

CS 240 – Data Structures and Data Management

Module 3: Sorting and Randomized Algorithms

T. Biedl É. Schost O. Veksler

Based on lecture notes by many previous cs240 instructors

David R. Cheriton School of Computer Science, University of Waterloo

Winter 2021

Outline

- 1 Sorting and Randomized Algorithms
 - QuickSelect
 - Randomized Algorithms
 - QuickSort
 - Lower Bound for Comparison-Based Sorting
 - Non-Comparison-Based Sorting

Outline

- 1 Sorting and Randomized Algorithms
 - QuickSelect
 - Randomized Algorithms
 - QuickSort
 - Lower Bound for Comparison-Based Sorting
 - Non-Comparison-Based Sorting

Selection vs. Sorting

The **selection problem**: Given an array A of n numbers, and $0 \leq k < n$, find the element that would be at position k of the sorted array.

0	1	2	3	4	5	6	7	8	9
30	60	10	0	50	80	90	10	40	70

select(3) should return 30.

Special case: **median finding** = selection with $k = \lfloor \frac{n}{2} \rfloor$.

Selection can be done with heaps in time $\Theta(n + k \log n)$.

Median-finding with this takes time $\Theta(n \log n)$.

This is the same cost as our best sorting algorithms.

Question: Can we do selection in linear time?

The *quick-select* algorithm answers this question in the affirmative.

The encountered sub-routines will also be useful otherwise.

Crucial Subroutines

quick-select and the related algorithm *quick-sort* rely on two subroutines:

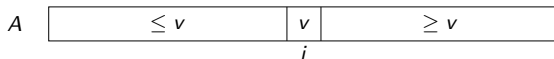
- *choose-pivot*(A): Return an index p in A . We will use the **pivot-value** $v \leftarrow A[p]$ to rearrange the array.

Simplest idea: Always select rightmost element in array

```
choose-pivot1( $A$ )  
1.  return  $A.size-1$ 
```

We will consider more sophisticated ideas later on.

- *partition*(A, p): Rearrange A and return **pivot-index** i so that
 - ▶ the pivot-value v is in $A[i]$,
 - ▶ all items in $A[0, \dots, i-1]$ are $\leq v$, and
 - ▶ all items in $A[i+1, \dots, n-1]$ are $\geq v$.



Partition Algorithm

Conceptually easy linear-time implementation:

partition(A, p)

A : array of size n , p : integer s.t. $0 \leq p < n$

1. Create empty lists *smaller*, *equal* and *larger*.
2. $v \leftarrow A[p]$
3. **for** each element x in A
4. **if** $x < v$ **then** *smaller.append*(x)
5. **else if** $x > v$ **then** *larger.append*(x)
6. **else** *equal.append*(x).
7. $i \leftarrow \text{smaller.size}$
8. $j \leftarrow \text{equal.size}$
9. Overwrite $A[0 \dots i-1]$ by elements in *smaller*
10. Overwrite $A[i \dots i+j-1]$ by elements in *equal*
11. Overwrite $A[i+j \dots n-1]$ by elements in *larger*
12. return i

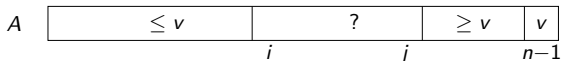
More challenging: partition **in place** (with $O(1)$ auxiliary space).

Efficient In-Place partition (Hoare)

$i=-1$	0	1	2	3	4	5	6	7	8	$j=9$
	30	60	10	0	50	80	90	20	40	$v=70$
	0	1	2	3	4	$i=5$	6	7	$j=8$	9
	30	60	10	0	50	80	90	20	40	$v=70$
	0	1	2	3	4	$i=5$	6	7	$j=8$	9
	30	60	10	0	50	40	90	20	80	$v=70$
	0	1	2	3	4	5	$i=6$	$j=7$	8	9
	30	60	10	0	50	40	90	20	80	$v=70$
	0	1	2	3	4	5	$i=6$	$j=7$	8	9
	30	60	10	0	50	40	20	90	80	$v=70$
	0	1	2	3	4	5	$j=6$	$i=7$	8	9
	30	60	10	0	50	40	20	90	80	$v=70$
	0	1	2	3	4	5	$j=6$	$i=7$	8	9
	30	60	10	0	50	40	20	70	80	90

Efficient In-Place partition (Hoare)

Idea: Keep swapping the outer-most wrongly-positioned pairs.



partition(A, p)

A : array of size n , p : integer s.t. $0 \leq p < n$

1. $swap(A[n-1], A[p])$
2. $i \leftarrow -1, j \leftarrow n-1, v \leftarrow A[n-1]$
3. **loop**
4. **do** $i \leftarrow i+1$ **while** $A[i] < v$
5. **do** $j \leftarrow j-1$ **while** $j \geq i$ and $A[j] > v$
6. **if** $i \geq j$ **then break** (goto 9)
7. **else** $swap(A[i], A[j])$
8. **end loop**
9. $swap(A[n-1], A[i])$
10. **return** i

Running time: $\Theta(n)$.

QuickSelect Algorithm

quick-select1(A, k)

A : array of size n , k : integer s.t. $0 \leq k < n$

1. $p \leftarrow \text{choose-pivot1}(A)$
2. $i \leftarrow \text{partition}(A, p)$
3. **if** $i = k$ **then**
4. **return** $A[i]$
5. **else if** $i > k$ **then**
6. **return** *quick-select1*($A[0, 1, \dots, i - 1], k$)
7. **else if** $i < k$ **then**
8. **return** *quick-select1*($A[i + 1, i + 2, \dots, n - 1], k - i - 1$)

Analysis of *quick-select1*

Worst-case analysis: Recursive call could always have size $n - 1$.

Recurrence given by $T(n) = \begin{cases} T(n - 1) + cn, & n \geq 2 \\ c, & n = 1 \end{cases}$

Solution: $T(n) = cn + c(n - 1) + c(n - 2) + \dots + c \cdot 2 + c \in \Theta(n^2)$

Best-case analysis: First chosen pivot could be the k th element
No recursive calls; total cost is $\Theta(n)$.

Average case analysis?

Sorting Permutations

- Need to take average running time over all inputs.
- How to characterize input of size n ?
(There are infinitely many sets of n numbers.)
- **Simplifying assumption:** All input numbers are *distinct*.
- Observe: quick-select1 would act the same on inputs
14, 2, 3, 6, 1, 11, 7 and
14, 2, 4, 6, 1, 12, 8
- The actual numbers do not matter, only their *relative order*.
- Characterize input via **sorting permutation**: the permutation π for which $A[\pi(0)] \leq A[\pi(1)] \leq \dots \leq A[\pi(n-1)]$.
- Assume all $n!$ sorting permutations are *equally likely*.

↪ Average cost is sum of costs for all permutations, divided by $n!$

Average-Case Analysis of *quick-select*1

$T^{\text{avg}}(n)$ = average-cost for selecting from size- n array
(Technically we should write $T^{\text{avg}}(n, k)$, but it turns out not to matter.)

$$= \frac{1}{n!} \sum_{l:\text{size}(l)=n} \text{running time for instance } l$$

(Use sorting-permutations instead and distinguish by pivot-index)

$$= \frac{1}{n!} \sum_{i=0}^{n-1} \sum_{\substack{\text{perm. } \pi \in \Pi_n \\ \text{has pivot-index } i}} \text{running time if sorting-permutation is } \pi$$

(Pivot-index $i \Rightarrow$ recurse in array of size $\leq \max\{i, n - i - 1\}$)

$$\stackrel{(*)}{\leq} \frac{1}{n!} \sum_{i=0}^{n-1} \left(\# \text{ such perm.} \right) \left(c \cdot n + T^{\text{avg}}(\max\{i, n - i - 1\}) \right)$$

(*) This clearly holds for T^{worst} . A non-trivial argument shows that over all permutations the run-time is the average. (No details.)

Average-Case Analysis of *quick-select*1

Claim: There are $(n - 1)!$ permutations for which the pivot-index is i .

Proof:

$$\begin{aligned} \text{So } T^{\text{avg}}(n) &\leq \frac{1}{n!} \sum_{i=0}^{n-1} (\# \text{ such perm.}) (c \cdot n + T^{\text{avg}}(\max\{i, n - i - 1\})) \\ &\leq \frac{1}{n!} \sum_{i=0}^{n-1} (n - 1)! (c \cdot n + T^{\text{avg}}(\max\{i, n - i - 1\})) \\ &= c \cdot n + \frac{1}{n} \sum_{i=0}^{n-1} T^{\text{avg}}(\max\{i, n - i - 1\}) \end{aligned}$$

Theorem: $T^{\text{avg}}(n) \in \Theta(n)$.

Proof:

Outline

1 Sorting and Randomized Algorithms

- QuickSelect
- Randomized Algorithms
- QuickSort
- Lower Bound for Comparison-Based Sorting
- Non-Comparison-Based Sorting

Randomized algorithms

A **randomized algorithm** is one which relies on some random numbers in addition to the input.

(Computers cannot generate randomness. We assume that there exists a *pseudo-random number generator (PRNG)*, a deterministic program that uses an initial value or *seed* to generate a sequence of seemingly random numbers. The quality of randomized algorithms depends on the quality of the PRNG!)

- The run-time will depend on the input and the random numbers used.
- **Goal:** Shift the dependency of run-time from what we can't control (the input) to what we *can* control (the random numbers).

No more bad instances, just unlucky numbers.

Expected running time

Define $T(I, R)$ to be the running time of a randomized algorithm \mathcal{A} for an instance I and the sequence of random numbers R .

The **expected running time** $T^{(\text{exp})}(I)$ for instance I is the expected value:

$$T^{(\text{exp})}(I) = \mathbf{E}[T(I, R)] = \sum_R T(I, R) \cdot \Pr[R]$$

- We could now take the *maximum* or the *average* over all instances of size n to define the **expected running time** of \mathcal{A} .
- But we usually design \mathcal{A} such that all instances of size n have the same expected run-time.
- Then maximum and average are the same, so we have

$$T^{(\text{exp})}(n) := \max_{\{I: \text{size}(I)=n\}} T^{(\text{exp})}(I) = \frac{\sum_{\{I: \text{size}(I)=n\}} T^{(\text{exp})}(I)}{|\{I: \text{size}(I)=n\}|}$$

We can still have good luck or bad luck, so occasionally we also discuss the worst that could happen, i.e., $\max_I \max_R T(I, R)$.

Randomized QuickSelect: Shuffle

Goal: Create a randomized version of *QuickSelect* for which all input has the same expected run-time.

First idea: Randomly permute the input first using *shuffle*:

```
shuffle(A)
```

```
A: array of size  $n$ 
```

1. **for** $i \leftarrow 1$ to $n - 1$ **do**
2. $swap(A[i], A[random(i + 1)])$

We assume the existence of a function *random*(n) that returns an integer uniformly from $\{0, 1, 2, \dots, n - 1\}$.

Expected cost becomes the same as the average cost: $\Theta(n)$.

Randomized QuickSelect: Random Pivot

Second idea: Change the pivot selection.

choose-pivot2(A)

1. **return** *random*(A.size)

quick-select2(A, k)

1. $p \leftarrow \text{choose-pivot2}(A)$
2. ...

With probability $\frac{1}{n}$ the random pivot has index i , so the analysis is just like that for the average-case. The expected running time is again $\Theta(n)$.

This is generally the fastest quick-select implementation.

There exists a variation that has worst-case running time $O(n)$, but it uses double recursion and is slower in practice. (\rightsquigarrow cs341)

Outline

1 Sorting and Randomized Algorithms

- QuickSelect
- Randomized Algorithms
- **QuickSort**
- Lower Bound for Comparison-Based Sorting
- Non-Comparison-Based Sorting

QuickSort

Hoare developed *partition* and *quick-select* in 1960.
He also used them to *sort* based on partitioning:

quick-sort1(A)

A : array of size n

1. **if** $n \leq 1$ **then return**
2. $p \leftarrow \text{choose-pivot1}(A)$
3. $i \leftarrow \text{partition}(A, p)$
4. *quick-sort1*($A[0, 1, \dots, i - 1]$)
5. *quick-sort1*($A[i + 1, \dots, n - 1]$)

QuickSort analysis

Define $T(n)$ to be the run-time for *quick-sort1* in a size- n array.

- $T(n)$ depends again on the pivot-index i .
- If we know i : $T(n) = \Theta(n) + T(i) + T(n - i - 1)$.
- **Worst-case analysis:** $i = 0$ or $n-1$ always. Then as for *quick-select*

$$T(n) = \begin{cases} T(n-1) + cn, & n \geq 2 \\ c, & n = 1 \end{cases}$$

for some constant $c > 0$. This resolves to $\Theta(n^2)$.

- **Best-case analysis:** $i = \lfloor \frac{n}{2} \rfloor$ or $\lceil \frac{n}{2} \rceil$ always. Then

$$T(n) = \begin{cases} T(\lfloor \frac{n-1}{2} \rfloor) + T(\lceil \frac{n-1}{2} \rceil) + cn & n \geq 2 \\ c, & n = 1 \end{cases}$$

Similar to *merge-sort*: This resolves to $\Theta(n \log n)$.

Average-case analysis of *quick-sort1*

Let $T^{\text{avg}}(n)$ be the *average-case* run-time for *quick-sort1* in a size- n array.

- As before, $(n - 1)!$ permutations have pivot-index i .
- As before, sub-arrays have size i and $n - i - 1$.
- As before, run-time for permutations average out.

$$\begin{aligned} \text{So } T^{\text{avg}}(n) &= \frac{1}{n!} \sum_{i=0}^{n-1} \sum_{\substack{\text{perm. } \pi \in \Pi_n \\ \text{has pivot-index } i}} \text{running time if sorting-perm. is } \pi \\ &\leq \frac{1}{n!} \sum_{i=0}^{n-1} (n-1)! \left(c \cdot n + T^{\text{avg}}(i) + T^{\text{avg}}(n-i-1) \right) \\ &= c \cdot n + \frac{1}{n} \sum_{i=0}^{n-1} (T^{\text{avg}}(i) + T^{\text{avg}}(n-i-1)) \end{aligned}$$

Theorem: $T^{\text{avg}}(n) \in \Theta(n \log n)$.

Proof:

Improvement ideas for QuickSort

- We can randomize by using *choose-pivot2*, giving $\Theta(n \log n)$ *expected time* for *quick-sort2*.
- The auxiliary space is $\Omega(\text{recursion depth})$.
 - ▶ This is $\Theta(n)$ in the worst-case.
 - ▶ It can be reduced to $\Theta(\log n)$ worst-case by recursing in smaller sub-array first and replacing the other recursion by a while-loop.
- One should stop recursing when $n \leq 10$.
Run InsertionSort at the end; this sorts everything in $O(n)$ time since all items are within 10 units of their required position.
- Arrays with many duplicates can be sorted faster by changing *partition* to produce three subsets

$\leq v$	$= v$	$\geq v$
----------	-------	----------
- Two programming tricks that apply in many situations:
 - ▶ Instead of passing full arrays, pass only the range of indices.
 - ▶ Avoid recursion altogether by keeping an explicit stack.

QuickSort with tricks

```
quick-sort3( $A, n$ )
1. Initialize a stack  $S$  of index-pairs with  $\{(0, n-1)\}$ 
2. while  $S$  is not empty
3.    $(l, r) \leftarrow S.pop()$ 
4.   while  $(r-l+1 > 10)$  do
5.      $p \leftarrow choose\_pivot2(A, l, r)$ 
6.      $i \leftarrow partition(A, l, r, p)$ 
7.     if  $(i-l > r-i)$  do
8.        $S.push((l, i-1))$ 
9.        $l \leftarrow i+1$ 
10.    else
11.       $S.push((i+1, r))$ 
12.       $r \leftarrow i-1$ 
13. InsertionSort( $A$ )
```

This is often the most efficient sorting algorithm in practice.

Outline

1 Sorting and Randomized Algorithms

- QuickSelect
- Randomized Algorithms
- QuickSort
- Lower Bound for Comparison-Based Sorting
- Non-Comparison-Based Sorting

Lower bounds for sorting

We have seen many sorting algorithms:

Sort	Running time	Analysis
Selection Sort	$\Theta(n^2)$	worst-case
Insertion Sort	$\Theta(n^2)$	worst-case
Merge Sort	$\Theta(n \log n)$	worst-case
Heap Sort	$\Theta(n \log n)$	worst-case
<i>quick-sort1</i>	$\Theta(n \log n)$	average-case
<i>quick-sort2</i>	$\Theta(n \log n)$	expected

Question: Can one do better than $\Theta(n \log n)$ running time?

Answer: Yes and no! *It depends on what we allow.*

- No: Comparison-based sorting lower bound is $\Omega(n \log n)$.
- Yes: Non-comparison-based sorting can achieve $O(n)$ (under restrictions!). → see below

The Comparison Model

In the **comparison model** data can only be accessed in two ways:

- comparing two elements
- moving elements around (e.g. copying, swapping)

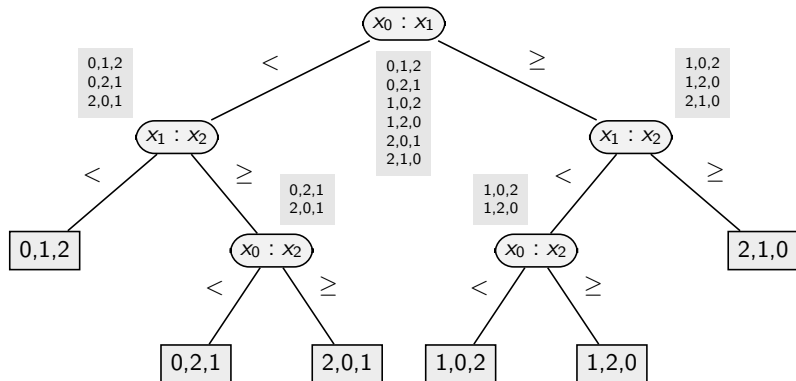
This makes very few assumptions on the kind of things we are sorting. We count the number of above operations.

All sorting algorithms seen so far are in the comparison model.

Decision trees

Comparison-based algorithms can be expressed as **decision tree**.

To sort $\{x_0, x_1, x_2\}$:



The permutations listed are the remaining possible *sorting permutations*, e.g. '0,1,2' means 'one possible remaining order is $x_0 \leq x_1 \leq x_2$ '.

Lower bound for sorting in the comparison model

Theorem. Any correct *comparison-based* sorting algorithm requires at least $\Omega(n \log n)$ comparison operations to sort n distinct items.

Proof.

Outline

1 Sorting and Randomized Algorithms

- QuickSelect
- Randomized Algorithms
- QuickSort
- Lower Bound for Comparison-Based Sorting
- Non-Comparison-Based Sorting

Non-Comparison-Based Sorting

- Assume keys are numbers in base R (R : **radix**)
 - ▶ $R = 2, 10, 128, 256$ are the most common.

Example ($R = 4$):

123	230	21	320	210	232	101
-----	-----	----	-----	-----	-----	-----

- Assume all keys have the same number m of digits.
 - ▶ Can achieve after padding with leading 0s.

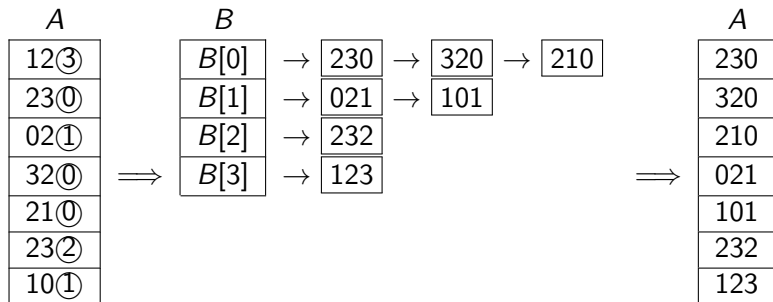
Example ($R = 4$):

123	230	021	320	210	232	101
-----	-----	-----	-----	-----	-----	-----

- Can sort based on individual digits.
 - ▶ How to sort 1-digit numbers?
 - ▶ How to sort multi-digit numbers based on this?

(Single-digit) Bucket Sort

Sort array A by last digit:



(Single-digit) Bucket Sort

Bucket-sort(A, d)

A : array of size n , contains numbers with digits in $\{0, \dots, R - 1\}$

d : index of digit by which we wish to sort

1. Initialize an array $B[0 \dots R - 1]$ of empty lists (**buckets**)
2. **for** $i \leftarrow 0$ to $n - 1$ **do**
3. Append $A[i]$ at end of $B[d^{\text{th}}$ digit of $A[i]$
4. $i \leftarrow 0$
5. **for** $j \leftarrow 0$ to $R - 1$ **do**
6. **while** $B[j]$ is non-empty **do**
7. move first element of $B[j]$ to $A[i++]$

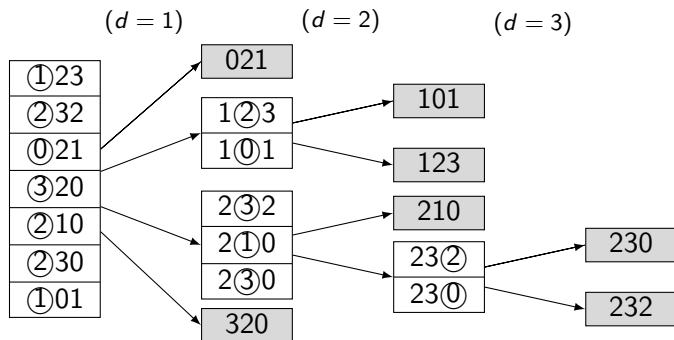
- Sorts numbers by single digit (specified by user).
- This is **stable**: equal items stay in original order.
- Run-time $\Theta(n + R)$, auxiliary space $\Theta(n + R)$
- It is possible to replace the lists by two auxiliary arrays of size R and $n \rightsquigarrow$ *count-sort* (no details).

MSD-Radix-Sort

Sorts array of m -digit radix- R numbers recursively:
sort by leading digit, then each group by next digit, etc.

```
MSD-Radix-sort( $A$ ,  $\ell \leftarrow 0$ ,  $r \leftarrow n-1$ ,  $d \leftarrow$  index of leading digit)
 $\ell, r$ : range of what we sort,  $0 \leq \ell, r \leq n-1$ 
1.   if  $\ell < r$ 
2.       bucket-sort( $A[\ell..r]$ ,  $d$ )
3.       if there are digits left // recurse in sub-arrays
4.            $\ell' \leftarrow \ell$ 
5.           while ( $\ell' < r$ ) do
6.               Let  $r' \geq \ell'$  be maximal s.t.  $A[\ell'..r']$  all have same  $d$ th digit
7.               MSD-Radix-sort( $A, \ell', r', d+1$ )
8.                $\ell' \leftarrow r' + 1$ 
```

MSD-Radix-Sort Example



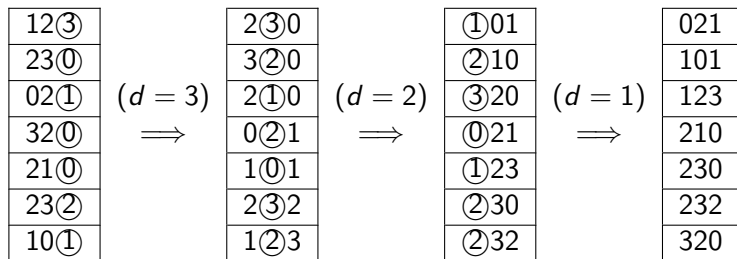
- Drawback of *MSD-Radix-Sort*: many recursions
- **Auxiliary space**: $\Theta(n + R + m)$ (for *bucket-sort* and recursion stack)
- **Run-time**: $\Theta(mnR)$ since we may have $\Theta(mn)$ subproblems.

LSD-Radix-Sort

LSD-radix-sort(A)

A : array of size n , contains m -digit radix- R numbers

1. **for** $d \leftarrow$ least significant to most significant digit **do**
2. *Bucket-sort*(A, d)



- Loop-invariant: A is sorted w.r.t. digits d, \dots, m of each entry.
- Time cost: $\Theta(m(n + R))$ Auxiliary space: $\Theta(n + R)$

Summary

- Sorting is an important and *very* well-studied problem
- Can be done in $\Theta(n \log n)$ time; faster is not possible for general input
- *HeapSort* is the only $\Theta(n \log n)$ -time algorithm we have seen with $O(1)$ auxiliary space.
- *MergeSort* is also $\Theta(n \log n)$, selection & insertion sorts are $\Theta(n^2)$.
- *QuickSort* is worst-case $\Theta(n^2)$, but often the fastest in practice
- *CountSort* and *RadixSort* achieve $o(n \log n)$ if the input is special

- Randomized algorithms can eliminate “bad cases”
- Best-case, worst-case, average-case, expected-case can all differ, but for well-design randomizations of algorithms, the expected case is the same as the average-case of the non-randomized algorithm.