

Greedy Algorithms

Lecture 4

Timothy Kim
GWU CSCI 6212

Sorting Lower Bound

- Theorem: Every comparison-based sorting algorithm has the worst-case running time of $\Omega(n \log n)$
- Comparison Sort: Sorting algorithm that only reads the list elements through a comparison operation that determines which of two elements should occur first in the final sorted list
- Non-comparison Sort: Bucket sort, Radix sort, and more.

Proof

- Fix a comparison-based sorting algorithm and an array of length n
- Given the input array, there are $n!$ possible input
- Suppose that the algorithm makes $\leq k$ comparisons to correctly sort any one of the $n!$ inputs.
- Across all $n!$ possible inputs algorithm exhibits $\leq 2^k$ distinct possible execution.

Proof continued

- By the pigeon hole principle
 - If $2^k < n!$, there must be two distinct input that share the same execution of the algorithm
 - Therefore $2^k \geq n! \geq (n/2)^{n/2}$
 - $k \geq (n/2) \log_2(n/2) = \Omega(n \log n)$

Greedy Algorithm

- Definition: iterative process where decision is made at the moment without access to additional information about the future nor the past.
- Fractional Knapsack
- Scheduling Algorithm
- Minimum Spanning Tree (Shortest Path)
- Maybe more

Typical Structure

```
1. Greedy( S ):
2.     P = []                                // processed
3.     while (solutionFound(P)): // usually is S empty
4.         i = selection(S)                // usually best choice
5.         P.add(i)
6.         S.remove(i)
```

D&C vs Greedy

| | Divide and Conquer | Greedy |
|-------------|----------------------------------|--------|
| Design | Hard (finding sub-problem) | Easy |
| Analysis | Hard (recursion) | Easy |
| Correctness | Easy to establish (induction) | Hard |

Proofs of Correctness

- By Induction - "greedy stays ahead"
- Exchange argument
 - By Contradiction
 - By Construction
- And more

Fractional Knapsack

- Input

$$\text{Items} = \{(w_1, v_1), (w_2, v_2), \dots, (w_n, v_n)\} | w_i > 0, v_i > 0$$

$$W = \text{weight limit}$$

- Output

$$\{x_1, x_2, \dots, x_n\} | 0 \leq x_i \leq 1$$

- such that we maximize $\sum_{i=1}^n x_i v_i$ while keeping $\sum_{i=1}^n x_i w_i \leq W$

Design

- Start with largest value
 - Sort by value
 - Gather until full
- Start with smallest weight
 - Sort by weight
 - Gather until full

Gather by value

```
1. GreedyByValue(items, W):  
2.     X = []  
3.     sort items by value  
4.     for item in items:  
5.         space = W - weight(X)  
6.         if (space > item.weight):  
7.             X += 1  
8.         else:  
9.             X += (space / item.weight)  
10.    return X
```

Is it really optimal?

1. $w = \{1, 3\}$
2. $v = \{1, 3\}$
3. $W = 2$
4. $X1 = \text{GreedyByValue}(w, v, W)$
5. $X2 = \text{GreedyByWeight}(w, v, W)$

1. $w = \{1, 3\}$
2. $v = \{3, 1\}$
3. $W = 2$
4. $X1 = \text{GreedyByValue}(w, v, W)$
5. $X2 = \text{GreedyByWeight}(w, v, W)$

Objective Function

- Ranking function that rewards low weight and high value?

Better method

- Use ratio
 - Score = $\frac{v_i}{w_i}$ (higher the better)

Proof

- Let X be our greedy solution by ratio
- Let Y be the optimal solution better than X
- Let j be the point at which X and Y must differ
- Then $X_j = 1$ and $Y_j < 1$
- By swapping X_j and Y_j , we create Y' more optimal than Y
 - *(Further justification needed)*
- Contradiction, thus there can't be any solution more optimal than X

Scheduling Problem

- Given a shared resource, there are many jobs that require it. Find the optimal sequence of jobs.

- Input

l_j = length of job j

w_j = weight or priority of job j

- Output

- Order the jobs to minimize $C = \sum_{j=1}^n w_j c_j$

- where

c_j = completion time of job j

Example

- 3 jobs with length $l_1 = 1$, $l_2 = 2$ and $l_3 = 3$
- With the weight of 1, 2, and 3
- Job order = $[1, 2, 3]$
- then $c_1 = 1$, $c_2 = 3$, $c_3 = 6$
- Total cost = ??

Objective Function

- What should we reward?
 - weight?
 - job length?

Scheduling Algorithm

1. `Schedule(L, W):`
2. `// Sort the job weight/length`
3. `// as the objective function`
4. `// Homework`

Proof of Correctness

- **Claim:** Greedy algorithm with w_j/l_j as the ordering produces the optimal job schedule.
- Proof by contradiction
 - Assume distinct score value for now.
 - Relabel the jobs in the order of our Greedy Algorithm's schedule:

$$\begin{array}{ccccccc} \frac{w_1}{l_1} & > & \frac{w_2}{l_2} & > & \dots & > & \frac{w_n}{l_n} \\ & & 1 & & 2 & & \dots & & n \end{array}$$

Proof of correctness

- Let X be the schedule that our Greedy algorithm outputs
- Let Y be the optimal schedule that's not X
- Then in Y , there must be a job schedule i, j where $i > j$.
- We swap the position of job i and j to produce schedule Y'

Proof of correctness

- By swapping,
 - c_i went up by value of l_j , therefore C went up by $w_i l_j$
 - c_j went down by value of l_i therefore C went down by $w_j l_i$
- Note
 - $i > j \rightarrow \frac{w_i}{l_i} < \frac{w_j}{l_j}$
$$w_i l_j < w_j l_i$$
 - This implies that Y' is better than Y, which is contradiction

Proof of correctness

- For the case where scores are NOT distinct
 - Proof will be the same until we compare Y and Y'
 - When transforming Y into Y' you get $w_i l_j \leq w_j l_i$
 - Which implies Y' is as good as Y or better.
- Note Y' is similar to X, so we can keep swapping to transform Y into X and it'll continue to be as good or better.
- Therefore X must be optimal.

Heap

- **max-heap:** A complete binary tree, where all nodes have key values that are greater than or equal to each of its children.
- **min-heap:** A complete binary tree, where all nodes have key values that are less than or equal to each of its children.

Heap Operations

- **Insert:** add a new object to a heap - $O(\log n)$
- **Extract Min:** remove a node in with a minimum key value - $O(\log n)$
- **Heapify:** n batched inserts - $O(n)$
- **Delete:** remove a node - $O(\log n)$

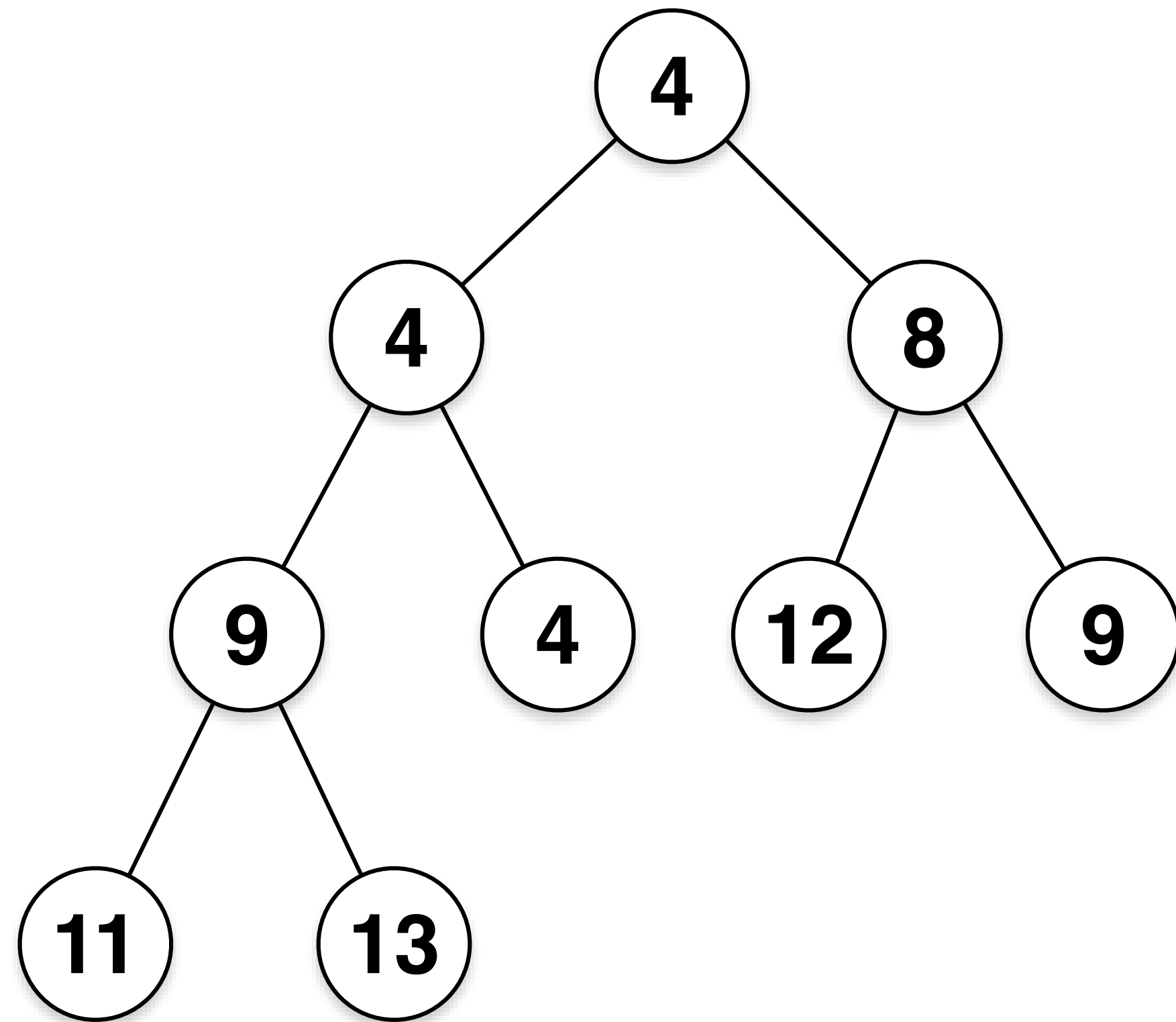
Application

- If your application requires minimum calculation repeatedly
 - **Selection Sort:** find min $O(n^2)$
 - **Heap Sort:** same as selection sort but $O(n \log n)$
 - Heapify input array - $O(n)$
 - Extract-min n times - $O(n \log n)$
- Event manager - "priority queue"
- Median maintenance - extract median

Implementation

- Usually trees are implemented using pointers and references. However because Heaps are complete binary tree, we can use an array.
- root node: 1st element
- $\text{child}(i) = (2i, 2i + 1)$
- $\text{parent}(i) = \begin{cases} i/2 & \text{if } n \text{ is even} \\ \lfloor i/2 \rfloor & \text{if } n \text{ is odd} \end{cases}$

Implementation



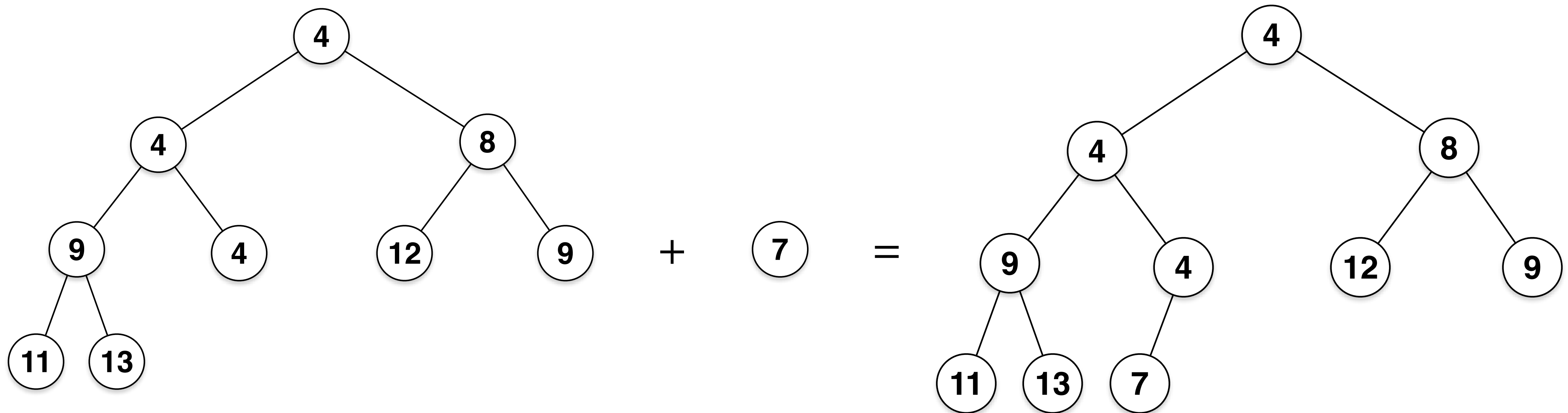
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|----|---|----|----|
| 4 | 4 | 8 | 9 | 4 | 12 | 9 | 11 | 13 |

Insert

1. **Insert(k):**
2. Add k to the next place (append)
3. if (parent(k).key < k.key):
4. done
5. else:
6. heapify-up (bubble up)

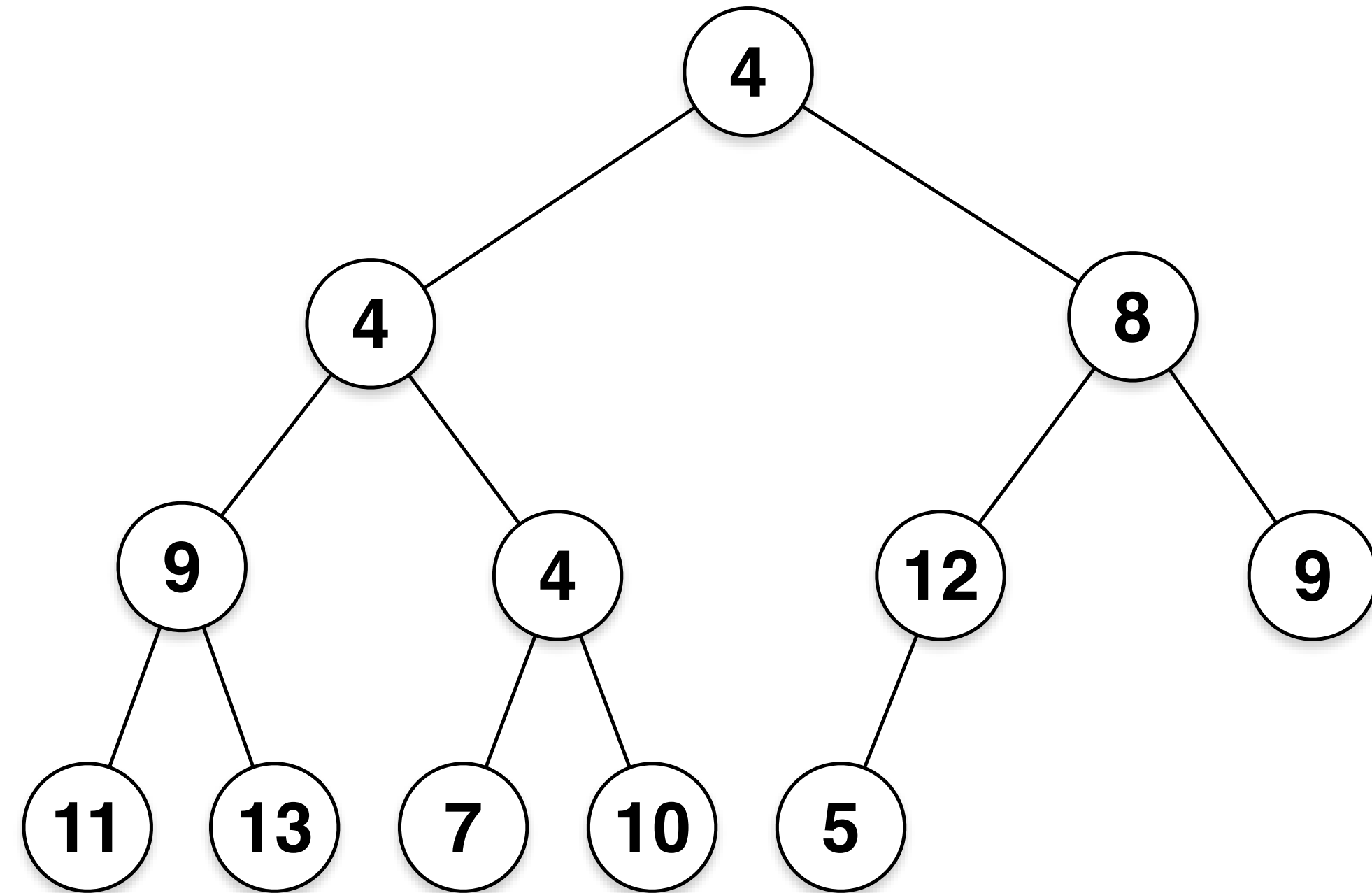
- **Invariant:** all nodes has key value that are less than or equal to each of its children

Example

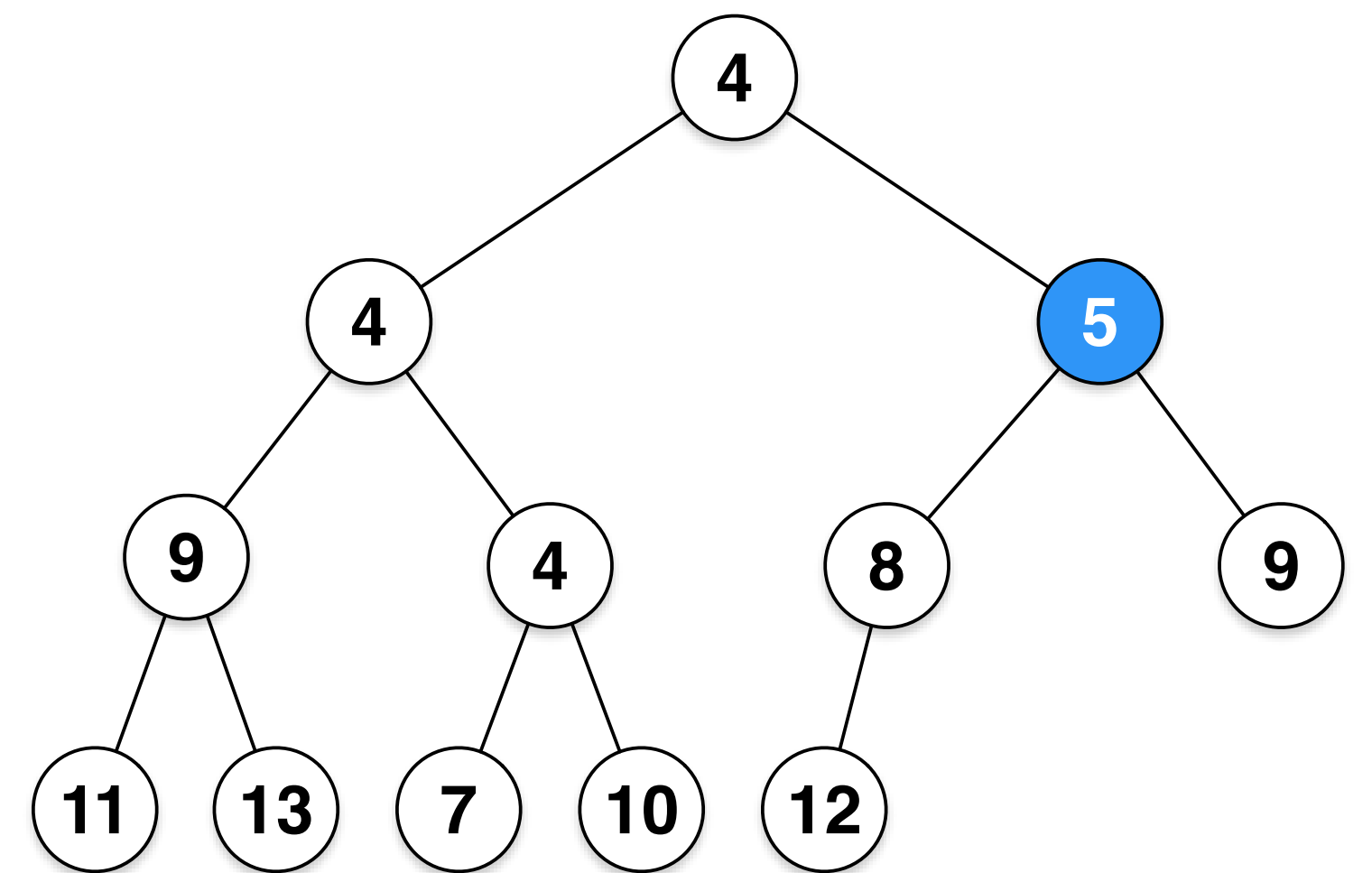
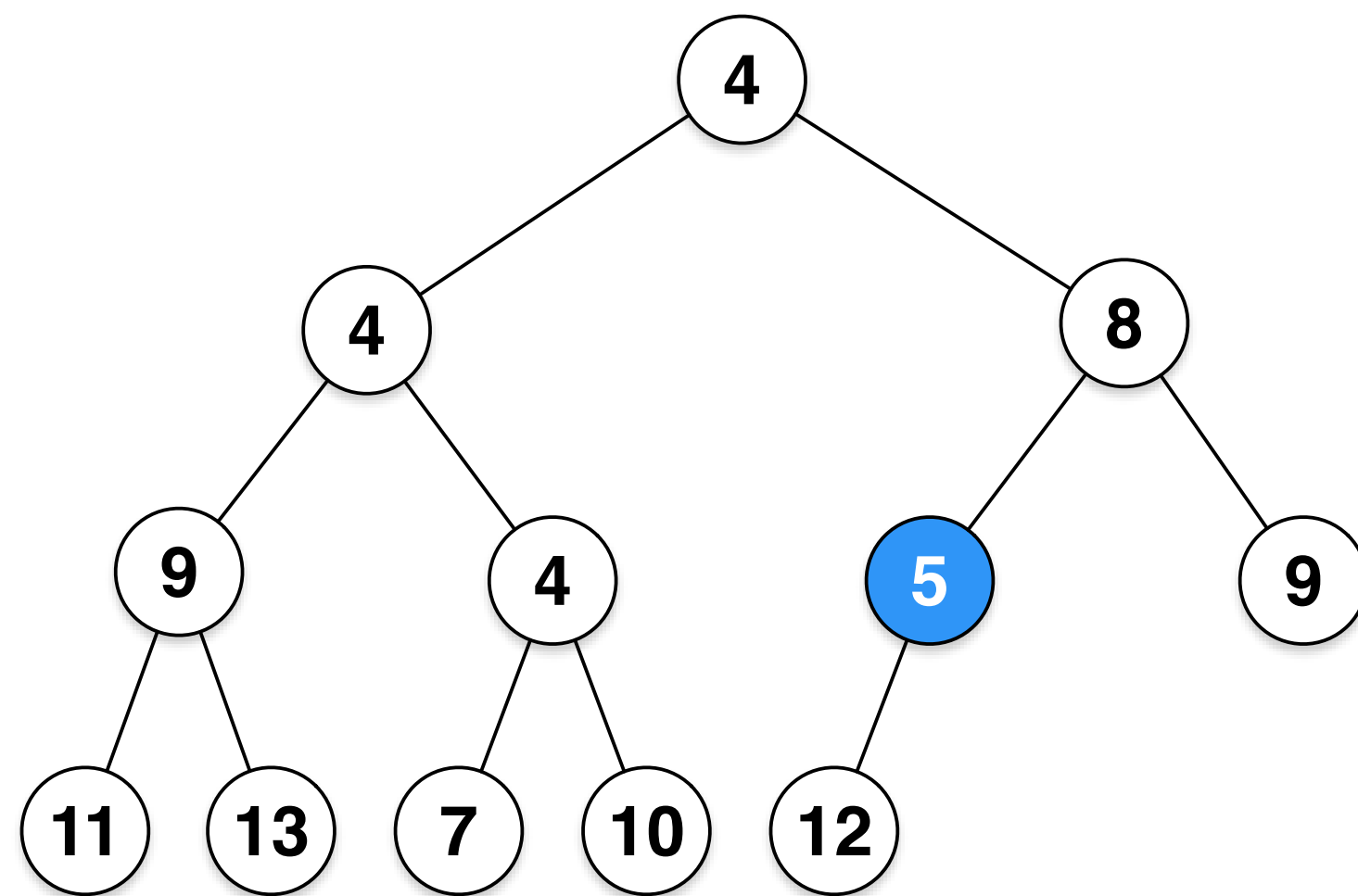
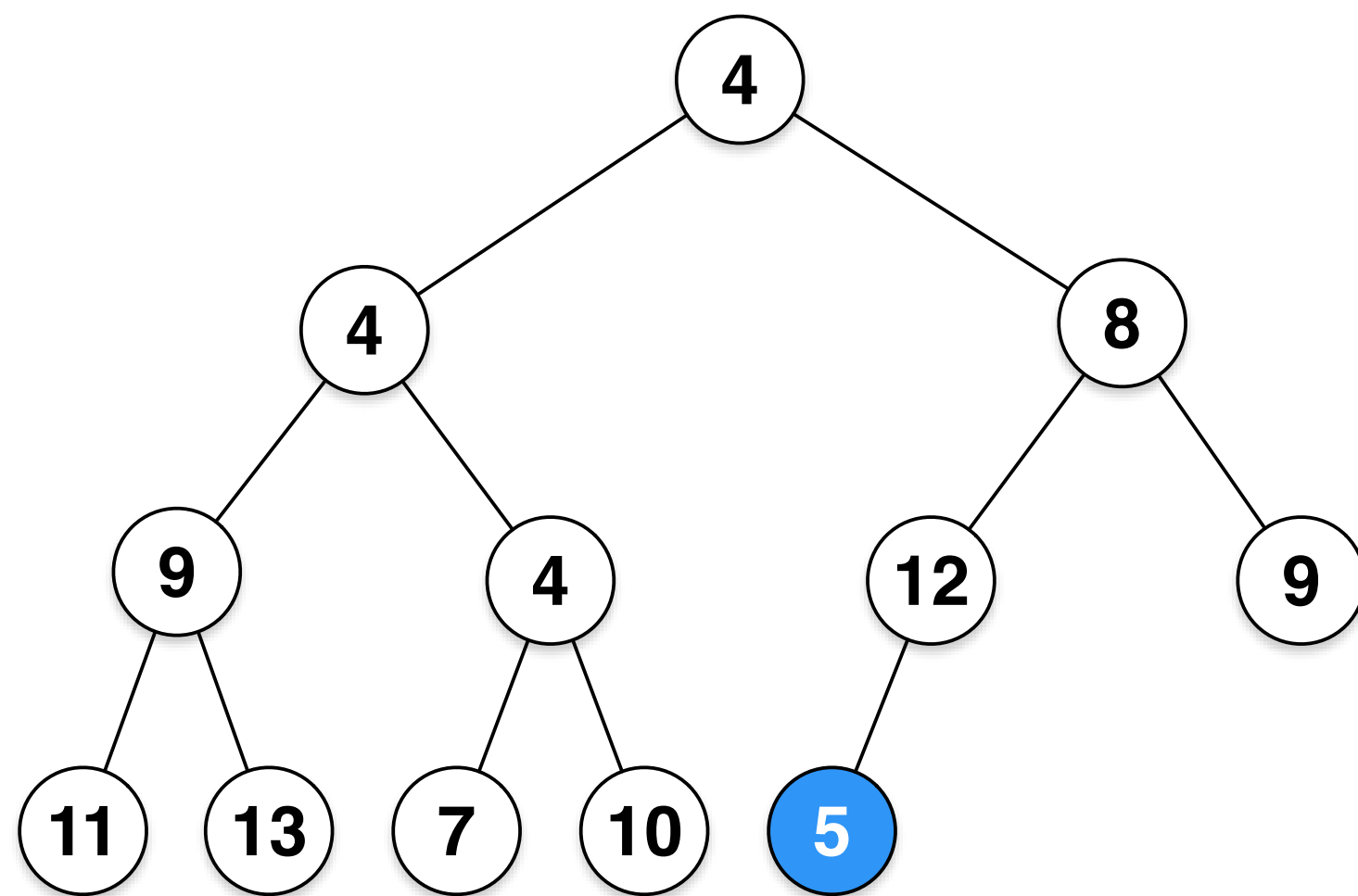


Heapify Up

1. **HeapifyUp**(A, k):
2. while(parent(k).key > k.key):
3. swap(parent(k), k)



Heapify Up



Insert Running Time

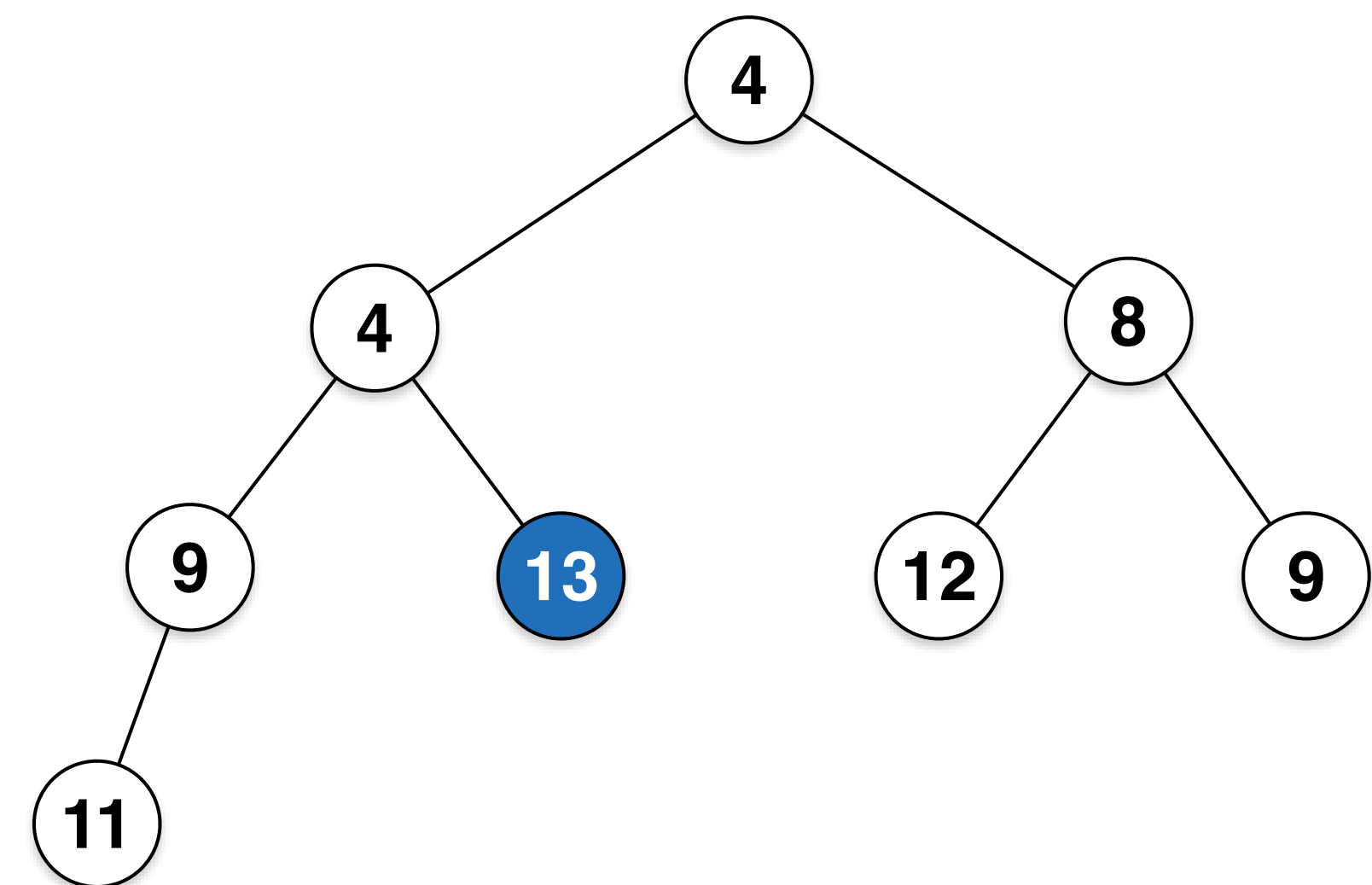
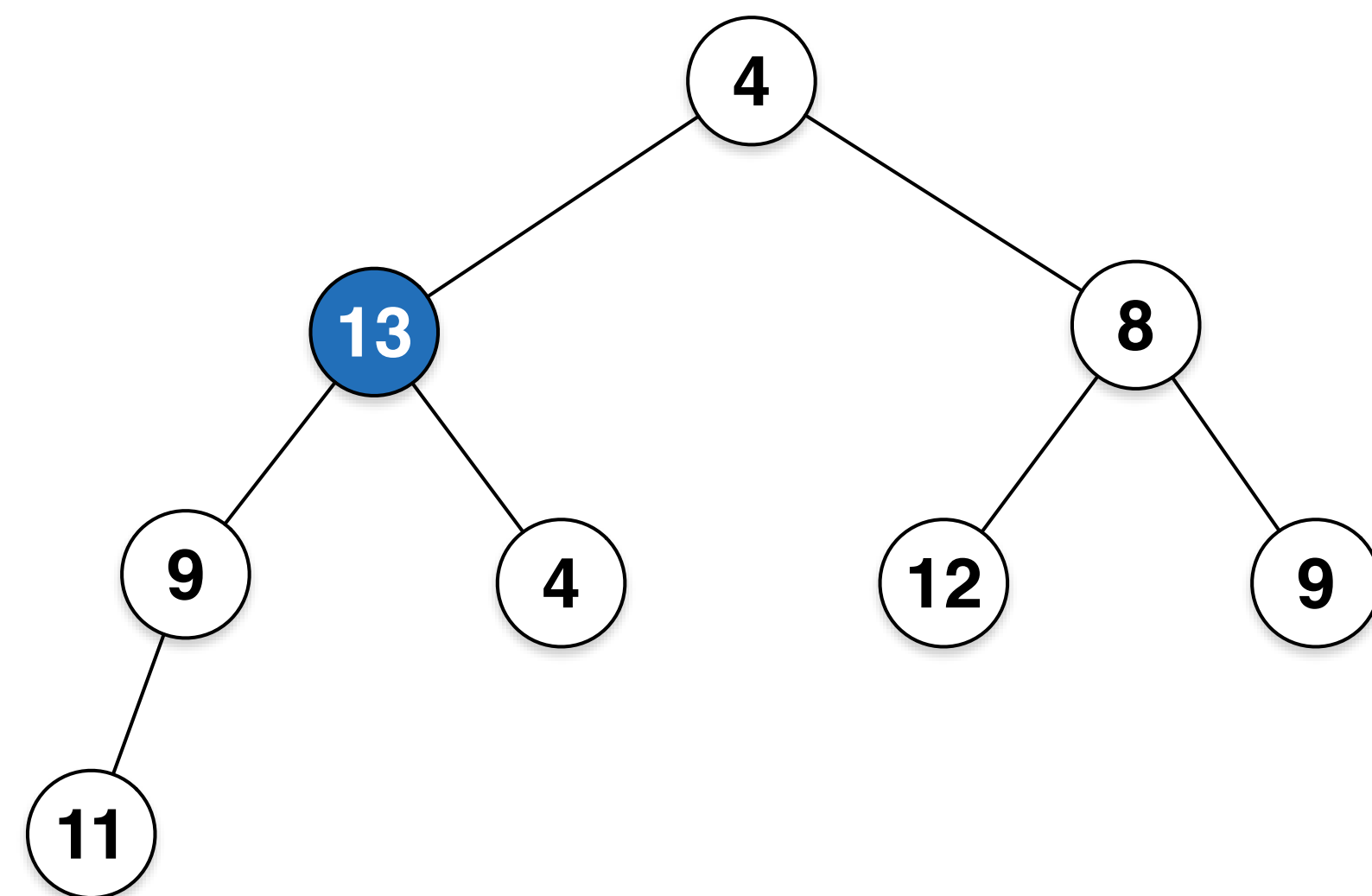
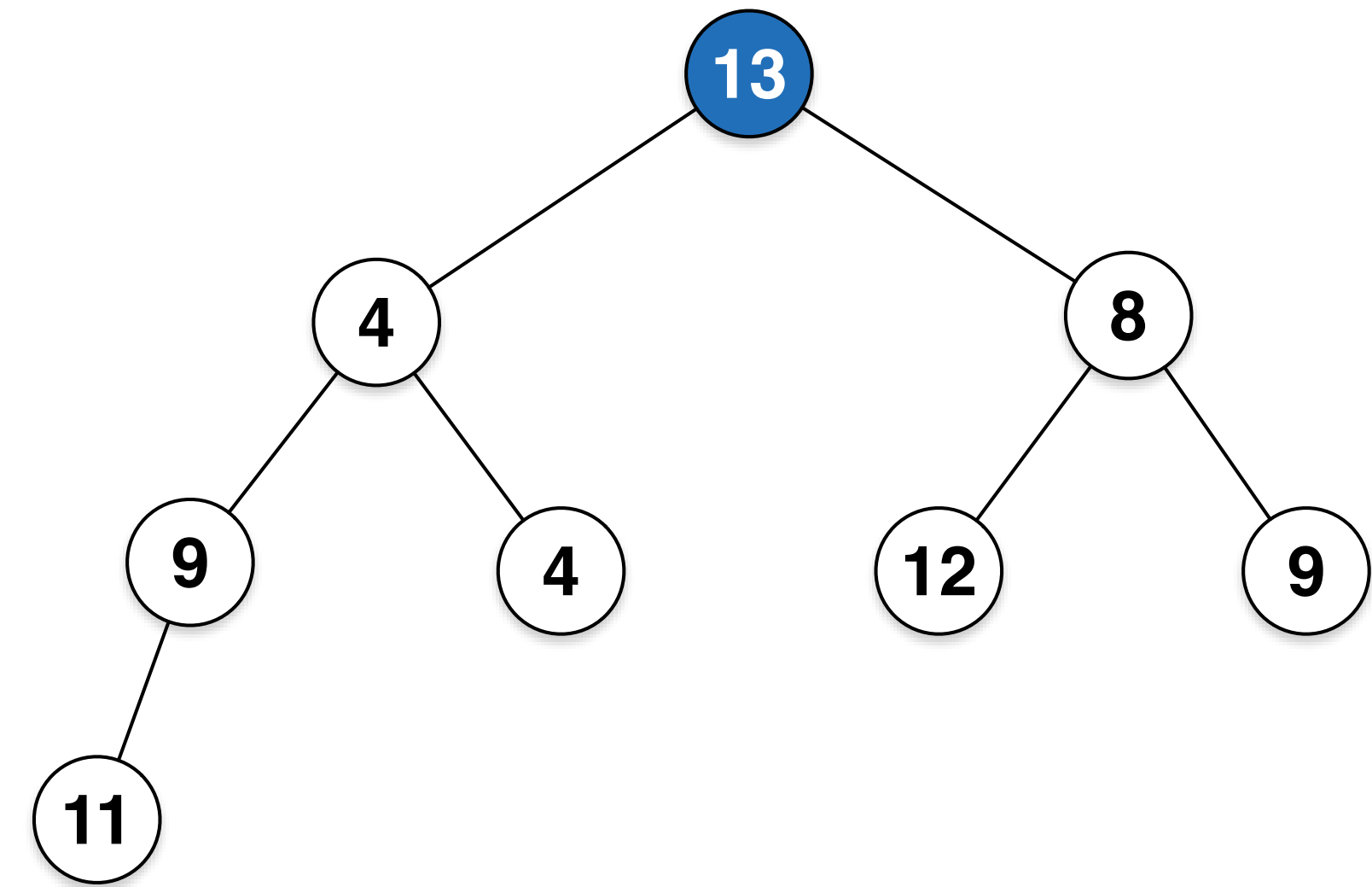
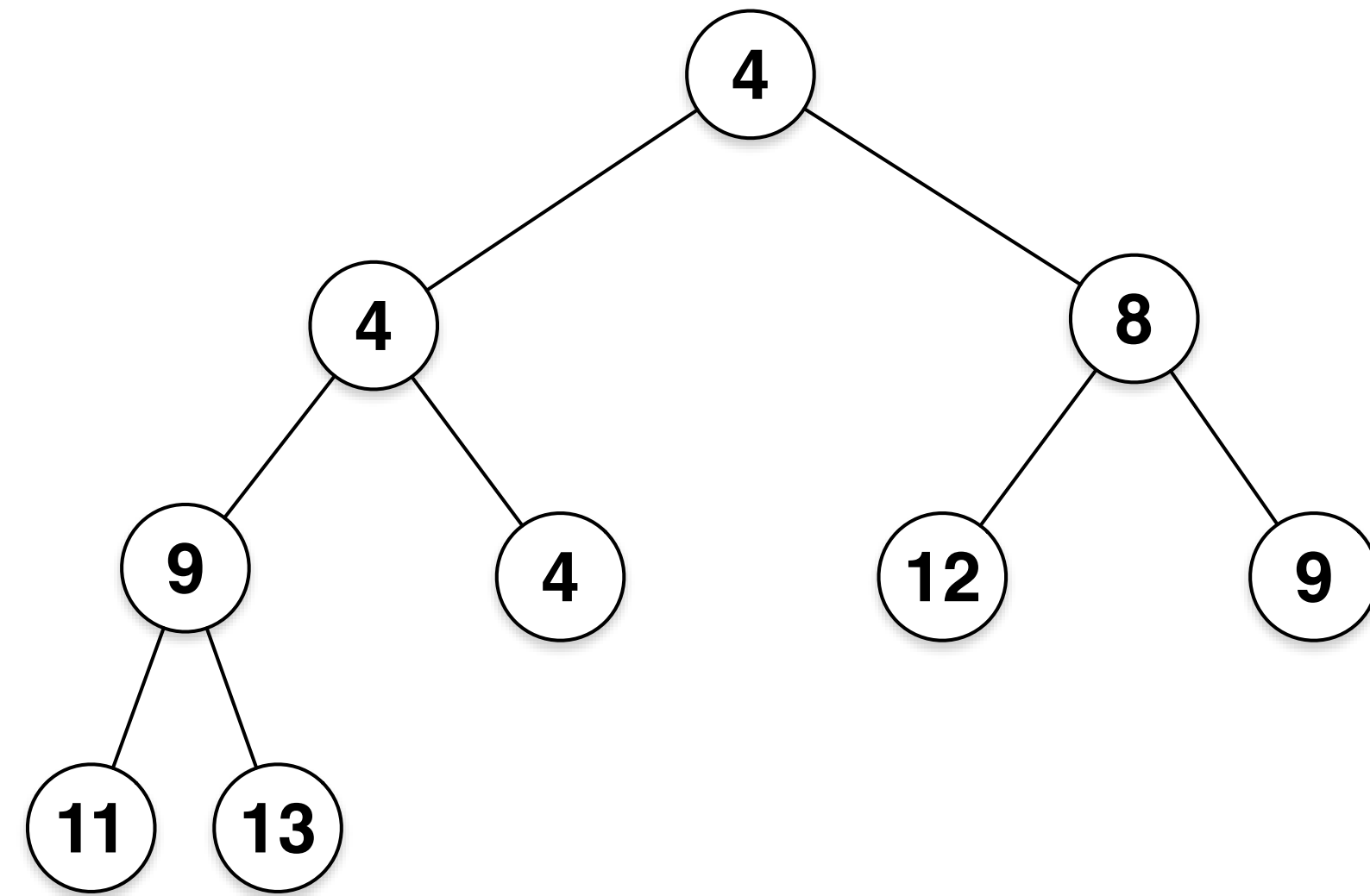
Extract Min

1. **ExtractMin(A):**
2. swap root with the last element
3. remove root
4. **MinHeapify(A, 1)**
5. return root

MinHeapify

```
1. MinHeapify(A, i):
2.     left = 2i
3.     right = 2i + 1
4.     min = i
5.     if (A[min] < A[left] && A[min] < A[right])
6.         return
7.     if (A[left] < A[right]):
8.         min = left
9.     else:
10.        min = right
11.    swap(A[i], A[min])
12.    MinHeapify(A, min)
```

MinHeapify



Heapify

1. **Heapify(A):**
2. Assume A is already a heap
3. start = floor($n/2$) // first root with children
4. for i = start to 1:
5. **MinHeapify(A, i)**

Heapify Running Time

- $O(\log n)$??
- If all the subtree at height h has been heapified, then heapifying the sub tree at $h+1$ level will only require bubbling down the root nodes.
 - $O(h)$ operations (swap) per node
 - Height is measured from bottom up starting at 0
- Notice how most of the heapifying happens at the bottom.

Heapify Running Time

- $O(h)$ operations (swap) per node
- $\text{NodeCount}(h) = \lceil 2^{(\log_2 n - h) - 1} \rceil = \left\lceil \frac{2^{\log n}}{2^{h+1}} \right\rceil = \left\lceil \frac{n}{2^{h+1}} \right\rceil$
- Cost of heapifying the entire tree:

$$\sum_{h=1}^{\lceil \log n \rceil} \frac{n}{2^{h+1}} O(h) = O \left(n \sum_{h=1}^{\lceil \log n \rceil} \frac{h}{2^{h+1}} \right)$$
$$\leq O \left(n \sum_{h=0}^{\infty} \frac{h}{2^h} \right) = O(2n) = O(n)$$

- http://www.symbolab.com/solver/series-calculator/%5Csum_%7Bn%3D0%7D%5E%7B%5Cinfty%7D%20%5Cfrac%7Bn%7D%7B2%5E%7Bn%7D%7D