

Week 3 - Lecture Notes

Topics: - Analysis of QuickSort
Randomized QuickSort
Heap
Heap Sort
Decision Tree

Pseudo-code for Quick Sort

QUICK SORT (A, p, r)

1. if $p < r$
2. then $q \leftarrow \text{PARTITION}(A, p, r)$
3. QUICK SORT (A, p, q)
4. QUICK SORT ($A, q+1, r$)

Initial call: QUICK SORT ($A, 1, n$)

Analysis of Quick sort

- Assume all input elements are distinct
- In practice, there are better partitioning algorithms for when duplicate input may exist.
- Let $T(n)$ = worst case running time on an array of n elements

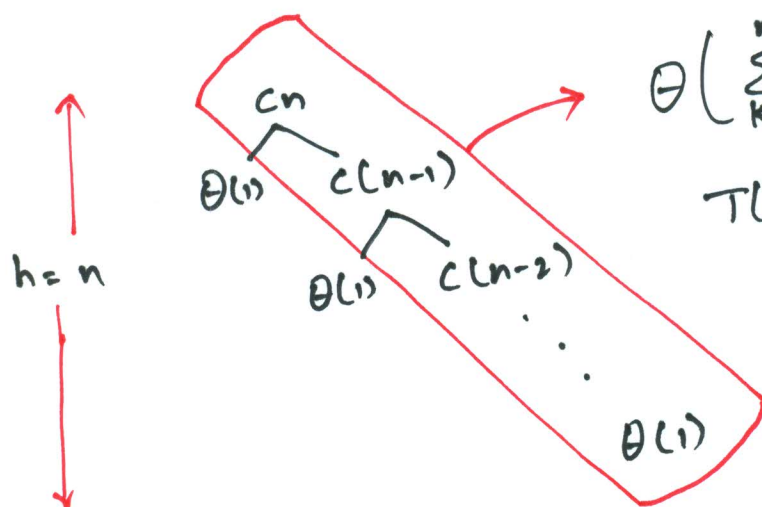
Worst case of Quick Sort

- Input is sorted or reverse sorted
- Partition around minimum or maximum element
- Split $\rightarrow 0 : n-1$,
one side of the partition always has no element.

$$\begin{aligned}T(n) &= T(0) + T(n-1) + \theta(n) \\&= \theta(1) + T(n-1) + \theta(n) \\&= \theta(n) + T(n-1) \\&= \theta(n^2)\end{aligned}$$

Worst Case Recursion Tree

$$T(n) = T(0) + T(n-1) + cn$$



$$\theta\left(\sum_{k=1}^n k\right) = \theta(n^2)$$

$$\begin{aligned}T(n) &= \theta(n) + \theta(n^2) \\&= \theta(n^2)\end{aligned}$$

Best - Case Analysis

If we are lucky, PARTITION splits the array evenly ($\frac{1}{2} : \frac{1}{2}$)

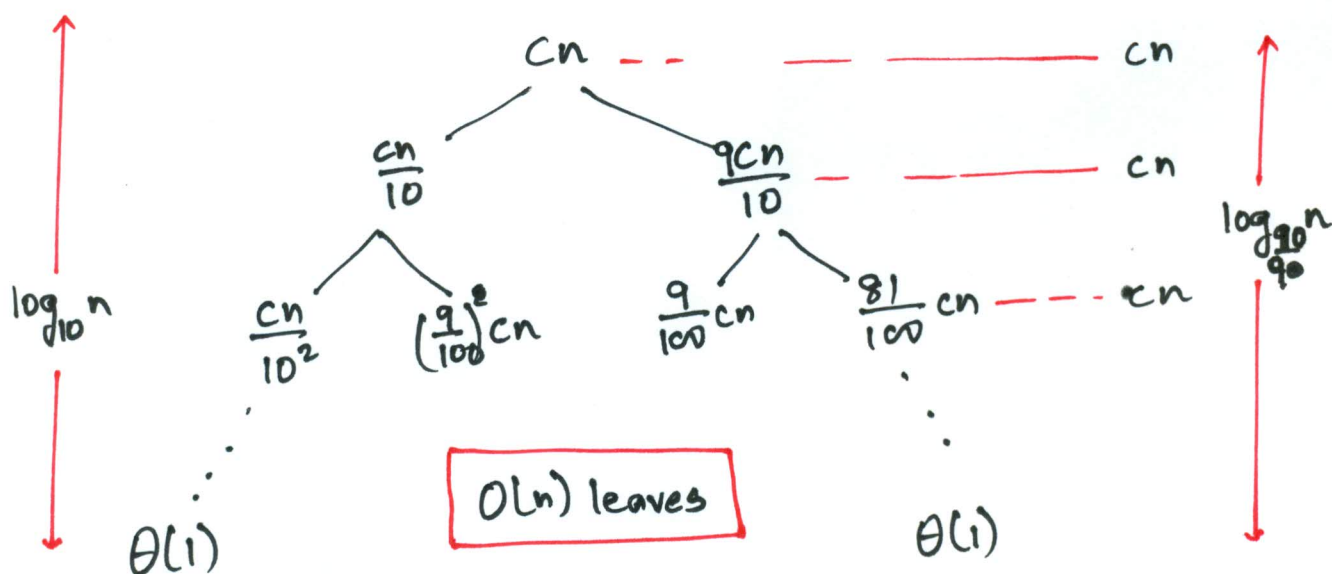
$$T(n) = 2T(n/2) + \theta(n)$$

$$= \theta(n \log n) \quad \left[\text{same as Merge Sort} \right]$$

Analysis of "almost - best" case

Consider the split is always $\frac{1}{10} : \frac{9}{10}$.

$$T(n) = T\left(\frac{1}{10}n\right) + T\left(\frac{9}{10}n\right) + \theta(n)$$



$$cn \log_{10} n \leq T(n) \leq cn \log_{\frac{10}{9}} n + O(n)$$

$$\theta(n \log n) \quad \text{Lucky}$$

More intuition

Let us consider a case in which we are alternate lucky, unlucky, lucky, unlucky, lucky, ...

$$L(n) = 2U(n/2) + \Theta(n) \quad \text{Lucky}$$

$$U(n) = L(n-1) + \Theta(n) \quad \text{unlucky}$$

Solving we get.

$$L(n) = 2\left(L(n/2 - 1) + \Theta(n/2)\right) + \Theta(n)$$

$$= 2L(n/2 - 1) + \Theta(n)$$

$$= \Theta(n \log n) \quad \underline{\text{Lucky}}$$

So, even in this case we are lucky.

How can we make sure we are usually lucky?

Randomized QuickSort

Idea: Partition around a random element

- Running order is independent of the input order.
- No assumptions need to be made about the input distribution.
- No specific input elicits the worst case behaviour.
- The worst case is determined only by the output of a random-number generator.

Randomized Quick Sort Analysis

Let $T(n)$ = the random variable for the running time of randomized quicksort on an input of size n , assuming random numbers are independent.

For $k = 0, 1, \dots, n-1$, define the indicator random variable as:

$$X_k = \begin{cases} 1 & \text{if PARTITION generates } k:n-k-1 \text{ split} \\ 0 & \text{otherwise} \end{cases}$$

$E[X_k] = \Pr\{X_k = 1\} = 1/n$, since all splits are equally likely, assuming elements are distinct.

$$T(n) = \begin{cases} T(0) + T(n-1) + \theta(n) & \text{if } 0:n-1 \text{ split} \\ T(1) + T(n-2) + \theta(n) & \text{if } 1:n-2 \text{ split} \\ \vdots \\ T(n-1) + T(0) + \theta(n) & \text{if } n-1:0 \text{ split.} \end{cases}$$

$$= \sum_{k=0}^{n-1} X_k (T(k) + T(n-k-1) + \theta(n))$$

Calculating expectation

$$E[T(n)] = E\left[\sum_{k=0}^{n-1} X_k (T(k) + T(n-k-1) + \theta(n))\right]$$

$$= \sum_{k=0}^{n-1} E[X_k (T(k) + T(n-k-1) + \theta(n))]$$

$$= \sum_{k=0}^{n-1} E[X_k] E[T(k) + T(n-k-1) + \theta(n)]$$

$$= \frac{1}{n} \sum_{k=0}^{n-1} E[T(k)] + \frac{1}{n} \sum_{k=0}^{n-1} E[T(n-k-1)]$$

$$+ \frac{1}{n} \sum_{k=0}^{n-1} \theta(n)$$

$$= \frac{2}{n} \sum_{k=1}^{n-1} E[T(k)] + \theta(n)$$

$$= \frac{2}{n} \sum_{k=2}^{n-1} E[T(k)] + \theta(n)$$

The $k=0,1$ terms can be absorbed in the $\theta(n)$

Prove: $E[T(n)] \leq an \lg n$ for constant $a > 0$

Choose 'a' large enough so that $an \lg n$ dominates $E[T(n)]$ for sufficiently small $n \geq 2$.

Use fact: $\sum_{k=2}^{n-1} k \lg k \leq \frac{1}{2} n^2 \lg n - \frac{1}{8} n^2$

$$\begin{aligned} E[T(n)] &\leq \frac{2}{n} \sum_{k=2}^{n-1} a k \lg k + \theta(n) \\ &\leq \frac{2a}{n} \left(\frac{1}{2} n^2 \lg n - \frac{1}{8} n^2 \right) + \theta(n) \\ &= an \lg n - \left(\frac{an}{4} - \theta(n) \right) \\ &\leq an \lg n \end{aligned}$$

if a is chosen large enough so
than $\frac{an}{4}$ dominates $\theta(n)$

Quicksort in Practice

- Quicksort is a great general-purpose sorting algorithm
- Quicksort can benefit substantially from code tuning
- Quicksort is typically over twice as fast as merge sort.

Priority Queue

A data structure implementing a set S of elements, each associated with a key, supporting the following operations.

$\text{insert}(S, x)$: insert element x into set S

$\text{max}(S)$: return element of S with largest key

$\text{extract-max}(S)$: return element of S with largest key and remove it from S

$\text{increase-key}(S, x, K)$: increase the value of element x 's key to new value K .

Heap

- Implementation of a priority queue
- An array, visualized as a nearly complete binary tree
- Max Heap Property: The key ~~node~~ of a node is \geq than the keys of its children
(Min Heap defined analogously)

Heap is a Tree

root of tree : first element in the array,
corresponding $i=1$

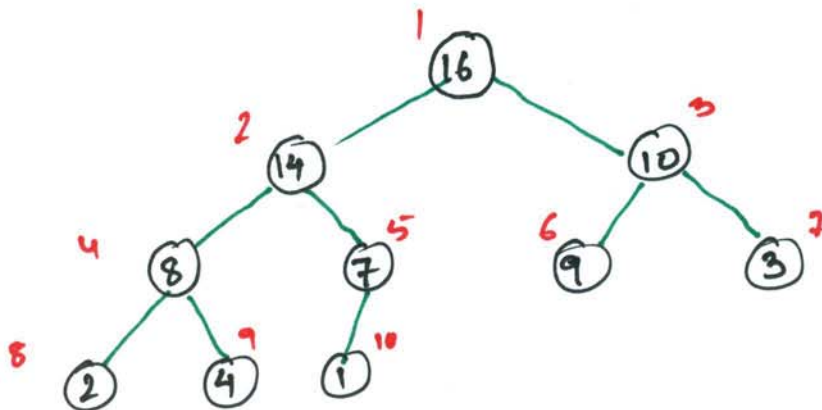
parent $(i) = i/2$: returns index of node's parent

left $(i) = 2i$: returns index of node's left child.

right $(i) = 2i+1$: returns index of node's right child.

Example:

16 14 10 8 7 9 3 2 4 1



No pointers required.

Height of a binary heap is $O(\lg n)$

Heap sort as a tree

Heap Operations:

max-heapify: correct a single violation of the heap property in a subtree at its root

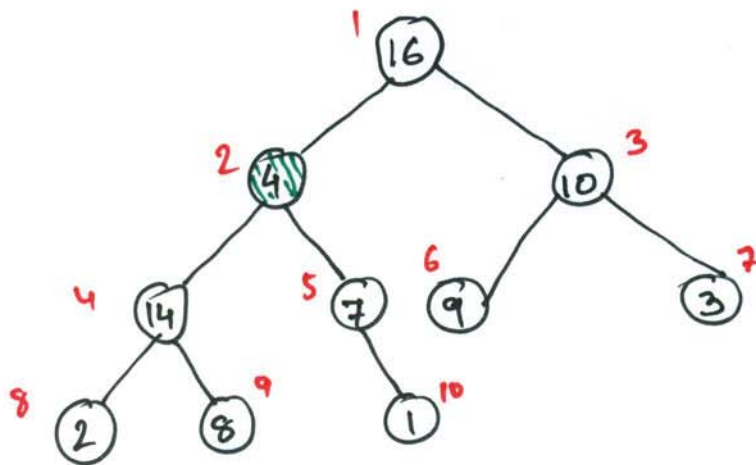
build-max-heap: produce a max-heap from an unordered array.

insert, extract-max, heapsort.

Max-heapify

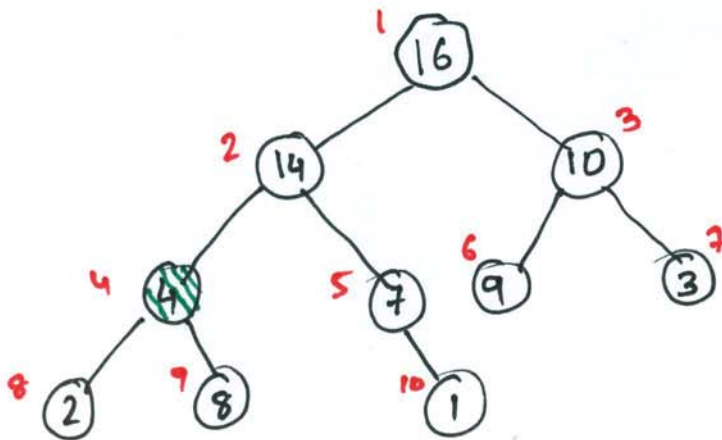
- Assume that the trees rooted at $\text{left}(i)$ and $\text{right}(i)$ are max-heaps
- If $A[i]$ violates the max-heap property, correct violation by "tricking" element $A[i]$ down the tree, making the subtree rooted at index i a max-heap.

Max-heapify (Example)



MAX-HEAPIFY(A, 2)
heap-size[A] = 10

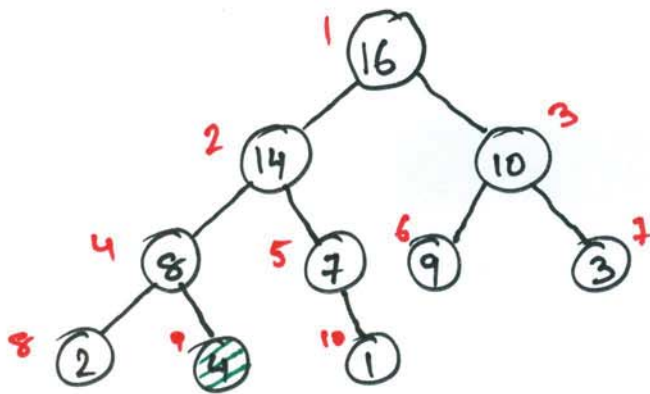
Node 10 is the left child of node 5 but is drawn to the right for convenience.



Exchange A[2] with A[4]

Call MAX-HEAPIFY(A, 4)

because max-heap property is
violated



Exchange $A[4]$ with $A[9]$

No more calls.

Time = $O(\log n)$

Max. Heapify Pseudocode

1. $l = \text{left}(i)$
2. $r = \text{right}(i)$
3. if $(l \leq \text{heap-size}(A) \text{ and } A[l] > A[i])$
4. then $\text{largest} = l$
5. else $\text{largest} = i$
6. if $(r \leq \text{heap-size}(A) \text{ and } A[r] > A[\text{largest}])$
7. then $\text{largest} = r$
8. if $\text{largest} \neq i$
9. then exchange $A[i]$ and $A[\text{largest}]$
10. Max-Heapify $(A, \text{largest})$

Build-Max-Heap (A)

Converts $A[1, \dots, n]$ to a max heap

Build-Max-Heap (A):

for $i = n/2$ down to 1

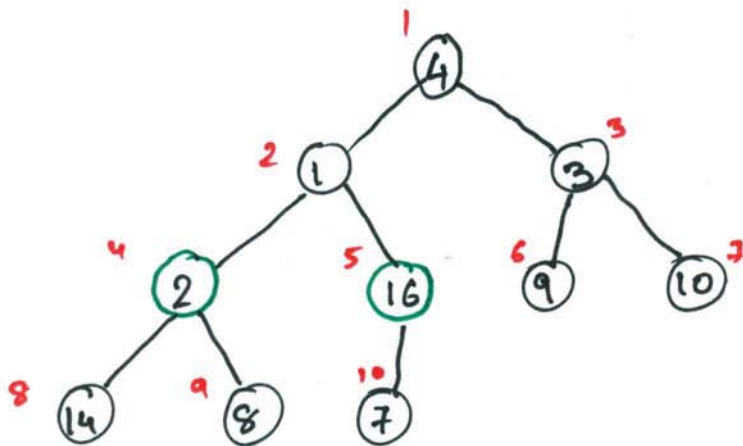
do Max-Heapify(A, i)

- We start at $i = n/2$ because elements $A[n/2+1, \dots, n]$ are all leaves of the tree

$2i > n$, for $i > n/2 + 1$

Build-Max-Heap Example

4 1 3 2 16 9 10 14 8 7

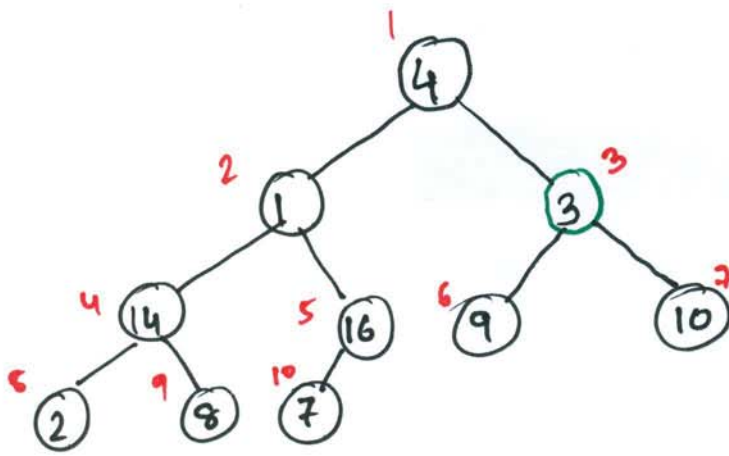


Max-Heapify(A, 5)

no change

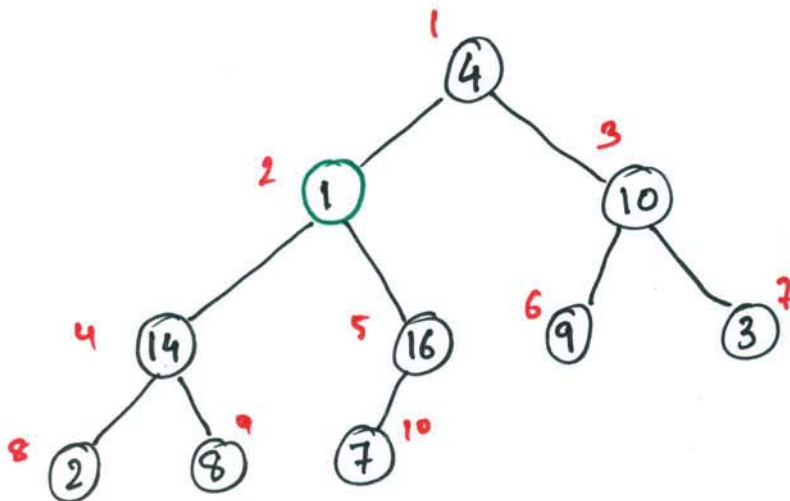
Max-Heapify(A, 4)

Swap $A[4]$ and $A[8]$



Max-Heapify (A, 3)

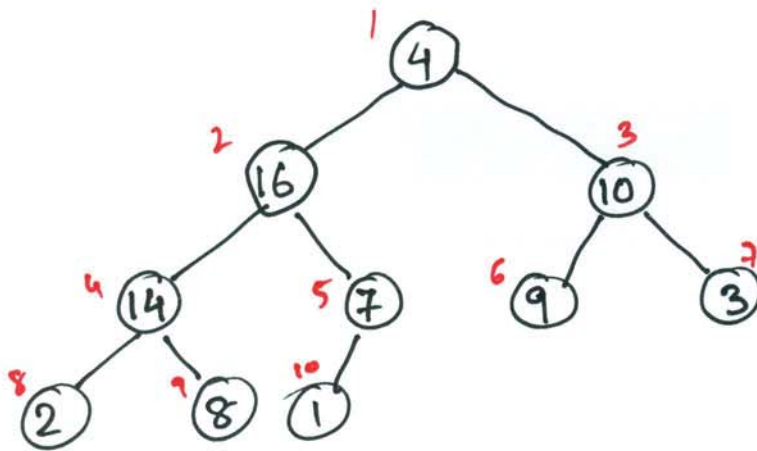
Swap A[3] and A[7]



Max-Heapify (A, 2)

Swap A[2] and A[5]

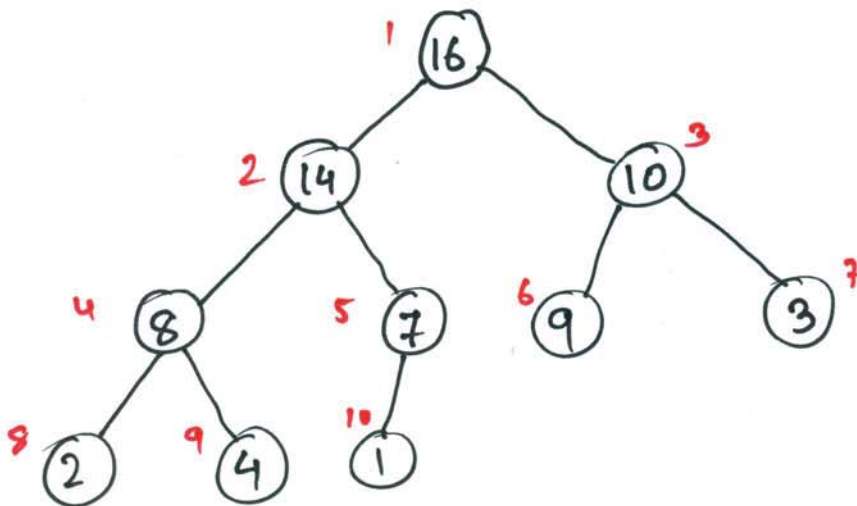
Swap A[5] and A[10]



Max-Heapify (A,1)
 Swap [A,1] with A[2]
 Swap [A,2] with A[4]
 Swap [A,4] with A[9]

So,

A: 4 1 3 2 16 9 10 14 8 7



Build-Max-Heap (A) Analysis

We can observe that Max-Heapify takes $O(1)$ for nodes that are one level above the leaves, and in general, $O(l)$ for the nodes that are l levels above the leaves.

We have $n/4$ nodes with level 1, $n/8$ with level 2, and so on till we have one root node that is $\lg n$ levels above the leaves.

So, total amount of work in the for loop can be summed as:

$$n/4(1c) + n/8(2c) + n/16(3c) + \dots + 1(\lg c)$$

Setting $n/4 = 2^k$ and simplifying we get

$$c 2^k \left(\frac{1}{2^0} + \frac{2}{2^1} + \frac{3}{2^2} + \dots + \frac{(k+1)}{2^k} \right)$$

The term in brackets is bounded by a constant.

This means that Build-Max-Heap is $O(n)$

Heap - Sort

Sorting Strategy

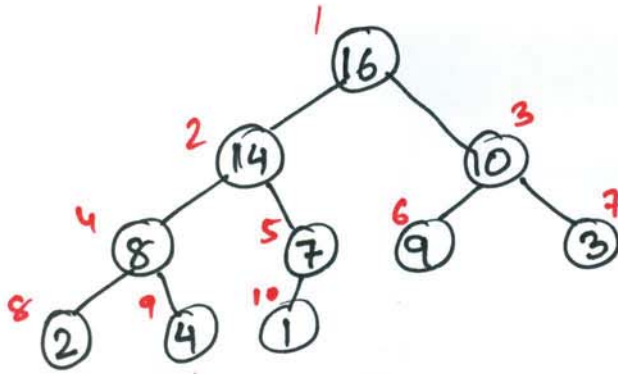
1. Build Max Heap from unordered array;
2. Find maximum element $A[i]$
3. Swap elements $A[n]$ and $A[i]$
now max element is at the end of array.
4. Discard node n from heap
(by decrementing heap-size variable)
5. New root may violate max heap property, but its children are max heaps. Run max-heap to fix this.
6. Go to Step 2 unless heap is empty.

Heap Sort Running Time

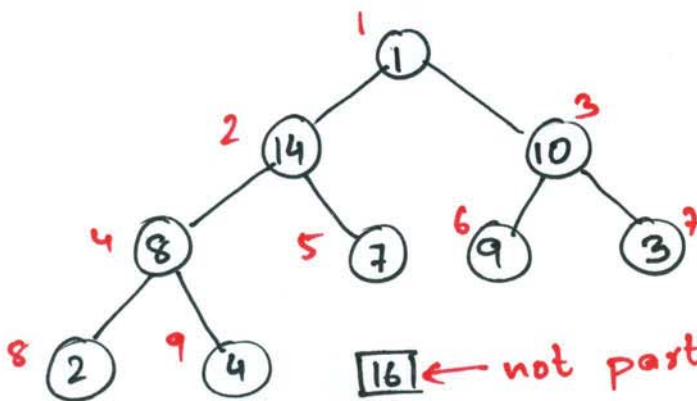
- after n iterations Heap is empty
- every iteration involves a swap and a max. heapify operation; $O(\log n)$ time.

Hence, overall : $O(n \log n)$.

Heap-sort Example

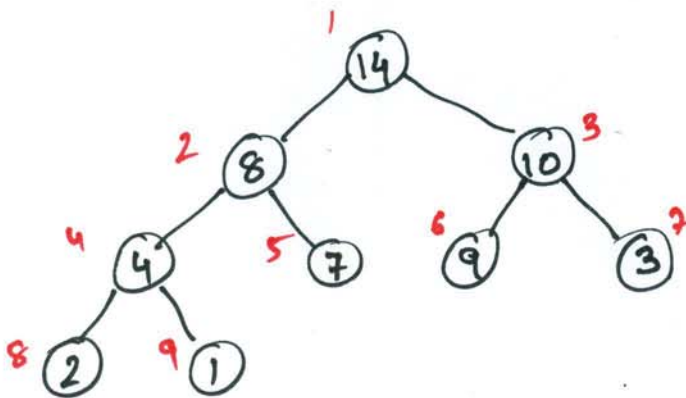


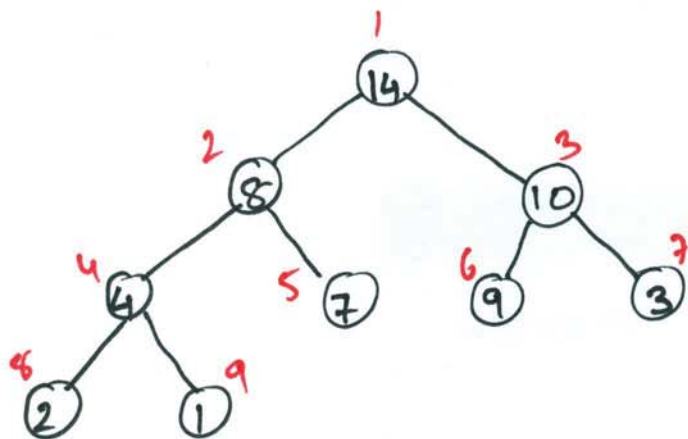
Swap $A[10]$ and $A[1]$



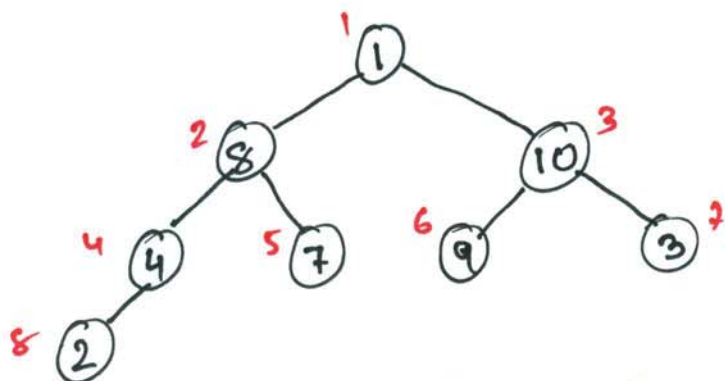
Heap size = 9

Max-heapify ($A, 1$)





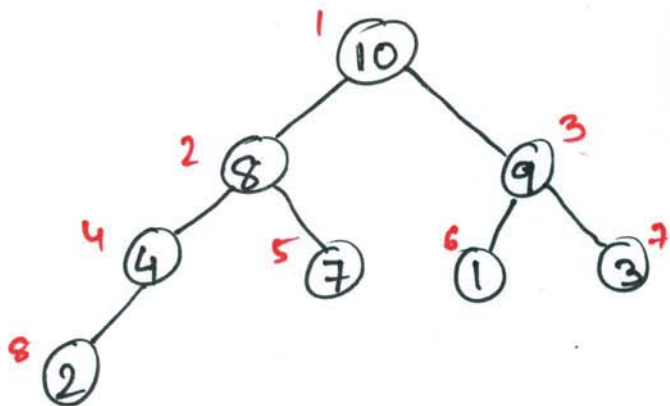
Swap $A[9]$ and $A[1]$

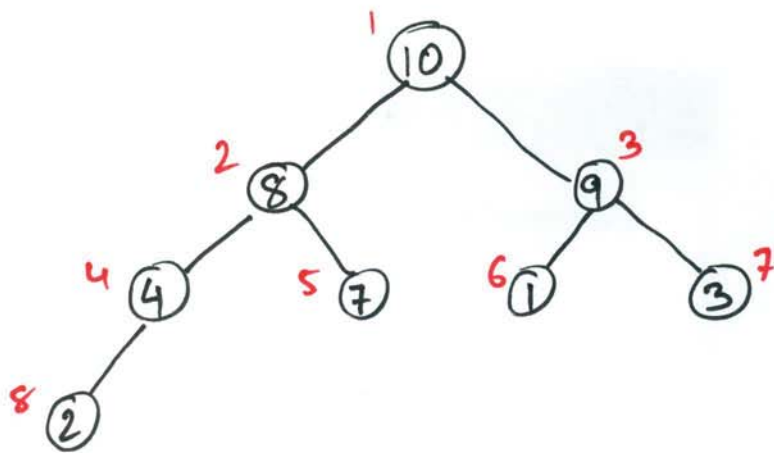


heap size = 8

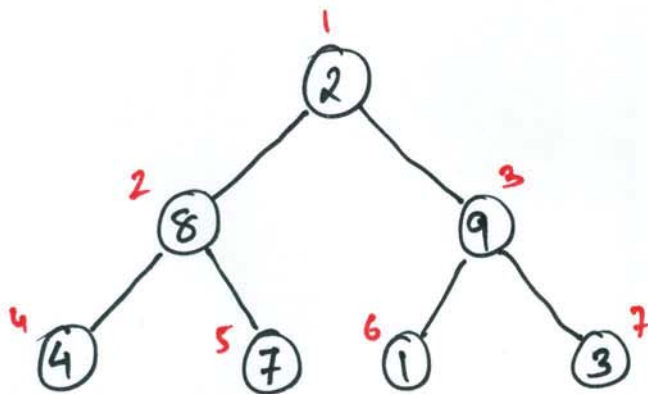
14 16 ← not part of heap

Max-heapify ($A, 1$)





Swap $A[8]$ and $A[1]$



$\begin{matrix} 8 \\ \boxed{10} \end{matrix}$
 $\begin{matrix} 9 \\ \boxed{14} \end{matrix}$
 $\begin{matrix} 10 \\ \boxed{16} \end{matrix}$
 ← not part of heap.

How fast can we Sort?

All the sorting algorithms we have seen so far are comparison sorts: only use comparisons to determine the relative order of elements.

- E.g.: insertion sort, mergesort, quicksort, heap sort.

The best worst case running time that we have seen for comparison sorting is $O(n \log n)$

Is $O(n \log n)$ the best we can do?

→ Decision Trees can help answer this question

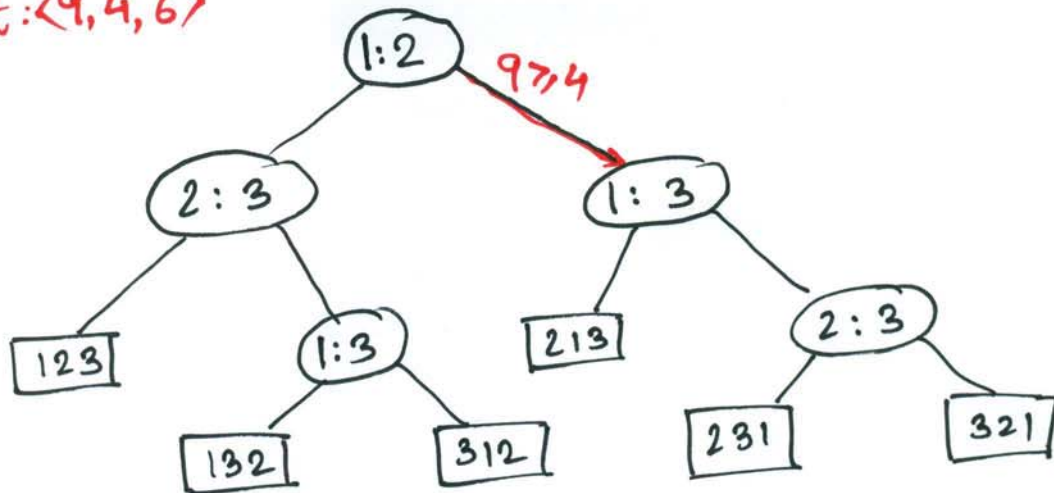
Decision Tree Model

A decision tree can model the execution of any comparison sort:

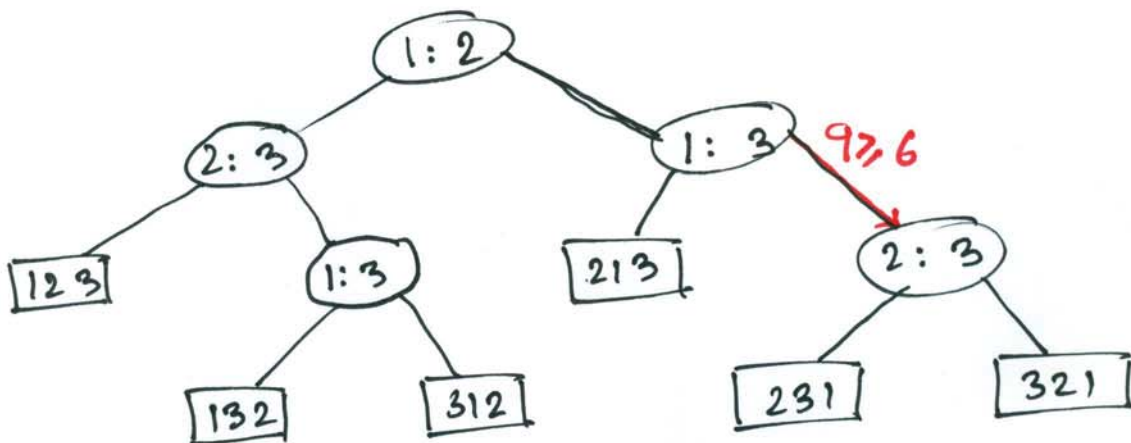
- One tree for each input size n .
- View the algorithm as splitting whenever it compares two elements.
- The tree contains the comparisons along all possible instruction traces.
- The running time of algorithm = the length of the path taken.
- Worst case running time = height of the tree.

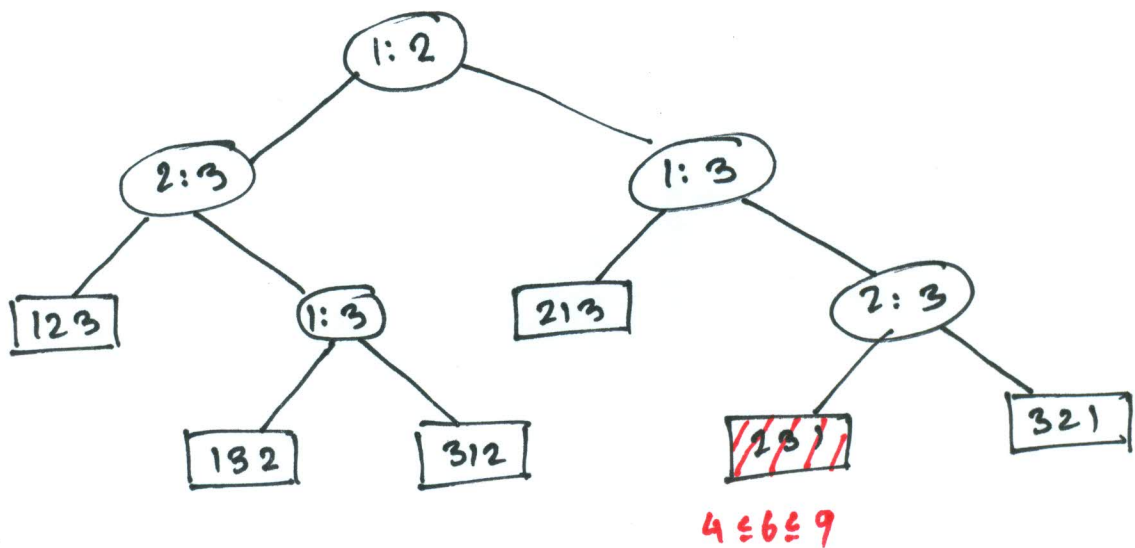
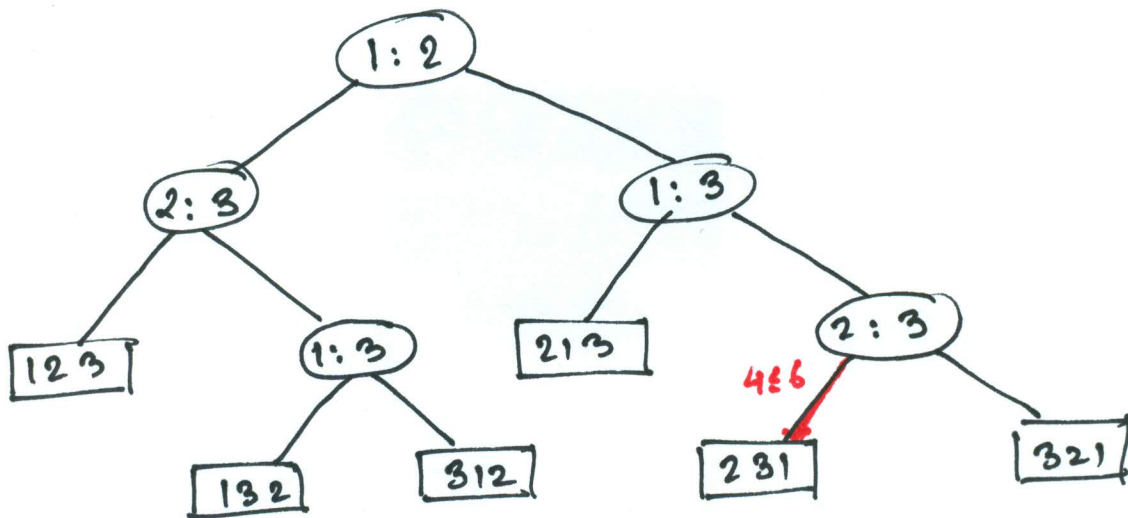
Decision Tree Example

Sort: $\langle 9, 4, 6 \rangle$



- Each internal node is labelled $i:j$ for $i, j \in \{1, 2, \dots, n\}$
- The left subtree shows subsequent comparisons if $a_i \leq a_j$
- The right subtree shows subsequent comparisons if $a_i > a_j$





Each leaf contains a permutation
 $\langle \pi(1), \pi(2), \dots, \pi(n) \rangle$ to indicate the ordering
 $a_{\pi(1)}, a_{\pi(2)}, \dots, a_{\pi(n)}$ has been
 established.

Lower bound for decision tree Sorting

Theorem:

Any decision tree that can sort n elements must have height $\Omega(n \lg n)$

Proof:

The tree must contain $\geq n!$ leaves, since there are $n!$ possible permutations.

A height- h binary tree has $\leq 2^h$ leaves.

Thus $n! \leq 2^h$

$$\therefore h \geq \log(n!)$$

$$\geq \log\left(\left(\frac{n}{e}\right)^n\right) \quad [\text{Stirling's formula}]$$

$$= n \lg n - n \lg e$$

$$= \Omega(n \lg n)$$

Corollary:

Heapsort and Merge sort are asymptotically optimal comparisons sorting algorithm.