two simply keeps the bottom p bits of the integer, which is bad. For p sufficiently large, this is just the same as just taking the value of x as the output of our hash function, which is fine if the keys are random, but not so great if they are not.
- Return a random integer. This does not hash the same key to the same place consistently, and so is not even a valid hash function.

Problem 7. Consider the following probing schemes and for each of them, explain whether they do or do not suffer from primary or secondary clustering. In each problem, h returns an index where the item k is to be

inserted when probing for the i^{th} time. h', h_1, h_2 are hash functions.

- (a) $h(k, i) = (h'(k) + 5i) \bmod m$
- (b) $h(k, i) = (h'(k) + i^{\frac{3}{2}}) \bmod m$
- (c) $h(k, i) = (h_1(k) \times (h_2(k)^i) \bmod m$
- (d) $h(k, i) = (h'(k) + 2^{i+1}) \bmod m$
- (e) $h(k, i) = (h_1(k) \times h_2(k) + i) \bmod m$

Solution

- (a) This is just linear probing, but in fixed steps of 5 instead of fixed steps of 1. This means it has both primary and secondary clustering.
- (b) This is almost quadratic probing, except instead of steps of size i^2 we have steps of size $i^{3/2}$. It will behave in a similar way to quadratic probing, namely that two keys having the same hash will have the same probe chain, since the step size does not depend on the key, but keys with different hashes will have non-overlapping chains, because the step size increases each step. So we do not have primary clustering, but we do have secondary.
- (c) This is almost double hashing, but multiplying by $h_2(k)^i$ instead of adding $h_2(k) * i$. It will display neither primary nor secondary clustering.
- (d) This hash function probes with ever increasing step sizes, in a similar way to quadratic probing, but the step sizes are increasing powers of a given number. Two items which have the same value for $h'(k, 0)$ will have the same probe chains, so this hash displays secondary clustering. It is possible for items with different values for $h'(k, 0)$ to have partially overlapping chains, for example when $h'(k_1) + 2 = 3, h'(k_2) + 2 = 9$. Then every second position in the probe chain of k_1 will overlap with an element from the chain from k_2 . This is not as bad as primary clustering (and will occur very rarely).
- (e) If we define a new hash function, $h_3(k, i) = h_1(k) * h_2(k) + i$, we notice that this is just linear probing. So it has both primary and secondary clustering.

Problem 8. You are given a set of distinct keys x_1, x_2, \dots, x_n .

- (a) Design an algorithm that creates a binary search tree of minimal height containing those keys in $O(n \log(n))$ time.
- (b) Prove that your algorithm produces a BST of height at most $\log(n)$.
- (c) Prove that the fastest algorithm for this problem in the comparison-based model has time complexity $\Theta(n \log(n))$.

Solution

- (a) First we sort the keys, this takes $O(n \log(n))$. We take the median element of the sequence and insert it into a BST, then recursively do the same thing with the resulting left and right sublists. In the below code, we assume we have a *Node* class which contains a key, a left child, and a right child. We must sort the list before calling `OPTIMAL_BST(x[1..n])`.

```

1: function OPTIMAL_BST( $x[1..n]$ )
2:   if  $x$  is empty then return null
3:    $mid = \lfloor n/2 \rfloor$ 
4:    $root = \text{Node}(x[mid])$ 
5:    $root.left = \text{OPTIMAL\_BST}(x[1..mid-1])$ 

```

```

6:   root.right = OPTIMAL_BST(x[mid+1..n])
7:   return root
8: end function

```

- (b) To prove that the height of the constructed tree is at most $\log(n)$, we write a recurrence for the height.

$$H(n) = 1 + H\left(\left\lceil \frac{n-1}{2} \right\rceil\right) \leq 1 + H\left(\frac{n}{2}\right)$$

Using telescoping,

$$\begin{aligned}
H(n) &\leq 1 + \left[1 + H\left(\frac{n}{4}\right)\right], \\
&\leq 2 + \left[1 + H\left(\frac{n}{8}\right)\right], \\
&\leq \dots, \\
&\leq k + H\left(\frac{n}{2^k}\right).
\end{aligned}$$

We know that $H(1) = 0$, so setting $k = \log(n)$, we obtain

$$H(n) \leq \log(n) + H(1) = \log(n)$$

as required.

- (c) If we could construct any BST from n keys in faster than $O(n \log(n))$ time, we could then do an in-order traversal of the resulting BST to obtain the keys in sorted order. This means we could sort faster than $O(n \log(n))$, but $\Omega(n \log(n))$ is a lower bound for comparison-based sorting, so such an algorithm is impossible.

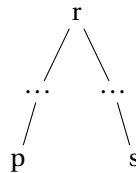
Problem 9. We know that when inserting an element into a binary search tree, there is only one valid place to put that item in the tree. Let's prove this fact rigorously. Let T be a binary search tree and let x be an integer not contained in T . Prove that exactly one of the following statements is true:

- The successor of x has no left child.
- The predecessor of x has no right child.

That is, prove that it cannot be the case that both of these statements are true simultaneously, nor that both of these statements are false simultaneously. This implies that there is a unique insertion point for x since upon insertion x must either become the left child of its successor or the right child of its predecessor.

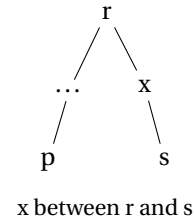
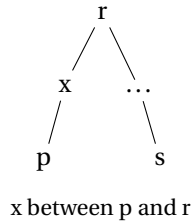
Solution

Let the predecessor of x be p , and the successor be s . We know that $p < x < s$, and that there is no key k in the tree such that $p < k < x$ or $x < k < s$. We claim that either s is an ancestor of p , or p is an ancestor of s . Suppose for contradiction that neither is an ancestor of the other, and therefore p and s have a common ancestor r .



Then we know that $p < r < s$. But that means that either $p < r < x < s$ or $p < x < r < s$. But by

the definition of successor and predecessor, both of those are impossible. So we have a contradiction, meaning that either s is an ancestor of p (in which case s has a left child, since p will be in its left subtree), or p is an ancestor of s (in which case p has a right child, since s will be in its right subtree). So we know that at least one of the statements is false, which means that at most, one statement can be true.



We still need to prove that both statements cannot be false, i.e. that at least one statement is true. Either p is an ancestor of s or vice versa, so we can consider one of the two symmetric cases. Assume p is an ancestor of s .

Since s is in the right subtree of p , p has a right child. Let's assume for contradiction that s has a left child. The value of this left child would be larger than p (since it is in the right subtree of p) but smaller than s by the definition of BSTs. This is a contradiction (as there can be no values in between the predecessor and successor of x), so both statements cannot be false, which means at least one is true.

Since at least one is true, at most one is true, exactly one must be true.

Supplementary Problems

Problem 10. Implement hash tables that use chaining, linear probing, and quadratic probing for collision resolution. You may reuse hashtable code from previous units if you have it. Compare these hash tables on randomly generated integers and compare their performance. Try adjusting the hash function, the table size, and the number of keys inserted and see how this affects the results.

Problem 11. In lectures we claimed that AVL trees are good because the balance property guarantees that the tree always has height $O(\log(n))$. Let's prove this.

- Write a recurrence relation for $n(h)$, the **minimum** number of nodes in an AVL tree of height h . [Hint: It should be related to the Fibonacci numbers]
- Find an exact solution to this recurrence relationship in terms of the Fibonacci numbers¹. [Hint: Compare the sequence with the Fibonacci sequence, find a pattern, and then prove that your pattern is right using induction]
- Prove using induction that the Fibonacci numbers satisfy $F(n) \geq 1.5^{n-1}$ for all $n \geq 0$
- Using (a), (b), and (c), prove that a valid AVL tree with n elements has height at most $O(\log(n))$

Solution

Consider an AVL tree of height h whose root node has the two subtrees L and R . Suppose that L and R have the same height. We could remove the nodes from the lowest level of one of the two subtrees, making them differ in height by one and it would still be a valid AVL tree with the same height. Therefore an AVL tree with the fewest possible nodes must have L and R differing in height by one, hence they must have

¹For this problem, index the Fibonacci numbers from zero with the base case $F(0) = F(1) = 1$.

heights $h-1$ and $h-2$. Lastly, note that L and R must also be AVL trees with the fewest possible nodes. Therefore the minimum number of nodes in an AVL tree satisfies the following recurrence relationship:

$$n(h) = \begin{cases} 1 & \text{if } h = 0, \\ 2 & \text{if } h = 1, \\ 1 + n(h-1) + n(h-2) & \text{if } h > 0. \end{cases}$$

This looks suspiciously similar to the Fibonacci sequence, so let's compare them:

h	0	1	2	3	4	5	6	7	8
$F(h)$	1	1	2	3	5	8	13	21	34
$n(h)$	1	2	4	7	12	20	33	54	88

This table strongly suggests the pattern

$$n(h) = F(h+2) - 1$$

Let's prove this by induction. We have $n(0) = F(2) - 1 = 1$ and $n(1) = F(3) - 1 = 2$ as required. Now suppose $n(h) = F(h+2) - 1$ for all $h \leq k$ for some value of k . We need to show that $n(k+1) = F(k+3) - 1$. From the recurrence, we have

$$n(k+1) = 1 + n(k) + n(k-1).$$

Invoking the inductive hypothesis, we can write

$$\begin{aligned} n(k+1) &= 1 + F(k+2) - 1 + F(k+1) - 1, \\ &= F(k+2) + F(k+1) - 1 \end{aligned}$$

By the definition of Fibonacci numbers, $F(k+2) + F(k+1) = F(k+3)$, hence we have

$$\begin{aligned} n(k+1) &= F(k+2) + F(k+1) - 1, \\ &= F(k+3) - 1, \end{aligned}$$

as required. Hence by strong induction on k , we have $n(h) = F(h+2) - 1$ for all h .

Next, we prove that the Fibonacci numbers have the desired bound. In the base case, we have

$$F(0) = 1 \geq 1.5^{-1}, \quad F(1) = 1 \geq 1.5^0.$$

Now, suppose that $F(n-1) \geq 1.5^{n-2}$ and $F(n) \geq 1.5^{n-1}$. We need to prove that $F(n+1) \geq 1.5^n$. From the definition of Fibonacci numbers, we have

$$F(n+1) = F(n) + F(n-1),$$

using the inductive hypothesis, we can write the bound

$$\begin{aligned} F(n+1) &\geq 1.5^{n-1} + 1.5^{n-2}, \\ &= 1.5 \times 1.5^{n-2} + 1.5^{n-2}, \\ &= (1.5 + 1)1.5^{n-2}, \\ &\geq 2.25 \times 1.5^{n-2}, \\ &= (1.5)^2 1.5^{n-2}, \\ &= 1.5^n. \end{aligned}$$

Therefore by induction on n , we have $F(n) \geq 1.5^{n-1}$.

Finally, combining the above, we have $n(h) = F(h+2) - 1 \geq 1.5^{h+1} - 1$. Rearranging, we find $1.5^{h+1} \leq n(h) + 1$, and hence

$$h+1 \leq \log_{1.5}(n(h)+1) \implies h = O(\log(n)),$$

as required.

Problem 12. (Advanced) Consider a hashtable implementing linear probing, with size $m = 17$ using the hash function

$$h(k) = (7k + 11) \bmod m.$$

Give a sequence of keys to insert into the table that will cause its worst-case behaviour.

Solution

We want a sequence of keys that will all hash to the same slot, which will cause linear probing to probe the entire cluster every time. Let's design a set of keys that all hash to slot 0. We want

$$7k + 11 \equiv 0 \bmod 17.$$

Adding 6 to both sides, this is equivalent to

$$7k \equiv 6 \bmod 17.$$

Now we want to divide out the 7 to get an expression for k . To do so, we need to multiply by some number x such that $7x \equiv 1 \bmod 17$ (in other words, we need the modular multiplicative inverse of 7 mod 17). We can simply find this by trial and error. Looking at multiples of 7, we find

$$\begin{aligned} 1 \times 7 &= 7 \equiv 7 \bmod 17, \\ 2 \times 7 &= 14 \equiv 14 \bmod 17, \\ 3 \times 7 &= 21 \equiv 4 \bmod 17, \\ 4 \times 7 &= 28 \equiv 11 \bmod 17, \\ 5 \times 7 &= 35 \equiv 1 \bmod 17. \end{aligned}$$

Therefore we find that 5 is the inverse we are looking for. We therefore can write

$$5 \times 7k \equiv 5 \times 6 \bmod 17,$$

from which we get

$$k \equiv 13 \bmod 17.$$

Therefore, a worst-case sequence of keys would be

$$k = 13 + 17i, \quad i \geq 0,$$

i.e. $k = 13, 30, 47, 64, 81, \dots$

Problem 13. (Advanced) Suppose that we insert n keys into a hashtable with m slots using a totally random hash function. What is the expected number of pairs of colliding elements? The pairs (k, k') and (k', k) are considered the same and should not be counted twice.

Solution

First, we define the following.

$$I_{i,j} = \begin{cases} 1 & \text{if } h(k_i) = h(k_j), \\ 0 & \text{otherwise.} \end{cases}$$

We call this an *indicator function* on the condition $h(k_i) = h(k_j)$. For each item k_i , the number of elements it collides with is given by

$$\# \text{ collisions with } k_i = \sum_{j=1}^{i-1} I_{i,j}.$$

The total number of collisions is therefore

$$\# \text{ collisions} = \sum_{i=1}^n \sum_{j=1}^{i-1} I_{i,j}$$

So, by linearity of expectation, the expected number of collisions is

$$\mathbb{E} \left[\sum_{i=1}^n \sum_{j=1}^{i-1} I_{i,j} \right] = \sum_{i=1}^n \sum_{j=1}^{i-1} \mathbb{E}[I_{i,j}].$$

We have by the definition of expectation

$$\mathbb{E}[I_{i,j}] = 1 \times \Pr[h(k_i) = h(k_j)] + 0 \times \Pr[h(k_i) \neq h(k_j)] = \Pr[h(k_i) = h(k_j)].$$

Since the hash is totally random,

$$\mathbb{E}[I_{i,j}] = \Pr[h(k_i) = h(k_j)] = \frac{1}{m}.$$

Therefore the expected number of colliding pairs is

$$\sum_{i=1}^n \sum_{j=1}^{i-1} \frac{1}{m} = \binom{n}{2} \frac{1}{m} = \frac{n(n-1)}{2m}.$$

Problem 14. (Advanced) Recall the algorithm for insertion into a hashtable using Cuckoo hashing. Assume that $m \geq 2n$ and $\text{MaxLoop} = O(\log(n))$. Given that the probability that an insertion triggers a rebuild is $O(n^{-2})$, and insertions succeed in expected constant time when a rebuild is not triggered:

- (a) Prove that the probability that a sequence of n insertions succeeds is $1 - O\left(\frac{1}{n}\right)$
- (b) Prove that the expected time complexity of insertion is $O(1)$

Solution

- (a) The probability that an insertion triggers a rebuild is $O(n^{-2})$. The probability that any of the n insertions triggers a rebuild is

$$\Pr\left(\bigcup_{i=1}^n \{\text{insertion } i \text{ triggers a rebuild}\}\right).$$

The events are not mutually exclusive, but the probability of their union is bounded above by the sum of their individual probabilities. Formally, this is called the *union bound*, or Boole's inequality.

$$\begin{aligned} \Pr\left(\bigcup_{i=1}^n \text{insertion } i \text{ triggers a rebuild}\right) &\leq \sum_{i=1}^n \Pr(\text{insertion } i \text{ triggers a rebuild}), \\ &\leq nO(n^{-2}), \\ &= O\left(\frac{1}{n}\right). \end{aligned}$$

So the probability that a rebuild is not triggered is $1 - O\left(\frac{1}{n}\right)$ as required.

- (b) Let X denote the time taken by an insertion. A successful insertion takes $O(1)$, and an unsuccessful insertion takes $O(\text{MaxLoop})$ to try, plus the cost of n subsequent insertions to perform the rebuild.

We have

$$E[X] = \begin{cases} O(1) & \text{if insertion succeeds,} \\ O(\text{MaxLoop}) + nE[X] & \text{if insertion fails.} \end{cases}$$

Therefore

$$\begin{aligned} E[X] &= \Pr(\text{Success}) \times O(1) + \Pr(\text{Failure}) \times (O(\log(n)) + nE[X]), \\ &= [1 - O(n^{-2})] \times O(1) + O(n^{-2}) (O(\log(n)) + nE[X]), \\ &= O(1) + O\left(\frac{\log(n)}{n^2}\right) + O\left(\frac{1}{n}\right)E[X], \\ &= O(1) + O\left(\frac{1}{n}\right)E[X]. \end{aligned}$$

Rearranging then gives

$$\begin{aligned} E[X] - O\left(\frac{1}{n}\right)E[X] &= O(1), \\ E[X] \left[1 - O\left(\frac{1}{n}\right)\right] &= O(1), \end{aligned}$$

and hence

$$E[X] = \frac{O(1)}{1 - O\left(\frac{1}{n}\right)} = O(1),$$

as required.

Problem 15. (Advanced) In this problem we consider a solution to the *static hashing* problem. We are given a set of n keys and want to build a hashtable on those keys that supports $O(1)$ worst-case lookup. No further keys will be added or removed.

A strategy similar to hashing with chaining is FKS hashing². In FKS hashing, we implement a two-level hashtable. The top-level hashtable of size $m = n$ stores colliding elements in inner hashtables. To support $O(1)$ lookup, we want the inner hashtables to have no collisions.

We start out by selecting a random hash function $h(k)$ for the top-level hashtable from a *universal family*. This means that we have for all $x \neq y$,

$$\Pr_{h \in \mathcal{H}} \{h(x) = h(y)\} \leq \frac{1}{m}.$$

We compute the hash values for all n keys and count the number of collisions in each slot. The inner hashtable for slot v is initialised with size n_v^2 where n_v is the number of keys that hashed to v . If the total size of all inner hashtables is too large, we pick a new hash function and try again.

Once we are happy with the total size of the hashtable, we select hash functions for the inner hashtables (also randomly from a universal family). If any of them have collisions, we select a new hash function for those until there are no longer any collisions.

It remains to prove that finding suitable hash functions can be done fast in expectation. Let's show that building the inner tables can be done fast.

- Prove that the probability that a collision occurs in a given inner hashtable of size $m_v = n_v^2$ containing n_v elements is less than $1/2$.
- Hence argue that we can find a hash function that yields no collisions on this inner table in a constant number of tries in expectation.

²Fredman, Komlós, Szemerédi, Storing a Sparse Table with $O(1)$ Worst Case Access Time, *Journal of the ACM* 1984

Next, let's show that we can find a hash function for the outer table fast.

(c) Argue that

$$\sum_{v=0}^{m-1} n_v^2 = \sum_{i=1}^n \sum_{j=1}^n I_{i,j},$$

where $I_{i,j} = 1$ if $h(k_i) = h(k_j)$ or 0 otherwise.

(d) Using (c), prove that

$$\mathbb{E} \left[\sum_{v=0}^{m-1} n_v^2 \right] \leq 2n.$$

(e) Using (d) and the fact that $\Pr[X \geq a] \leq \frac{\mathbb{E}[X]}{a}$ (this is called *Markov's inequality*), deduce that

$$\Pr \left[\sum_{v=0}^{m-1} n_v^2 \geq 4n \right] \leq \frac{1}{2}.$$

(f) Use (f) to deduce that we can find a hash function that yields a total table size of $O(n)$ in a constant number of tries in expectation.

Now we have all of the ingredients to conclude the analysis.

(g) Finally, deduce that we can build the entire hashtable in $O(n)$ expected time.

Solution

(a) We want to bound the probability that a collision occurs in one of the inner tables. Recall the concept of an indicator random variable from problem 13.

$$I_{i,j} = \begin{cases} 1 & \text{if } h_v(k_i) = h_v(k_j), \\ 0 & \text{otherwise.} \end{cases}$$

Using a union-bound (see Problem 14), and the fact that h_v is from a universal family, the probability that any collision occurs is

$$\begin{aligned} \Pr \left[\bigcup_{i < j} \{I_{i,j} = 1\} \right] &\leq \sum_{i < j} \Pr[I_{i,j} = 1], \\ &= \sum_{i < j} \frac{1}{n_v^2}, \\ &= \binom{n_v}{2} \frac{1}{n_v^2}, \\ &= \frac{n_v(n_v - 1)}{2} \frac{1}{n_v^2}, \\ &\leq \frac{1}{2} \end{aligned}$$

as required.

(b) In a sequence of Bernoulli trials, if we have a probability p of success and $1 - p$ of failure, then the expected number of trials to obtain a success is $\frac{1}{p}$. By (a), the probability that no collisions occur is at least $\frac{1}{2}$, so the expected number of trials is at most $1/\frac{1}{2} = 2$, which is constant.

(c) First, note that

$$n_v = \sum_{\substack{i=1 \\ h(k_i)=v}}^n 1,$$

and hence

$$\begin{aligned}
n_v^2 &= \left(\sum_{\substack{i=1 \\ h(k_i)=v}}^n 1 \right) \left(\sum_{\substack{j=1 \\ h(k_j)=v}}^n 1 \right), \\
&= \sum_{\substack{i=1 \\ h(k_i)=v}}^n \sum_{\substack{j=1 \\ h(k_j)=v}}^n 1, \\
&= \sum_{i=1}^n \sum_{j=1}^n \left(\begin{cases} 1 & \text{if } h(k_i) = h(k_j) = v \\ 0 & \text{otherwise} \end{cases} \right).
\end{aligned}$$

Therefore

$$\begin{aligned}
\sum_{v=0}^{m-1} n_v^2 &= \sum_{v=0}^{m-1} \left(\sum_{i=1}^n \sum_{j=1}^n \left(\begin{cases} 1 & \text{if } h(k_i) = h(k_j) = v \\ 0 & \text{otherwise} \end{cases} \right) \right), \\
&= \sum_{i=1}^n \sum_{j=1}^n \left(\begin{cases} 1 & \text{if } h(k_i) = h(k_j) \\ 0 & \text{otherwise} \end{cases} \right), \\
&= \sum_{i=1}^n \sum_{j=1}^n I_{i,j}.
\end{aligned}$$

(d) Using (c), by linearity of expectation,

$$\begin{aligned}
\mathbb{E} \left[\sum_{v=0}^{m-1} n_v^2 \right] &= \mathbb{E} \left[\sum_{i=1}^n \sum_{j=1}^n I_{i,j} \right], \\
&= \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[I_{i,j}], \\
&= \sum_{i \neq j} \mathbb{E}[I_{i,j}] + \sum_{i=1}^n \mathbb{E}[I_{i,i}], \\
&= \sum_{i \neq j} \Pr(h(k_i) = h(k_j)) + \sum_{i=1}^n \Pr(h(k_i) = h(k_i)), \\
&\leq \sum_{i \neq j} \frac{1}{m} + \sum_{i=1}^n 1, \\
&= \frac{n(n-1)}{m} + n, \\
&\leq 2n,
\end{aligned}$$

as required. Note that we used the fact that $\Pr(h(k_i) = h(k_j)) \leq \frac{1}{m}$ for $i \neq j$, and otherwise $\Pr(h(k_i) = h(k_i)) = 1$.

(e) Using (d), by Markov's inequality

$$\begin{aligned}
\Pr \left[\sum_{v=0}^{m-1} n_v^2 \geq 4n \right] &\leq \frac{\mathbb{E} \left[\sum_{v=0}^{m-1} n_v^2 \right]}{4n}, \\
&\leq \frac{2n}{4n}, \\
&= \frac{1}{2},
\end{aligned}$$

as required.

- (f) We again use the fact that if we have Bernoulli trials with probability of success p , then the expected number of trials until a success is $\frac{1}{p}$. From (e), we have that

$$\Pr \left[\sum_{v=0}^{m-1} n_v^2 \geq 4n \right] \leq \frac{1}{2}.$$

which means that the probability that the total table size is **less** than $4n$ is at least $\frac{1}{2}$. So $p \geq \frac{1}{2}$, hence we can find a hash function that yields $O(n)$ space in two tries in expectation, which is constant.

- (g) Each internal table can be built with a constant number of selections of hash functions, in expectation. So we will need to insert each of the n items a constant number of times, in expectation, which is $O(n)$ time. The total size of the outer table will be $O(n)$ after a constant number of tries, in expectation. So the expected total amount of work is $O(n)$, as required.