

Chapter 2. Sorting & Order Statistics

We discuss :

- Sorting algor. (HeapSort , QuickSort)
- Median & Order Statistics
- A Lower Bound on Sorting .

Selection Sort.

We've discussed Insertion Sort : it's in-place and has time complexity $O(n^2)$.

Idea of Selection Sort:

- Remove largest element and add it to output sequence until input sequence is exhausted.

SELECTION (A[1..n])

$$j = n$$

while $j > 1$

$$k = j ; \max = j ; S = A[j]$$

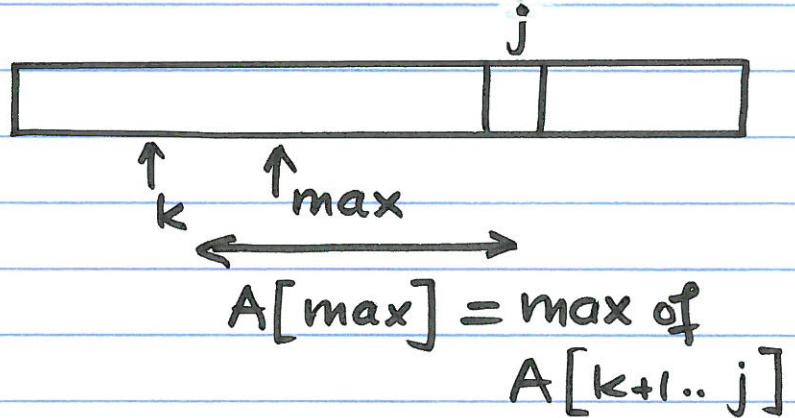
while $k > 1$

$$k = k - 1$$

if $A[k] > S$ then
 $\max = k$
 $S = A[k]$

swap $A[\max]$ with $A[j]$
 $j = j - 1$

Illustration



Complexity of SELECTION :

- Searching for max in $A[1..j]$ takes $O(j)$ steps
- Outer loop is executed for values $j = n, n-1, \dots, 2.$

$$\begin{aligned} \text{Thus, } T(n) &= O\left(\sum_{j=2}^n O(j)\right) \\ &= O(n^2). \end{aligned}$$

Heapsort

Observation: If we can find max faster, we can improve on SELECTION.

→ We organize $A[1..n]$ as a heap

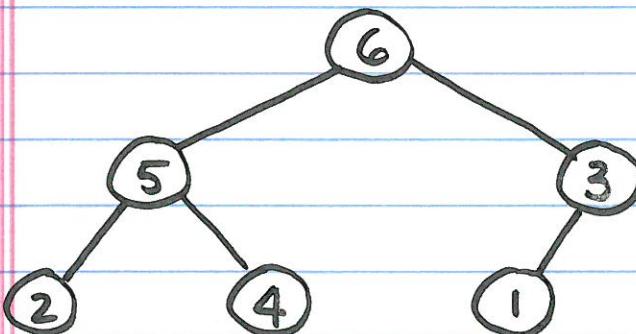
A heap is an array object that can be identified as an almost complete binary tree : $\text{parent}(i) = \lfloor i/2 \rfloor$
 $\text{left}(i) = 2i$ and $\text{right}(i) = 2i+1$

A heap is a max heap if

$$A[\text{parent}(i)] \geq A[i]$$

(min heap if $A[\text{parent}(i)] \leq A[i]$)

Example: (max heap)



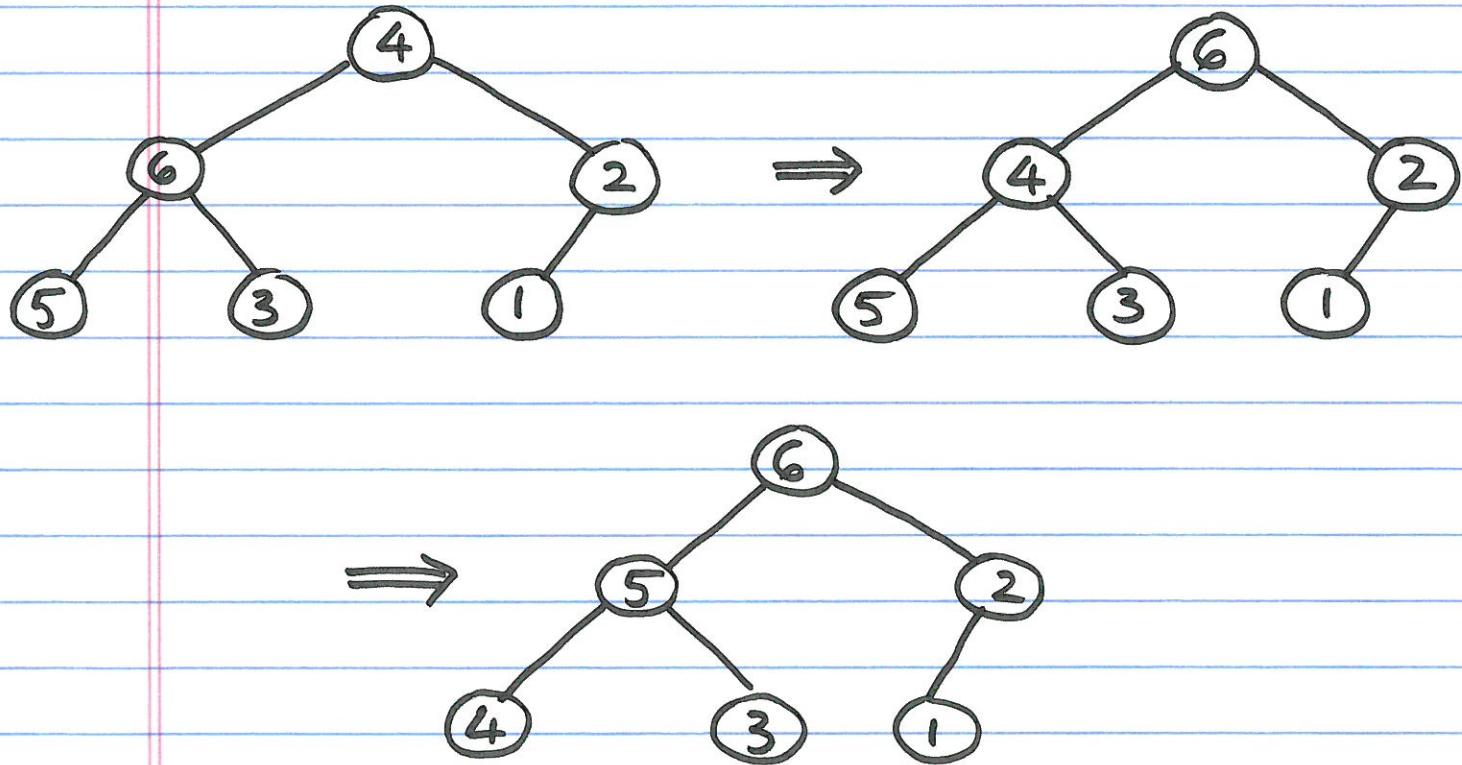
6	5	3	2	4	1
---	---	---	---	---	---

Max-Heapify

Max-Heapify is a procedure on (A, i) :

- Assumes subtrees rooted at $\text{left}(i)$, $\text{right}(i)$ are max-heaps but $A[i]$ may not sat. max.heap property
- Transform subtree rooted at $A[i]$ into a max.heap by "floating $A[i]$ down".

Example.



Max-Heapify (A, i)

$l = \text{left}(i)$; $r = \text{right}(i)$

if $l \leq \text{heap-size}[A]$ and $A[l] > A[i]$

then largest = l

else largest = i ;

if $r \leq \text{heap-size}[A]$ and

$A[r] > A[\text{largest}]$

then largest = r ;

if largest $\neq i$

then

exchange $A[i] \leftrightarrow A[\text{largest}]$

Max-Heapify (A, largest)

Complexity of Max-Heapify

- Subtree rooted at $A[\text{largest}]$ for recursive call is of size at most $2n/3$, where $n = \text{size of subtree rooted at } A[i]$.

(Note left subtree is larger when bottom level is half-full.)

$$\text{Thus, } T(n) \leq T(\lfloor 2n/3 \rfloor) + \Theta(1)$$

$$= O(\lg n)$$

Remark. If $h = \text{height of Subtree}$ rooted at $A[i]$, then $T(n) = O(h) = O(\lg n)$.

Building a Heap

Build-Max. Heap (A)

$\text{heap-size}(A) = \text{length}(A)$

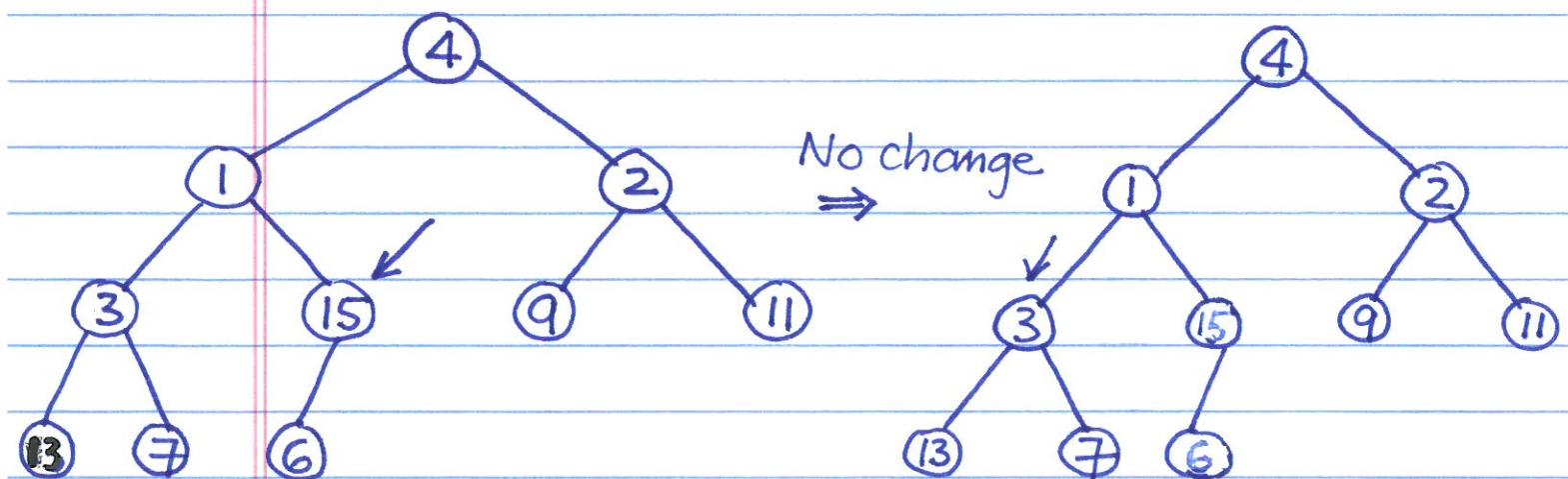
for $i = \lfloor \text{length}(A)/2 \rfloor$ down to 1

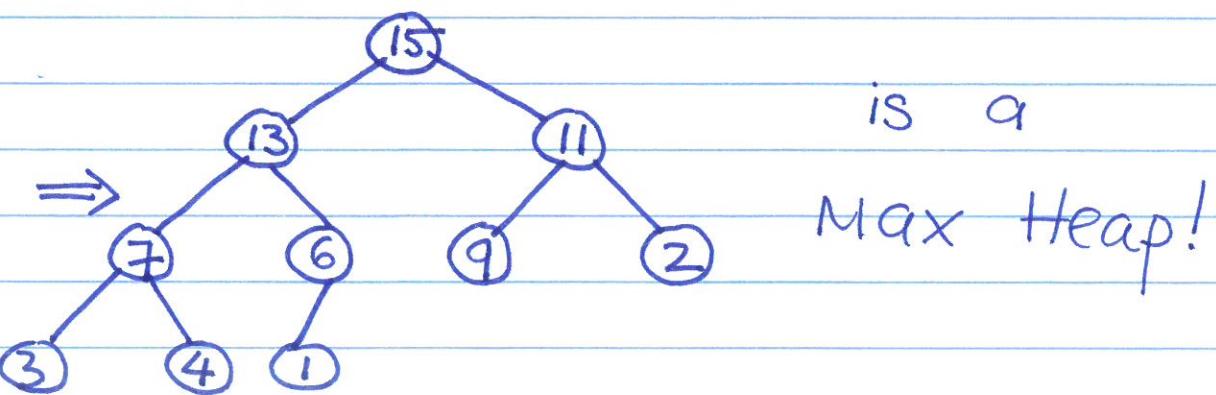
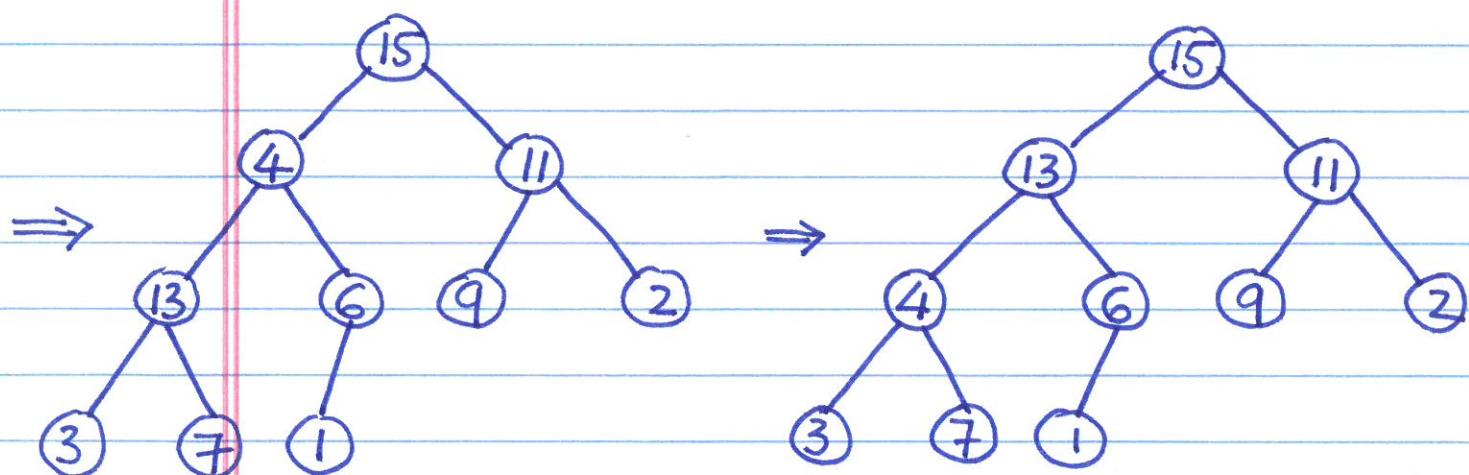
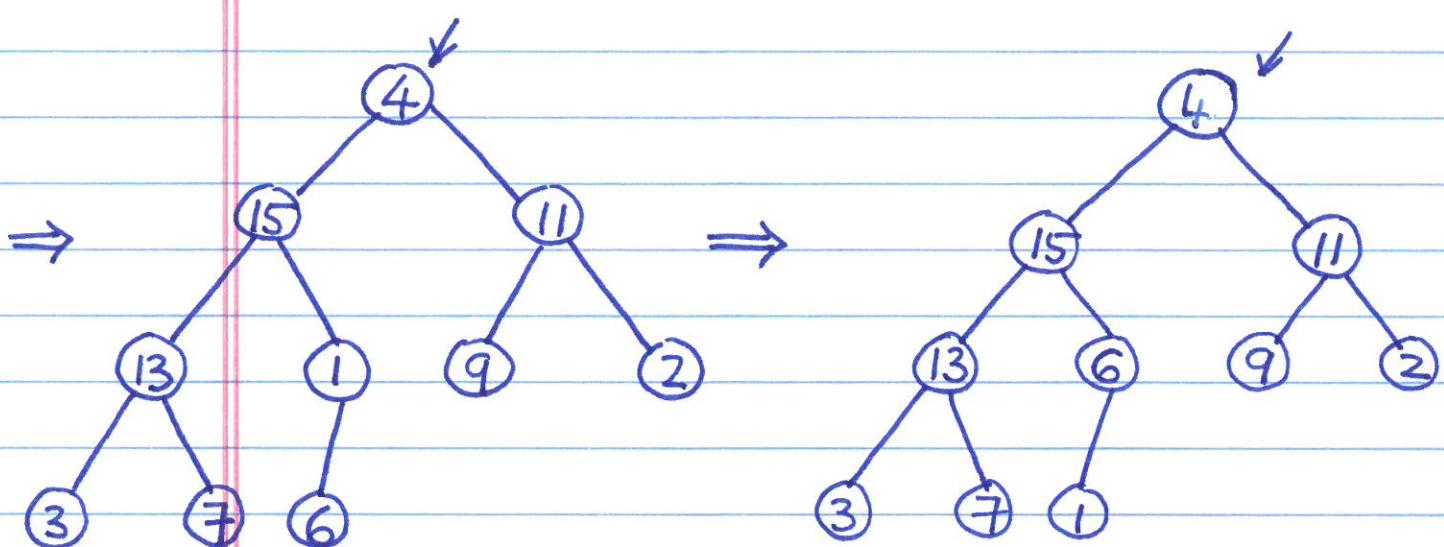
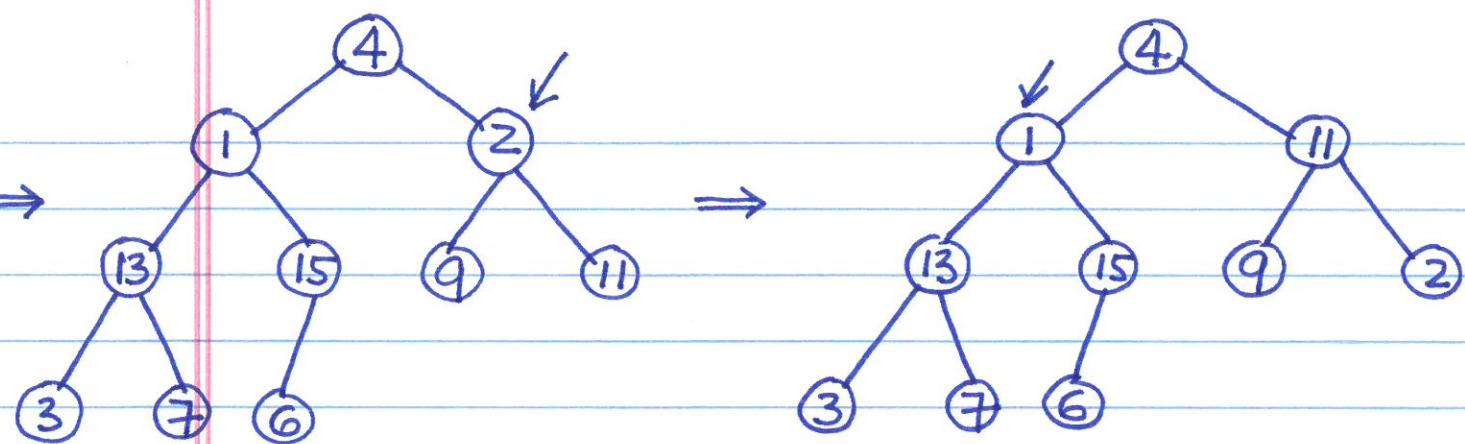
Max-Heapify (A, i)

Note: Building a heap bottom-up

Example:

4	1	2	3	15	9	11	13	7	6
---	---	---	---	----	---	----	----	---	---





Analysis of Build. Max. Heap.

Fact. An n -element heap has :

$$(1) \text{ height} = \lfloor \lg n \rfloor$$

$$(2) \leq \lceil n/2^{h+1} \rceil \text{ nodes at height } h.$$

Pf (1) ✓

(2) Ind. on h .

Observe:

$$\# \text{ internal nodes} \leq \# \text{ leaves}$$

$$\leq \# \text{ internal nodes} + 1$$

Thus,

$$\# \text{ nodes at height } 0 \leq \lceil n/2 \rceil$$

Now, assume $\# \text{ nodes at height } h$ is $\leq \lceil n/2^{h+1} \rceil$. Then

$\# \text{ nodes at height } h+1$ is

$$\leq \lceil \lceil n/2^{h+1} \rceil / 2 \rceil = \lceil \frac{n}{2^{h+2}} \rceil$$

□

Thus, running time of Build-Max-Heap:

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

$$= O(n)$$

since $\sum_{h=0}^{\infty} \frac{h}{2^h} = 2$

Note: $\sum_{i=0}^{\infty} x^i = \frac{1}{1-x}$
 diff $\Rightarrow \sum_{i=1}^{\infty} i x^{i-1} = \frac{1}{(1-x)^2}$
 mult. by $x \Rightarrow$

The Heapsort algor.

Heapsort (A)

$O(n)$ Build-Max-Heap (A);

$O(n)$ for $i = A.length \text{ downto } 2$
 exchange $A[1]$ with $A[i]$
 $\text{heap-size}(A) = \text{heap-size}(A)-1$

$O(\lg n)$ Max-Heapify (A, 1)

Heapsort's running time is:

$O(n \lg n)$

Quicksort.

Quicksort :

- has $\Theta(n^2)$ worst-case running time
- has $O(n \lg n)$ expected running time
- is in place
- is method of choice

Quicksort is based on divide-and-conquer technique :

- Divide : Partition $A[p..r]$ into subarrays $A[p..q-1], A[q..r]$ s.t.

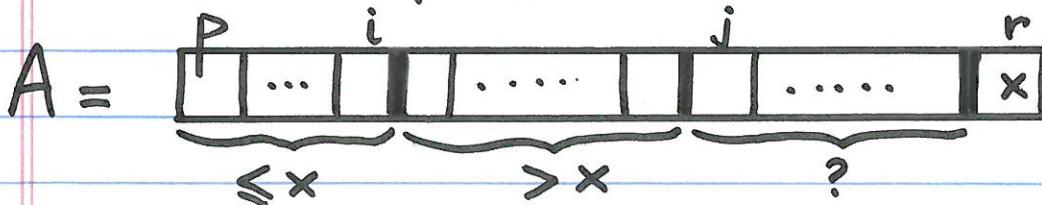
$$A[i] \leq A[q] \quad \forall p \leq i \leq q-1$$

$$A[q] \leq A[j] \quad \forall q+1 \leq j \leq r$$

- Conquer : Recursively sort $A[p..q-1], A[q+1..r]$
 (Now A is sorted)

Question : How to do partition ?

For $A[p..r]$ let $x = A[r]$ be pivot



- Search for j with $A[j] \leq x$
- Exchange $A[i+1]$ with $A[j]$
(this increases first region $\leq x$ by one element)

Partition(A, p, r)

$$x = A[r]$$

$$i = p - 1$$

for $j = p$ to $r-1$

if $A[j] \leq x$

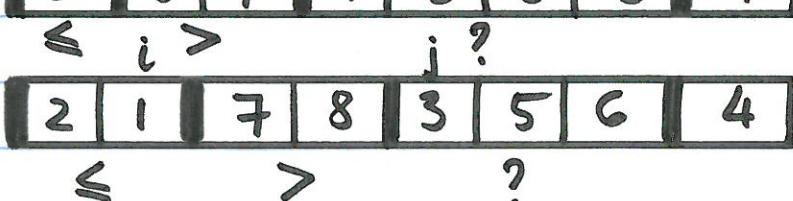
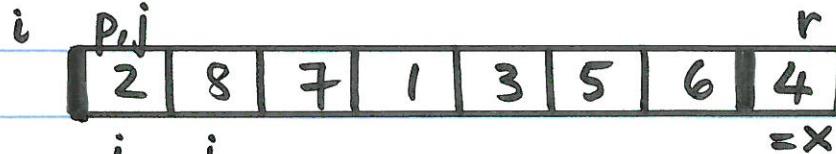
$$i = i + 1$$

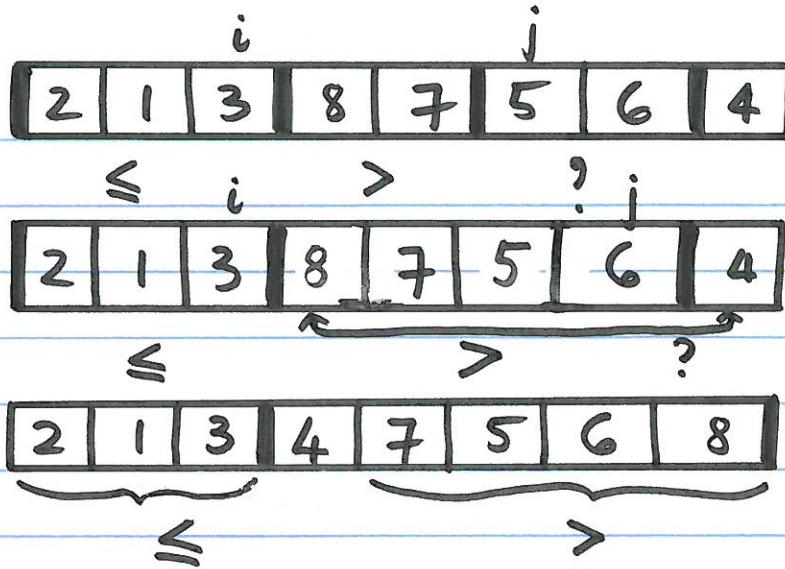
exchange $A[i]$ with $A[j]$

exchange $A[i+1]$ with $A[r]$

return $i+1$

Ex.





Complexity of Partition : $\Theta(n)$

Quicksort algorithm:

Quicksort (A, p, r)

if $p < r$

$q = \text{PARTITION}(A, p, r)$

$\text{QUICKSORT}(A, p, q-1)$

$\text{QUICKSORT}(A, q+1, r)$

Worst-case complexity of Quicksort

$$T(n) \leq \max_{0 \leq q \leq n-1} (T(q) + T(n-q-1)) + \Theta(n)$$

Guess : $T(n) \leq cn^2$ for some c .

We have:

$$T(n) \leq \max_{0 \leq q \leq n-1} (cq^2 + c(n-q-1)^2) + \Theta(n)$$

$$= c \cdot \max_{0 \leq q \leq n-1} (q^2 + (n-q-1)^2) + \Theta(n)$$

which achieves max. when $q=0, n-1$.

Thus,

$$\begin{aligned} T(n) &\leq c \cdot (n-1)^2 + \Theta(n) \\ &= c \cdot (n^2 - 2n + 1) + \Theta(n) \\ &= c \cdot n^2 - c(2n-1) + \Theta(n) \\ &\leq cn^2 \end{aligned}$$

for sufficiently large c . \square

On sorted arrays Quicksort requires $\Omega(n^2)$ steps. Thus,

$$T(n) = \Theta(n^2) \quad \square \square$$

Expected complexity of Quicksort

When partitioning is unbalanced Quicksort has worse performance.

Observation: Expected partition is balanced!

Assume:

- Array elements are distinct
- All permutations are equally likely.

$\Pr(\text{last element has rank } i) = \frac{1}{n}$

$\Rightarrow \text{Expected size of subarray is}$

$$\frac{1}{n} (0 + 1 + \dots + n-1) = \frac{n-1}{2} \quad \checkmark$$

Guess: $E[T(n)] \leq c \cdot n \cdot \lg n$

We have:

$$\begin{aligned} E[T(n)] &\leq \frac{1}{n} (E[T(0)] + E[T(n-1)] + c, n) \\ &\quad + \frac{1}{n} (E[T(1)] + E[T(n-2)] + c, n) \\ &\quad \vdots \\ &\quad + \frac{1}{n} (E[T(n-1)] + E[T(0)] + c, n) \\ &\leq c, n + \frac{2c}{n} \sum_{i=1}^{n-1} i \lg i \end{aligned}$$

Fact. $\sum_{i=1}^{n-1} i \lg i \leq \frac{1}{2} n^2 \lg n - \frac{1}{8} n^2$

Thus,

$$\begin{aligned} E[T(n)] &\leq \frac{2c}{n} \left(\frac{1}{2} n^2 \lg n - \frac{1}{8} n^2 \right) + c, n \\ &\leq cn \lg n - \frac{c}{4} n + c, n \\ &\leq cn \lg n \end{aligned}$$

for sufficiently large c . \square

Proof of Fact.

$$\begin{aligned}
 \sum_{i=1}^{n-1} i \lg i &= \sum_{i=1}^{\lceil n/2 \rceil - 1} i \lg i + \sum_{i=\lceil n/2 \rceil}^{n-1} i \lg i \\
 &\leq (\lg n - 1) \sum_{i=1}^{\lceil n/2 \rceil - 1} i + \lg n \sum_{i=\lceil n/2 \rceil}^{n-1} i \\
 &= \lg n \sum_{i=1}^{n-1} i - \sum_{i=1}^{\lceil n/2 \rceil - 1} i \\
 &\leq \frac{1}{2} n(n-1) \lg n - \frac{1}{2} \left(\frac{n}{2} - 1\right) \frac{n}{2} \\
 &\leq \frac{1}{2} n^2 \lg n - \frac{1}{8} n^2
 \end{aligned}$$

for $n \geq 2$.

□

A Lower Bound for Sorting

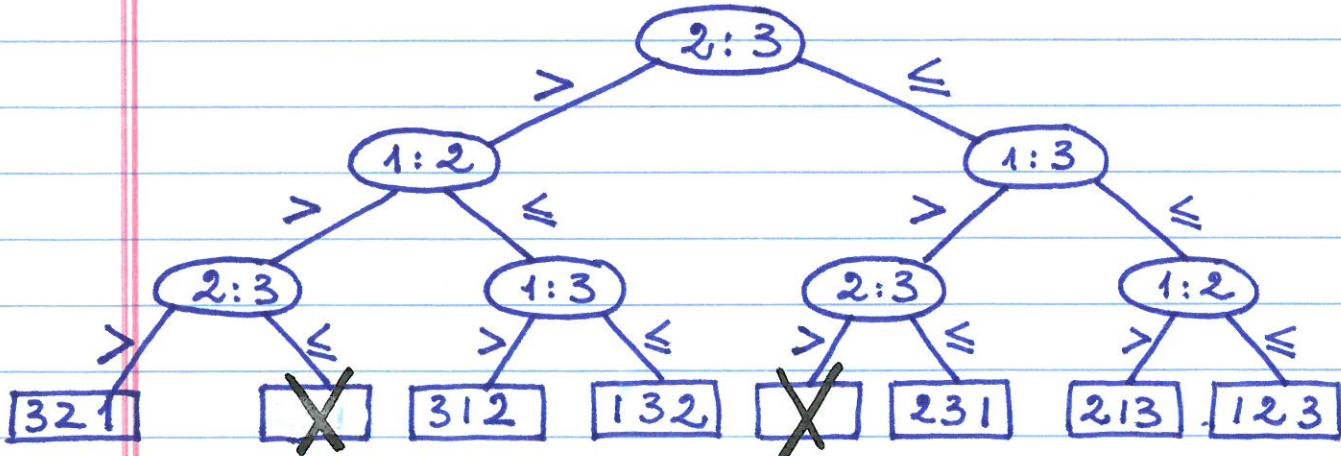
Observ. Sorting algor. are based on comparisons in general!

Question. Are there comparison-based sorting algor. with time complexity $O(n \cdot \log n)$?

Consider for example sorting on $A[1..3]$ by maximum selection!

- Suppose that we first compare $A[2]$ with $A[3]$: **2:3**

if $A[2] > A[3]$ then **1:2**
 else **1:3**



This is an example of a decision tree

A **decision tree** is a bin. tree whose nodes have labels of form $i:j$.

The 2 outgoing edges are labeled $>, \leq$.

Observ. On input arrays $A[1..n]$ comparison-based sorting algor. produce decision trees in a natural way.

The leaves correspond to permutations on $A[1..n]$

Def. A decision tree T solves the sorting problem of size n if there is a labeling of the leaves of T by permutations on $\{1, 2, \dots, n\}$ s.t:

For every $A[1..n]$ if a leaf labeled π is reached, then

$$A[\pi(1)] \leq A[\pi(2)] \leq \dots \leq A[\pi(n)]$$

For sorting problem of size n def.

$S(n)$ and $A(n)$ by:

Let T be a decision tree, π a permut.

Let $l_{\pi}^T =$ length of path from root to leaf labeled π in T .

$$S(n) := \min_T \max_{\underbrace{\pi}_{\text{depth of } T}} l_{\pi}^T$$

depth of decision tree of smallest depth

($S(n) \triangleq$ # of comparisons required by best comparison-based sort. alg.)

$$A(n) := \min_T \left(\frac{1}{n!} \sum_{\pi} l_{\pi}^T \right)$$

($A(n) \triangleq$ depth of decision tree with minimum average depth)

Goal: To obtain lower bounds for $S(n)$ and $A(n)$!

Claim: $S(n) = \Theta(n \lg n)$

Pf. Since there are $n!$ permutations of n elements, any decision tree for sorting n elements has $\geq n!$ leaves.

$$\text{Thus, } 2^{S(n)} \geq n!$$

$$\Rightarrow S(n) \geq \lceil \log n! \rceil$$

Stirling's approximation for $n!$ gives:

$$\log_2 n! = \left(n + \frac{1}{2}\right) \log_2 n - \frac{n}{\ln 2} + O(1)$$

$$\text{Hence, } S(n) = \Omega(n \lg n)$$

On the other hand, known sorting algorithms imply the existence of decision trees of depth $O(n \lg n)$ □

Next: A lower bound for $A(n)$.

Question: What is average depth $D(k)$ of a binary tree with k leaves?

Claim: $D(k) \geq \log_2 k \quad \forall k \geq 2$

Pf. By ind. on k .

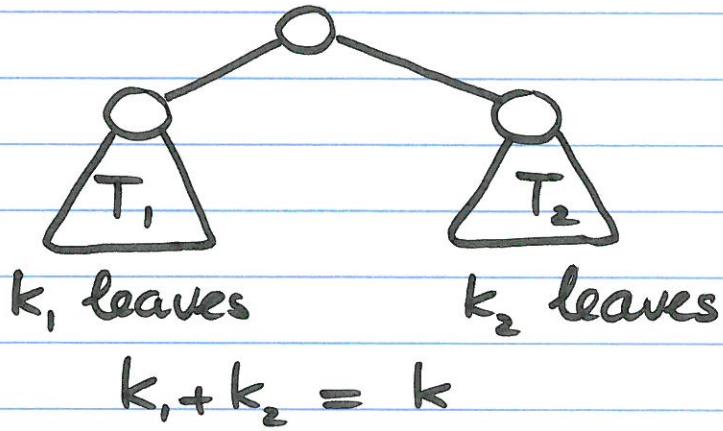
Basis. $k = 2 \checkmark$

Ind. Step. Suppose Claim holds true for for binary trees with $\leq k-1$ leaves. We show Claim for k .

Let T be a bin. tree with k leaves.

Without loss of generality (w.l.o.g.)

T is of form:



Average depth of T =

$$1 + \frac{k_1}{k_1 + k_2} (\text{aver. depth of } T_1) + \frac{k_2}{k_1 + k_2} (\text{av. depth of } T_2)$$

$$\geq \frac{k_1}{k} \log_2 k_1 + \frac{k_2}{k} \log_2 k_2 + 1 \quad (\text{by I.H.})$$

$$= \frac{1}{k} (k_1 \log_2 k_1 + k_2 \log_2 k_2 + k_1 + k_2)$$

$$= \frac{1}{k} \underbrace{(k_1 \log_2 2k_1 + k_2 \log_2 2k_2)}$$

$$f(k_1, k_2)$$

$f(k_1, k_2)$ achieves a ~~max.~~ ⁱⁿ max. when $k_1 = k_2 = k/2$. Thus,

Average depth of $T \geq \log_2 k \quad \square$

Theorem. Every comparison-based sorting algor. has worst-case and average-case complexities $\Omega(n \lg n)$ \square

Remark This lower bound doesn't hold for general models of computation such as RAMs or Turing machines.

Question. Is there some other type of sorting that can be solved faster?

Of course such algorithms cannot be based on comparisons!

Counting Sort.

Array to be sorted : $A[1..n]$

$A[1], \dots, A[n] \in \{0, 1, \dots, k\}$, $k = O(n)$

Counting-Sort (A, n, B, k) ($B[1..n] = \text{output}$)

- Let $C[0..k]$ be new array

- for $i=0$ to k
 $C[i] = 0$

1	2	3	4	5	6	7	8
2	5	3	0	2	3	0	3

$$k = 5$$

- for $j=1$ to n

$$C[A[j]] = C[A[j]] + 1$$

/* $C[i] = \# \{A[j] = i\}$ */

- for $i=1$ to k

$$C[i] = C[i] + C[i-1]$$

/* $C[i] = \# \{A[j] \leq i\}$ */

/* Now place $A[j]$ into its correct position in B */

$C \downarrow$	0	1	2	3	4	5
	2	0	2	3	0	1

$C \downarrow$	0	1	2	3	4	5
	2	2	4	7	7	8

for $j = n$ down to 1

$$B[C[A[j]]] = A[j]$$

$$C[A[j]] = C[A[j]] - 1$$

$B =$						3	
-------	--	--	--	--	--	---	--

$$B[7] = 3$$

$C \downarrow$	2	2	4	6	7	8
----------------	---	---	---	---	---	---

	1	2	3	4	5	6	7	8
A =	2	5	3	0	2	3	0	3

	1	2	3	4	5	6	7	8
B =						3		

	0	1	2	3	4	5
C =	2	2	4	6	7	8

Where to place $A[7]$? ($A[7]=0$)

	1	2	3	4	5	6	7	8
B =		0				3		

	0	1	2	3	4	5
C =	1	2	4	6	7	8

Where to place $A[6]=3$?

	1	2	3	4	5	6	7	8
B =		0			3	3		

	0	1	2	3	4	5
C =	1	2	4	5	7	8

Next we place $A[5]=2$ into B

	1	2	3	4	5	6	7	8
B =		0	1	2	3	3		

	0	1	2	3	4	5
C =	1	2	3	5	7	8

We now place $A[4]=0$ into B

	1	2	3	4	5	6	7	8
B =	0	0	1	2	3	3		

	0	1	2	3	4	5
C =	0	2	3	5	7	8

Counting Sort

for $i=0$ to k do $C[i]=0$

$O(k)$

for $j=1$ to n do $C[A[j]]=C[A[j]]+1$

$O(n)$

for $i=1$ to k do $C[i]=C[i]+C[i-1]$

$O(k)$

for $j=n$ down to 1 do

$O(n)$

$B[C[A[j]]] = A[j]$

$C[A[j]] = C[A[j]] - 1$

Total = $O(n+k)$

Remark: not in-place, stable.

Radix-Sort.

Elements of input array $A[1..n]$ are d -digit numbers

Radix-Sort (A, d)

for $i = 1$ to d do

use a stable sort to sort A on digit i

Example.

329	436	329	329
457	457	436	436
657 \Rightarrow	657 \Rightarrow	839 \Rightarrow	457
839	329	457	657
436	839	657	839

Proposition. If each digit can take on up to k values, Radix-Sort correctly sorts $A[1..n]$ in time $\Theta(d(n+k))$ if the stable sort used takes $\Theta(n+k)$ time. \square

Selection Problem

Input. An array $A[1..n]$ of distinct elements and an integer i , $1 \leq i \leq n$.

Output. The i .th smallest element of A .

Randomized Selection

Randomized-Select (A, p, r, i)

if $p=r$ then return $A[p]$

$q = \text{Random.-Partition}(A, p, r)$

$k = q - p + 1$; if $i = k$ then return $A[q]$;

if $i < k$ then return Random-Select ($A, p, q-1, i$)

else return Random-Select ($A, q+1, r, i-k$)

Claim. Expected running time of Randomized-Select is $O(n)$.

Pf. Let $T(n)$ denote running time of Randomized-Select. We show:

$$E[T(n)] \leq c.n \text{ for some } c.$$

Assume by induction that

$$E[T(k)] \leq c.k \text{ for } k < n.$$

We have:

$$\begin{aligned}
 E[T(n)] &\leq \frac{1}{n} \left(\sum_{k=1}^{n-1} E[T(\max(k, n-k))] \right) + cn \\
 &\leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + cn \\
 &\leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} c \cdot k + cn \quad (\text{by I.H.}) \\
 &= \frac{2c}{n} \left(\sum_{k=1}^{\lfloor n/2 \rfloor - 1} k - \sum_{k=1}^{\lfloor n/2 \rfloor} k \right) + cn \\
 &\leq \frac{2c}{n} \left(\frac{(n-1)n}{2} - \frac{(\lfloor n/2 \rfloor - 2)(\lfloor n/2 \rfloor - 1)}{2} \right) + cn \\
 &= \frac{2c}{n} \left(\frac{n^2 - n - \frac{n^2}{4} + \frac{3n}{2} - 2}{2} \right) + cn \\
 &= c \left(\frac{3n}{4} + \frac{1}{2} - \frac{2}{n} \right) + cn \\
 &\leq \frac{3cn}{4} + \frac{c}{2} + cn \\
 &= cn - \underbrace{\left(\frac{cn}{4} - \frac{c}{2} - cn \right)}_{\geq 0 \text{ for suitable } c} \\
 &\leq cn. \quad \square
 \end{aligned}$$

Linear Time Selection (Worst-Case)

Observation: Worst-case running time

of Select is good if partition is "balanced".

⇒ Choose pivot element for partition close to median, and find it fast!

Select: $A[1..n]$, $1 \leq i \leq n$.

(1) Divide $A[1..n]$ into $\lceil n/5 \rceil$ groups of 5 elements each. (Last group may have < 5 elements)

(2) Find median of each group

(3) Recursively use Select ~~the median~~^{compute} × of the $\lceil n/5 \rceil$ medians found in (2)

(4) Using x as pivot, partition A
Let k = size of left subarray

(5) Recursively use Select to find i -th smallest element

- in left subarray if $i < k$

- in right subarray if $i > k$

(if $i = k$ then return x)

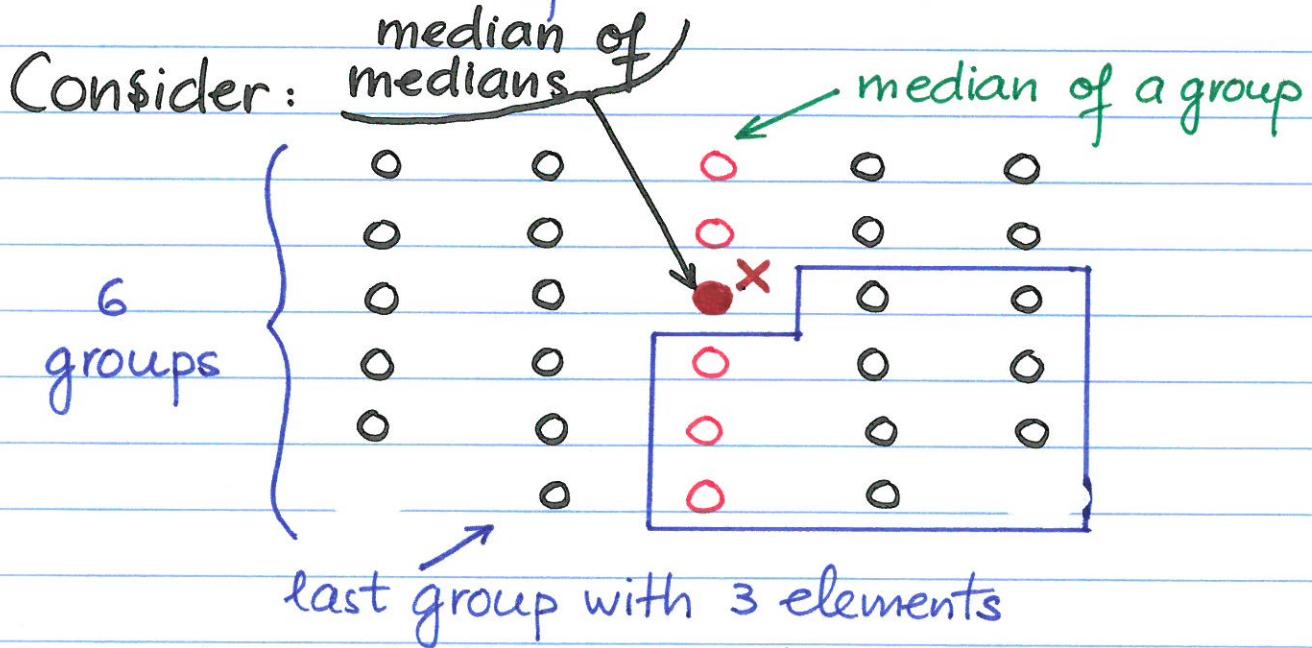
(Remark. If there are 2 possible medians,
take the smaller one)

Observation. Recursive call in Step (3)

applies to array of size $\lceil \frac{n}{5} \rceil$.

Question. What is size of subarrays
in recursive calls in Step (5) ?

What is # of elements $> x$?



The box containing elements $> x$
must have at least

$$3\left(\left\lceil \frac{1}{2} \left\lceil \frac{n}{5} \right\rceil \right\rceil - 2\right) \geq \frac{3n}{10} - 6$$

Half of
groups

- 1 group containing x
- last group with < 5
elements.

Similarly, the # of elements $< x$
is also $\geq \frac{3n}{10} - 6$.

Thus, subarrays in rec. call in Step(5)
must have size $\leq \frac{7n}{10} + 6$.

The recurrence for Select is:

$$T(n) \leq \begin{cases} O(1) & \text{if } n < 140 \\ T\left(\lceil \frac{n}{5} \rceil\right) + T\left(\frac{7n}{10} + 6\right) + O(n) & \text{otherwise} \end{cases}$$

Claim. $T(n) \leq cn$ for some c .

$$\begin{aligned} \text{Pf. } T(n) &\leq c\lceil \frac{n}{5} \rceil + c\left(\frac{7n}{10} + 6\right) + c_1 n \\ &\leq c\frac{n}{5} + c + \frac{7cn}{10} + 6c + c_1 n \\ &= \frac{9cn}{10} + 7c + c_1 n \\ &= cn + \left(-\frac{cn}{10} + 7c + c_1 n\right) \\ &\leq cn \end{aligned}$$

$$\text{if } -\frac{cn}{10} + 7c + c_1 n \leq 0$$

$$\text{or } c \geq 10c_1 \frac{n}{n-70} \quad \text{when } n > 70.$$

$$\text{Now if } n \geq 140 \Rightarrow \frac{n}{n-70} \leq 2.$$

Thus, if $c \geq 20c_1$, then

$$c \geq 10 \cdot c_1 \cdot 2 \geq 10c_1 \frac{n}{n-70} \quad \square$$