

Final Review

Heaps

- Motivation: What if we always want to find the minimum or maximum element?
- Keep high priority items at the top
 - Min heap: high priority corresponds to low priority value
 - Max heap: high priority corresponds to high priority value
 - Notice the difference between priority and priority value!
- Represented as a binary tree with two more properties:
 - Complete: no empty spaces other than on the right-hand side of the bottommost level; height will be $\Theta(\log N)$ where N is the number of nodes
 - (Min/Max)-heap property: for a particular node n , the children of n must have (greater/-lesser) priority value than n ; the root will always contain the (lowest/highest) priority value element
- Methods
 - `peek()`: returns (but does not remove) the item with the highest priority; runtime is $\Theta(1)$
 - `removeMin()`: returns (and does remove) the item with the highest priority; runtime is $O(\log N)$
 - * Take the item in the bottom-rightmost position and replace the value at the root
 - * Bubble down the new root value
 - `insert(T item, int priorityVal)`: Insert the item with priority value of `priorityVal` into the heap; runtime is $O(\log N)$
 - * Insert the item in the first open bottom-left position
 - * Bubble up the new inserted value
- Bubbling (in a **min-heap**)
 - Bubble up: while the priority value of a particular node n is less than the priority value of its parent, swap the two
 - Bubble down: while the priority value of a particular node n is greater than the priority value of its child/children, swap the two (pick the lesser of the children if both have priority value less than n)

- Representation
 - Number each element in the heap, starting from 1, left to right top to bottom, this will represent the index of the item in the array
 - For a particular node at index i :
 - * Parent is at index $\frac{i}{2}$
 - * Left child is at index $2i$
 - * Right child is at index $2i + 1$
- PriorityQueue<T>
 - Implemented with a min heap, methods include T poll(), T peek(), void push(T item)
 - Can use own Comparator object to change how the PriorityQueue organizes elements

Practice Problems

1. What is the size of the largest binary min heap that is also a valid BST? Draw an example assuming all keys are letters from the English alphabet. **Assume the heap has no duplicate elements**
2. What is the size of the largest binary max heap that is also a valid BST? Draw an example assuming all keys are letters from the English alphabet. **Assume the heap has no duplicate elements.**
3. What is the size of the largest binary min heap in which the fifth largest element is a child of the root? You do not need to draw an example. **Assume the heap has no duplicate elements.**

Graph Traversals

BFS/DFS

Idea: We must look through all the values of our graph. So, given some starting point, we do a BFS/DFS traversal with the caveat that we now track what vertices we've visited before (to avoid cycles). Below is the code for DFS. For BFS, we need only replace that `Stack` fringe with a `Queue` fringe.

```
public void dfs() {
    Stack fringe = new Stack();
    Set visited = new Set();
    fringe.push(startVertex);
    while (!fringe.isEmpty()) {
        Vertex v = fringe.pop();
        if (!visited.contains(v)) {
            process(v); //Do something with v
            for (Vertex neighbor: v.neighbors) {
                fringe.push(neighbor);
            }
            visited.add(v);
        }
    }
}
```

Topological Sort

Idea: Given a directed, acyclic graph, how do we 'sort' the vertices based on their dependencies on one another?

```
public void topologicalSort() {
    Stack fringe = new Stack();
    Map currentInDegree = new Map<Vertex, Integer>();

    while (!fringe.isEmpty()) {
        Vertex v = fringe.pop();
        process(v); //Do something with v
        for (Vertex neighbor: v.neighbors) {
            currentInDegree(neighbor) -= 1; //Not actual Java code
            if (currentInDegree(neighbor) == 0) {
                fringe.push(neighbor);
            }
        }
    }
}
```

Dijkstra's Algorithm

Runtime: $O((|V| + |E|) \log(|V|))$

Important Assumption: When a vertex is removed from the fringe, we assume we have then found the shortest path to that vertex. Thus, we will no longer try to update the shortest path to that vertex; the path is finalized.

Pseudocode

- **Initializing data structures**

1. Initialize the following structures:
 - **fringe**, a priority queue ordered by a vertex's distance from the start vertex
 - **distanceMap**, a mapping between vertex and distance from the start vertex
 - **predecessorMap**, a mapping between vertex and previous vertex
2. Add the start vertex to **fringe** and **distanceMap** with distance 0
3. For all other vertices, add them to **fringe** and the **distanceMap** with distance infinity

- **While-Loop (processing the shortest paths)**

1. Pop off a vertex v from **fringe**
2. Loop over each neighbor n of v :
 - (a) Let $\text{newDistance} = \text{distanceMap}[v] + \text{edge}(v, n)$
 - (b) If $\text{newDistance} < \text{distanceMap}[n]$, then
 - (1) Update the priority value of n in **fringe** and distance in **distanceMap** to be **newDistance**
 - (2) Update **predecessorMap** so that the previous vertex of n is v

Practice Problems

1. Consider a weighted, directed graph G with distinct and positive edge weights, a start vertex s and a target vertex t . Assume that G contains at least 3 vertices and that every vertex is reachable from s .
 - (T/F) Any shortest $s \rightarrow t$ path must include the lightest edge.
 - (T/F) Any shortest $s \rightarrow t$ path must include the second lightest edge.
 - (T/F) Any shortest $s \rightarrow t$ path must exclude the heaviest edge.
 - (T/F) The shortest $s \rightarrow t$ path is unique.

2. The runtime for Dijkstra's algorithm is $O((V + E) \log V)$. However, this is specific only to binary heaps. Let's provide a more general runtime bound. Suppose we have a priority queue backed by some arbitrary heap implementation. Given that this unknown heap contains N elements, suppose the runtime of remove-min is $f(N)$ and the runtime of change-priority is $g(N)$.

(a) What is the runtime of Dijkstra's in terms of $f(V)$ and $g(V)$?

(b) Turns out the optimal version of Dijkstra's algorithm uses something called a fibonacci heap. The fibonacci heap has amortized $\Theta(1)$ change-priority, and $\Theta(\log(N))$ remove-min. What is the runtime of Dijkstra's algorithm?

Sorting Algorithms

Comparison-Based Sorts

- Selection Sort
 - Keep a sorted and unsorted portion, initialize the sorted portion to nothing and the unsorted portion to the entire list
 - ‘Select’ the minimum value from the unsorted portion and append it to the end of the sorted portion
 - Tips to identify: The minimum item will always move to the sorted portion after an iteration. The sorted portion (AKA, the beginning of the list) will hold the final ordering of the elements at every iteration.
- Insertion Sort
 - Inversion: any two elements that are in the wrong place relative to each other
 - Insertion sort sorts elements by removing all inversions so the list is sorted
 - Keep a sorted and unsorted portion, initialize the sorted portion to nothing and the unsorted portion to the entire list
 - For the first element of the unsorted portion, ‘insert’ it into the correct position in the sorted list by continuously comparing and swapping with the element before it
 - Tip to identify: In intermediate steps, the front of the list will contain the items in sorted order, though it may not be the elements in the final sorted order. In other words, the sorted portion of the list grows, but does not always contain the final sorted elements.
- Merge Sort
 - Split the list in half, recursively `mergeSort` each half, then merge the halves together
 - For this class, assume is stable
 - Tip to identify: Items will not cross the middle boundary/partition until the last step.
- Quicksort
 - Choose a pivot element (pivot choice is key here!) and partition the list into elements less than the pivot and elements greater than the pivot, recursively `quickSort` each partition
 - For this class, assume is unstable
 - Tip to identify: Use the pivot specified to see if the partitions are what is expected.
- Heapsort
 - Throw all elements into a max heap, remove the root and swap places with the newly freed spot in the heap
 - Can use bottom-up heapification to heapify the elements in linear time
 - Tip to identify: Will become sorted starting from the end of the list.

Counting-Based Sorts

Radix Sort

- Sort the elements digit by digit based on the radix by running counting sort on each digit
- LSD (least significant digit):
 - Run counting sort on the rightmost digit and move towards the left
 - Used more in practice, stable sort, can be done iteratively
- MSD (most significant digit):
 - Run counting sort on the leftmost digit, and repeat towards the right
 - Harder to implement, usually done recursively which can use a lot of space on the stack, unstable if done in place due to swapping elements when reordering

Sorting Runtimes

Algorithm	Best-case Scenario	Best-case Runtime	Worst-case Scenario	Worst-case Runtime
Heap Sort				
Quicksort				
Merge Sort				
Selection Sort				
Insertion Sort				
LSD Radix Sort				
MSD Radix Sort				

Practice Problems

1. Give an example of when insertion sort is more efficient than merge sort.
2. When would you use merge sort over quicksort?
3. When would you use radix sort over a comparison sort?
4. You have an array of integers with a length of N , where the maximum number of digits is given by K . Suppose that $K \gg N$ (that is, K is much larger than N). What sort could you use?
5. Assume for the previous question that you used LSD radix sort. What is the runtime?
6. Suppose you are given an integer array with length N such that there are $N - \sqrt{N}$ copies of the minimum element located in the front of the array. What is the worst case runtime of insertion sort on this list?
7. For the same construction, what is the worst case runtime of selection sort?

Disjoint Sets

Disjoint sets are used to determine what sets, or groups, elements belong in to. Our main goal with implementing a disjoint set is to be able to check whether or not two elements are in the same set, and to be able to merge two sets together. You can also extend this feature to checking the connectivity of two elements (think of Kruskal's algorithm, for instance). Thus, there are three operations we want our disjoint set to support: `makeSet`, `find`, and `union`.

Weighted Quick-Union trees are a way to implement disjoint sets. In WQU trees, we will represent sets as a tree-like structure. Each element is itself a node, that may point to another node. The root of any tree acts as the 'representative' of the set. Let us examine how the disjoint set operations are implemented in WQU trees:

- `find(n)`: find and return the root node of n
 - **Path Compression** optimization:
As you follow the parent pointers from n to the root, reassign the parent pointers of n and all other nodes you see along the way to be the root.
- `union(a, b)`: merge the set of a with the set of b
 1. Let `root1 = find(a)`
 2. Let `root2 = find(b)`
 3. **Union by size** optimization:
If `root2` has a smaller set size than `root1`, then set the parent pointer of `root2` to point to `root1`.
Else, set the parent pointer of `root1` to point to `root2`.

Notice that both of the optimizations have the goal of preventing the tree from getting too tall.

Runtime (including both optimizations):

- One call to `find`:
Worst case: $\Theta(\log(n))$
Best case: $\Theta(1)$
- One call to `union`:
Worst case: $\Theta(\log(n))$
Best case: $\Theta(1)$
- For f calls to `find` and u calls to `union`:
Essentially $\Theta(u + f * \alpha(f + u, u))$

Where α is the inverse Ackerman function, a function that grows veeeeeeeeeeeeerrrrrrrrrryyyyyyyyyy slowly.

Minimum Spanning Trees

A *spanning tree* of a graph is a tree that connects all vertices of the graph. We define the weight of a spanning tree as the sum of the weights of all edges in the tree. A **minimum spanning tree** is a spanning tree of a graph that has the least weight. Note that any tree that connects V vertices will have $V - 1$ edges.

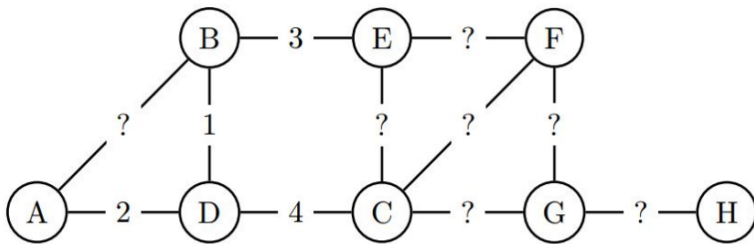
We can use Kruskal's algorithm to find the MST. The main idea of Kruskal's algorithm is to successively select edges with the smallest weight to add to our tree, until a MST is formed. Thus, given a graph G , Kruskal's algorithm works as follows:

1. Let T be an empty list of edges. This is where you will store the edges for your final MST.
2. Make a sorted list `sortedEdges` of all the edges in G , from smallest to greatest.
3. For each edge $(u \rightarrow w)$ in `sortedEdges`:
 - (a) If u and w are not already connected by the edges in T , then add the edge $(u \rightarrow w)$ to T .
 - (b) If u and w are already connected by the edges in T , then continue.

We can see that the runtime of Kruskal's algorithm is dependent on the sorting step in (2), and the check in (3a/b) to see if two vertices are already connected by T . The sorting step in (2) takes $O(E \log E)$. By using WQU trees with path compression to accomplish steps (3a/b), a single check to see if two vertices are already connected by T will take time proportional to $\alpha(V)$. Since this factor is extremely small, the total runtime of Kruskal's ends up being $O(E \log E)$ due to the sorting step.

Practice Problems

1. In the graph below, some of the edge weights are known, while the rest are unknown. List all edges that must belong to a minimum spanning tree, regardless of what the unknown edge weights turn out to be. Justify your answers briefly.



$$\text{cost}(A, D) = 2, \text{cost}(B, D) = 1, \text{cost}(C, D) = 4, \text{cost}(B, E) = 3$$

2. Given a minimum spanning tree T of a weighted graph G , describe an $O(V)$ algorithm for determining whether or not T remains a MST after an edge (x, y) of weight w is added.

3. True or False: Suppose G is a graph and T is a MST of G . If we decrease the weight of any edge in T , then T is still a MST.

4. Suppose instead of decreasing the weight of an edge in the MST, we decrease the weight of any random edge not in the MST. Give an efficient algorithm for finding the MST in the modified graph.

5. Consider a weighted, undirected graph G with distinct and positive edge weights. Assume that G contains at least 3 vertices and is connected.
 - (T/F) Any MST must include the lightest edge.
 - (T/F) Any MST must include the second lightest edge.
 - (T/F) Any MST must exclude the heaviest edge.
 - (T/F) The MST is unique.

Regular Expressions

- Special Characters

- `+` matches one or more of the preceding character (greedy operator)
- `*` matches zero or more of the preceding character (greedy operator)
- `?` matches zero or one of the preceding character (greedy operator), `?` after a greedy operator will make the operator reluctant
- `{n}` matches exactly n instances of the preceding character
- `{n, m}` matches at least n but no more than m instances of the preceding character
- `.` matches any character
- `[...]` denotes a character class, matches one character that is described within the brackets
- `(...)` denotes a capturing group, can be used to group characters together
- To match special characters (**Strings** including periods or question marks), must escape them with a backslash

- Pattern

- `Pattern compile(String s)`: returns a `Pattern` that is represented by `s`
- `Matcher matcher(String phraseToMatch)`: returns a `Matcher` that checks if `phraseToMatch` matches the `Pattern`

- Matcher

- `boolean matches()`: returns true if `Matcher`'s `Pattern` matches the given `phraseToMatch` entirely
- `boolean find()`: returns true if `Matcher`'s `Pattern` matches a substring of the given `phraseToMatch` (will only match up to the substring, subsequent calls to `find()` may result in more matches in the same string)
- `String group(int index)`: returns the `String` corresponding to the group indicated by `index`, `index 0` will correspond to the entire matched phrase

- Greedy vs. Reluctant

- Some operators (`+`, `*`, `?`) will match as much as possible
- Want to match the `'lo000l'` in `'lo000l halo'`? If our pattern is `'l.*l'`, will actually match `'lo000l hal'` because the `*` operator will match as much as possible
- Use the `?` operator after a greedy operator to make it reluctant (match as little as possible while still following the pattern described, in above example, will match `'lo000l'`)

- Capturing

- Can use `(...)` and `group` method to return certain parts of a matched substring
- `group(1)` will refer to the first grouped elements in the `String`, same idea with 2, 3, etc.

Practice Problems

1. Write the regex that matches any string that has an odd number of a's and ends with b.
2. Write a regex that matches any binary string of length exactly 8, 16, 32, or 64.
3. Write the regex that matches any string with the word "banana" in it.
4. A certain type of bracketed list consists of words composed of one or more lower-case letters alternating with unsigned decimal numerals and separated with commas:
[]
[wolf, 12, cat, 41]
[dog, 9001, cat]

As illustrated, commas may be followed (but not preceded) by blanks, and the list may have odd length (ending in a word with no following numeral.) There are no spaces around the [] braces. Write the **Java pattern** that matches an alternating list