

Skip Lists

This is yet another data structure that allows for inserts, searches and deletes in an average of $O(\log n)$ time, where n is the total number of items in the skip list. Interestingly enough, this data structure is realized with a mesh of linked lists.

Unlike some of the other data structures we have looked at, a skip list uses randomization. Thus, if you give the exact same input to a skip list twice, you **MAY** actually get different behavior both times. (An example of an algorithm you have seen before that is randomized is Quick Sort when the partition element is picked randomly.)

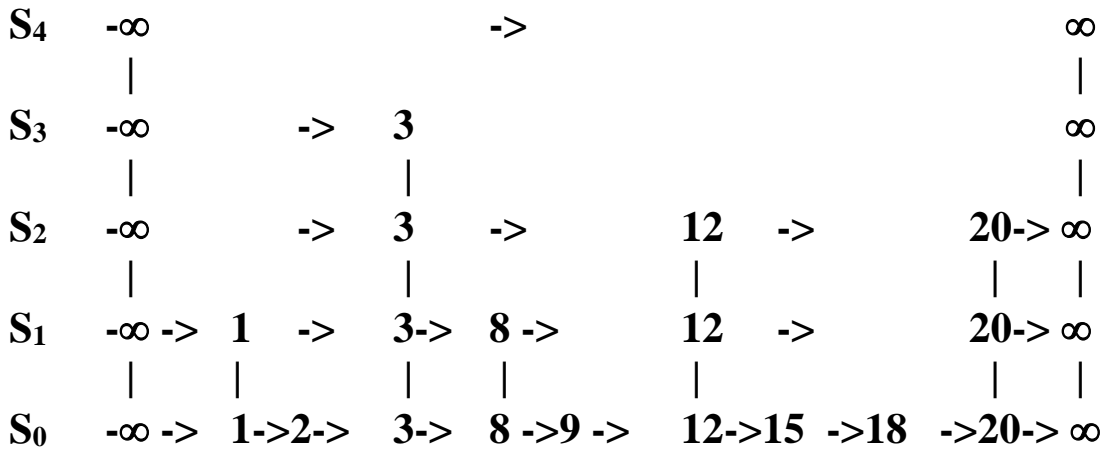
Skip List Definition

A skip list consists of a series of lists $\{S_0, S_1, \dots, S_h\}$. The list S_0 starts with the special value $-\infty$ and continues with all the numbers in numerical order and is terminated by the special value $+\infty$. (Note: These two special values, $-\infty$ and $+\infty$ are always stored in the first and last nodes respectively of every list.)

Each list S_i stores a subset of the values stored in S_{i-1} . You can visualize list S_i to appear directly above list S_{i-1} . Furthermore, each list will be connected to adjacent lists through vertical links between lists. In particular, if a value appears in both in S_i and S_{i-1} , then there will be a link between these two nodes.

Finally, in S_h will only contain $-\infty$ and $+\infty$.

Here is a visual representation of a skip list:



So how do we decide what elements get repeated in higher lists?

Whenever we insert an element, we will first insert the element in the list S_0 . (Although it may look like it, this will NOT take $\theta(n)$ time.) Once this is done, we will randomly pick a real number in between 0 and 1. If it is greater than .5, we will copy this value in S_1 . This involved creating a new node, linking that new node to the node with the same value in S_0 , and then linking the node in between the two appropriate values in S_1 . Now, we will continue and create another random real number in between 0 and 1. If it is greater than .5 we'll repeat this process for S_2 . We continue until we get a random number less than .5 and stop at that point. So, in the worst case, this algorithm could go forever, but technically speaking, the probability of that occurring is 0.

So what does this do?

It leaves about $1/2$ as many nodes in S_i as in S_{i-1} . Thus, if the skip list is storing n values, there should roughly be $n/2$ values in S_1 . Using the repeated halving principle, we see that the expected

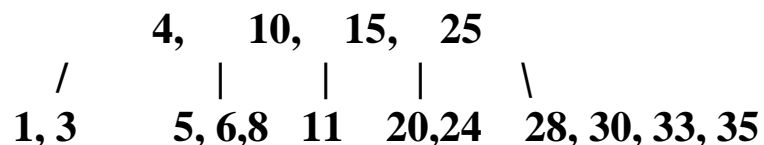
“height” of the skip tree is logarithmic in n . In order to maintain a skip list, we will have to provide four operations for any node in the list in $O(1)$ time. These are:

- 1) after(p)
- 2) before(p)
- 3) below(p)
- 4) above(p)

The first two provide a way to navigate between elements in a list, and the last two provide a way to navigate between lists. In order to implement a skip list with these features, a connected mesh of nodes with four separate references would suffice.

Searching

We need to efficiently use the setup of each list to help us find a value. To solve this problem, consider searching in a Multi-way Tree. We will go through the list of values at a node until we have gone too far. When we have, we know to traverse the link “right before” the node that stored the value that was too big:



If we are searching for 23, we will look at 4, then 10, then 15. When we hit 25, we know if 23 is in the tree, it must be in the subtree to the right of the 15 and the left of the 25. Follow this link down. Now, after you look at the 20 and then the 24, you would search down the subtree to the right of 20 and left of 24, but this is null. Thus, the value isn’t in the tree.

We will utilize this same idea in searching for a value in a skip list.

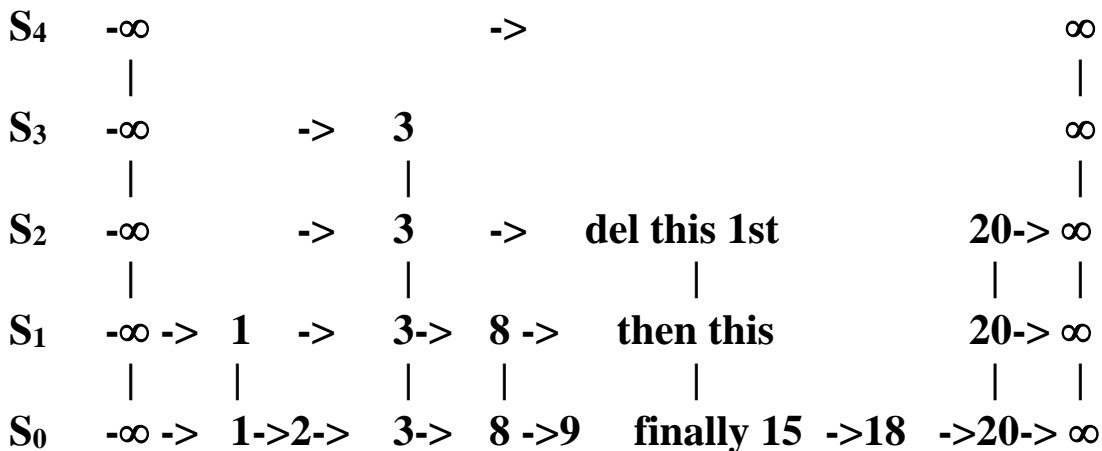
$2(n+\log n)$ using the sum of an infinite geometric series. (Note: the $2\log n$ accounts for the $-\infty$ and $+\infty$ nodes in each of the approximately $\log n$ lists.)

Removal

Once again, implement the search procedure. In the book they make the search procedure go down all the way to S_0 , but this isn't entirely necessary. (The only situation this would be necessary is if you only stored entire records in S_0 and just stored keys to those records in the list above. This idea has the potential of saving some space, but also is more complicated since you'd have a linked mesh structure with different types of nodes.)

Instead, you will find a node in the highest list it appears first if you implement the search algorithm. (In order to stop the search algorithm at this point, change the pseudocode in the book to read "while (below(p) != null && key(p) != k)")

If you do this, you can simply delete each desired node, one by one until you've deleted the node storing the value to be deleted in S_0 . For example, if you were to remove 12, you would do the following steps:



Maintaining the Top-most Level

You must always have a reference to the top-most left-most node. (This is the $-\infty$ on list S_h .) One way to control this level is to fix the maximum height of the skip list based on the number of elements in the list. You could do something like $\max(10, 2\log n)$, where n is the number of elements stored. This keeps the height logarithmic. (Thus, if you inserted an element, and got $2\log n$ consecutive random numbers over $.5$, which would indicate to "grow" a list beyond $S_{2\log n}$, you don't do it. In essence, when you get to the top row, you don't flip a coin to see if you will add a new row.)

Or, you could simply allow the height of the tree to expand as inserted elements "build" towers. (Though it does seem silly to have a single tower that's more than one element taller than all the rest.) Either way, the expected amount of time for the three operations is $O(\log n)$.

Bounding the height of a Skip List

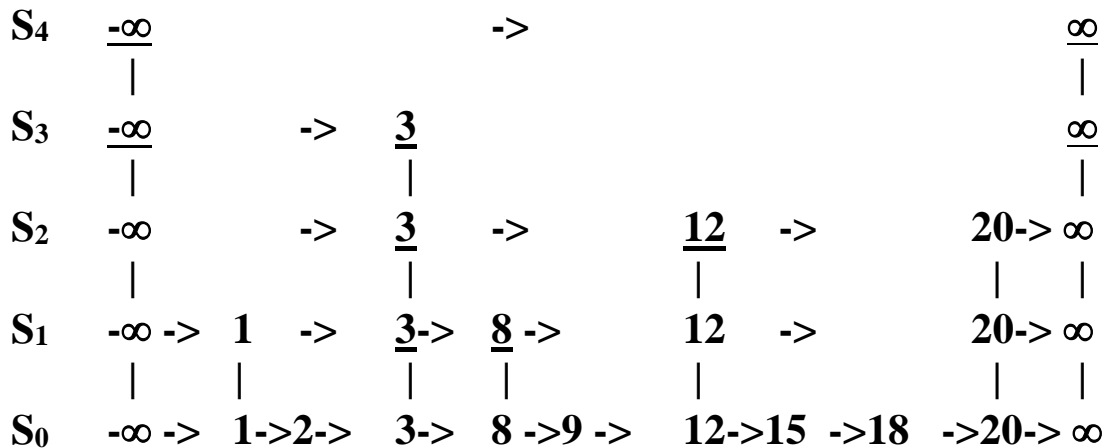
The probability that an element gets to S_i upon insertion is $(1/2)^i$. Given that there are n items in the list, the maximum probability that at least one of those items is on level i is $n/2^i$. (To see this, consider the inclusion-exclusion principle for probability. We can simply add the probabilities of each of the events, but when we do this, we have "overcounted" by counting the probability that multiple items end up in S_i several times. Thus, adding the probabilities as we have done is an overestimate of the actual probability, since the events are NOT mutually exclusive.)

Thus, the height of the skip list does NOT exceed $c \log n$ with probability $n/n^c = 1/n^{c-1}$. The book plugs in $c=3$ to show that the probability of a skip list exceeding the height $3 \log n$ is at most $1/n^2$.

Analyzing Search Time in a Skip List

Basically, each "move" in the search algorithm takes several steps down and several steps forward. The number of steps down is bounded by the height of the list, which is expected to be $O(\log n)$ as shown above. Now, we must calculate the expected number of forward steps.

Note that each new element examined at level i of the list could not have existed in level $i+1$. This is because we drop levels in the skip list if we've "gone too far" in the level above. When we drop down, the new nodes we will traverse going forward will NOT be in the level above, because we would have stopped then.



If we are searching for 8 above, we know that since 12 was too big in S_2 , it's not one of the new elements we'll see in S_2 . In particular, each of the underlined elements above denotes the "search ranges" for each list. We you drop down a list, you are confining yourself to the search region between adjacent elements in the list above. What is the expected number of elements on S_{i-1} in between adjacent elements from S_i ? Since an element is twice as likely to be in S_{i-1} as S_i , the density of elements is twice as much. This would infer an average of one extra element in between two adjacent ones from the previous list. Certainly this value can be bounded by 2 new elements. Thus the number of elements searched in each level of the skip list is $O(1)$. It follows that the expected search time in a skip list is $O(\log n)$, where the list stores n values.

The expected number of physical nodes used in a skip list is at most $2n$. Can you prove this?