

Order Statistics

David Kauchak
cs140
Fall 2024



1

Administrative

Group assignment

Assignment 2

Midterm 1

Assignment 3

- Available early next week
- Due on Sunday 9/29 (two weeks)



2

Administrative

Midterm 1

- Take-home midterm
 - Download from Gradescope
 - 1 hour and 15 minutes to take
 - Upload answers to Gradescope (in any reasonable format)
- Available Thursday (9/19) morning
- Must submit by Friday (9/20) 11:59pm
- If you want to take it during the normal class time, I'll be in the classroom
- Can use 2 pages of notes (and nothing else)
- Cover material up through order statistics (today)
- Post sample problems soon



3

Midterm 1 topics

Math foundations

- log properties
- properties of exponentials

Proofs by induction (weak, strong)

Big-O (theta and omega)

- Proving and disproving
- Categories and function ordering



4

Midterm 1 topics



Recurrences

- Generating (i.e., given a function/algorithm, write the recurrence)
- Solving: recurrence tree, substitution, master method

Sorting

- Insertion sort, Selection sort, Mergesort, Quicksort
- Runtimes

Order statistics

- median/selection
- run-time

Divide and conquer algorithms

5

Medians



The median of a set of numbers is the number such that half of the numbers are larger and half smaller

$A = [50, 12, 1, 97, 30]$

How might we calculate the median of a set?

Sort the numbers, then pick the $n/2$ element

$A = [1, 12, 30, 50, 97]$

runtime?

6

Medians



The median of a set of numbers is the number such that half of the numbers are larger and half smaller

$A = [50, 12, 1, 97, 30]$

How might we calculate the median of a set?

Sort the numbers, then pick the $n/2$ element

$A = [1, 12, 30, 50, 97]$

$\Theta(n \log n)$

7

Selection



More general problem:

find the k -th smallest element in an array

- i.e. element where exactly $k-1$ things are smaller than it
- aka the "selection" problem
- can use this to find the median if we want

Can we solve this in a similar way?

- Yes, sort the data and take the k th element
- $\Theta(n \log n)$

8

Can we do better?

Are we doing more work than we need to?

To get the k-th element (or the median) by sorting, we're finding *all* the k-th elements at once

We just want the one!

Often when you find yourself doing more work than you need to, there is a faster way (though not always)



9

selection problem

Our tools

- divide and conquer
- sorting algorithms
- other functions
 - merge
 - partition
 - binary search



10

Partition

Partition takes $\Theta(n)$ time and performs a similar operation

given an element $A[q]$, Partition can be seen as dividing the array into two sets:

- $\leq A[q]$
- $> A[q]$

Ideas?



11

An example

We're looking for the 5th smallest

5 2 34 9 17 2 1 34 18 5 3 2 1 6 5

If we called partition, what would be the in three sets?

≤ 5 :

> 5 :



12


An example

We're looking for the 5th smallest

5 2 34 9 17 2 1 34 18 5 3 2 1 6 5

≤ 5: 2 5 2 5 1 5 3 2 1
 > 5: 34 9 17 34 18 6

Does this help us?




13

An example

We're looking for the 5th smallest

5 2 34 9 17 2 1 34 18 5 3 2 1 6 5

≤ 5: 2 5 2 5 1 5 3 2 1 ← We know the 5th smallest has to be in this set
 > 5: 34 9 17 34 18 6




14

```

Selection(A, k, p, r)
  q ← Partition(A,p,r)
  if k = q
    Return A[q]
  else if k < q
    Return Selection(A, k, p, q-1)
  else // k > q
    Return Selection(A, k, q+1, r)
    
```

A: array of data
 k: find the kth smallest
 p,r: current span we're exploring (initially 1, len(A))




15

Selection(A, 3, 1, 8)

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

```

Selection(A, k, p, r)
  q ← Partition(A,p,r)
  if k = q
    Return A[q]
  else if k < q
    Selection(A, k, p, q-1)
  else // k > q
    Selection(A, k, q+1, r)
    
```



23

Selection(A, 3, 1, 8)

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

↑

```

Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
  Return A[q]
else if k < q
  Selection(A, k, p, q-1)
else // k > q
  Selection(A, k, q+1, r)
    
```

24

Selection(A, 3, 1, 8)

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

↑

```

Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
  Return A[q]
else if k < q
  Selection(A, k, p, q-1)
else // k > q
  Selection(A, k, q+1, r)
    
```

25

Selection(A, 3, 1, 5)

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

}

At each call, discard

part of the array

```

Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
  Return A[q]
else if k < q
  Selection(A, k, p, q-1)
else // k > q
  Selection(A, k, q+1, r)
    
```

26

Selection(A, 3, 1, 5)

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

↑

```

Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
  Return A[q]
else if k < q
  Selection(A, k, p, q-1)
else // k > q
  Selection(A, k, q+1, r)
    
```

27

Selection(A, 3, 3, 5)

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

↑

```

Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
    Return A[q]
else if k < q
    Selection(A, k, p, q-1)
else // k > q
    Selection(A, k, q+1, r)
    
```

28

Selection(A, 3, 3, 5)

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

```

Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
    Return A[q]
else if k < q
    Selection(A, k, p, q-1)
else // k > q
    Selection(A, k, q+1, r)
    
```

29

Selection(A, 3, 3, 5)

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

↑

```

Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
    Return A[q]
else if k < q
    Selection(A, k, p, q-1)
else // k > q
    Selection(A, k, q+1, r)
    
```

30

Selection(A, 3, 3, 4)

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

↑


```

Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
    Return A[q]
else if k < q
    Selection(A, k, p, q-1)
else // k > q
    Selection(A, k, q+1, r)
    
```

31

Selection(A, 3, 3, 4)

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



```


Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
    Return A[q]
else if k < q
    Selection(A, k, p, q-1)
else // k > q
    Selection(A, k, q+1, r)
    
```

32

Selection(A, 3, 3, 4)

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

↑



```


Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
    Return A[q]
else if k < q
    Selection(A, k, p, q-1)
else // k > q
    Selection(A, k, q+1, r)
    
```

33

Selection(A, 3, 3, 4)

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

↑



```


Selection(A, k, p, r)
q ← Partition(A,p,r)
if k = q
    Return A[q]
else if k < q
    Selection(A, k, p, q-1)
else // k > q
    Selection(A, k, q+1, r)
    
```

34

Running time of Selection?

Best case?
 We get lucky and the element at the end of the list is the kth smallest element!

One call to partition: $\theta(n)$



35

Running time of Selection?

Worst case?

Each call to Partition only reduces our search by 1



Recurrence?

$$T(n) = T(n-1) + \Theta(n)$$

$\Theta(n^2)$

36

Running time of Selection?

Worst case?

Each call to Partition only reduces our search by 1



When does this happen?

- sorted
- reverse sorted
- others...

37

How can randomness help us?

```

RSelection(A, k, p, r)
  q ← RPartition(A,p,r)
  if k = q
    Return A[q]
  else if k < q
    Return Selection(A, k, p, q-1)
  else // k > q
    Return Selection(A, k, q+1, r)
  
```

38

Running time of RSelection?

Best case

- $\Theta(n)$

Worst case

- Still $\Theta(n^2)$
- As with Quicksort, we can get unlucky

Average case?

39

Average case

Depends on how much data we throw away at each step

40

Average case

We'll call a partition "good" if the pivot falls within within the 25th and 75th percentile

- a "good" partition throws away at least a quarter of the data
- Or, each of the partitions contains at least 25% of the data

What is the probability of a "good" partition?

Half of the elements lie within this range and half outside, so 50% chance

41

Average case

Recall, that like Quicksort, we can absorb the cost of a constant number of "bad" partitions

42

Average case

On average, how many times will Partition need to be called before we get a good partition?

Let E be the number of times

Recurrence:

$$E = 1 + \frac{1}{2}E$$

half the time we get a good partition on the first try and half of the time, we have to try again

$$= 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} + \dots$$

$$= 2$$

43

Mathematicians and beer

An infinite number of mathematicians walk into a bar. The first one orders a beer. The second orders half a beer. The third, a quarter of a beer. The bartender says "You're all idiots", and pours two beers.



44

Average case

If on average we can get a "good" partition ever other time, what is the recurrence?

- recall the pivot of a "good" partition falls in the 25th and 75th percentile

46

Average case

If on average we can get a "good" partition ever other time, what is the recurrence?

- recall the pivot of a "good" partition falls in the 25th and 75th percentile

$$T(n) = T\left(\frac{3}{4}n\right) + \theta(n)$$

We throw away at least 1/4 of the data

roll in the cost of the "bad" partitions

47

Which is?

$$T(n) = T(3/4n) + \theta(n)$$

48

$T(n) = T(3/4n) + \Theta(n)$

if $f(n) = O(n^{\log_b a - \epsilon})$ for $\epsilon > 0$, then $T(n) = \Theta(n^{\log_b a})$
 if $f(n) = \Theta(n^{\log_b a})$, then $T(n) = \Theta(n^{\log_b a} \log n)$
 if $f(n) = \Omega(n^{\log_b a + \epsilon})$ for $\epsilon > 0$ and $af(n/b) \leq cf(n)$ for $c < 1$
 then $T(n) = \Theta(f(n))$

$a = 1$
 $b = 4/3$
 $f(n) = n$

$n^{\log_b a} = n^{\log_{4/3} 1} = n^0$

is $n = O(n^{0-\epsilon})$? **Case 3: $\Theta(n)$**
 is $n = \Theta(n^0)$? **Average case running time!**
 is $n = \Omega(n^{0+\epsilon})$?

49

Selection

Worst case: $\Theta(n^2)$

Best case: $\Theta(n)$

Average case: $\Theta(n)$

50

An aside...

Notice a trend?

$T(n) = T(n/2) + \Theta(n) \quad \Theta(n)$

$T(n) = T(3/4n) + \Theta(n) \quad \Theta(n)$

51

$T(n) = T(pn) + f(n)$

for $0 < p < 1$ and

$f(n) \notin \Theta(1)$ if $f(n) = O(n^{\log_b a - \epsilon})$ for $\epsilon > 0$, then $T(n) = \Theta(n^{\log_b a})$
 if $f(n) = \Theta(n^{\log_b a})$, then $T(n) = \Theta(n^{\log_b a} \log n)$
 if $f(n) = \Omega(n^{\log_b a + \epsilon})$ for $\epsilon > 0$ and $af(n/b) \leq cf(n)$ for $c < 1$
 then $T(n) = \Theta(f(n))$

$a = 1$
 $b = 1/p$
 $f(n) = f(n)$

$n^{\log_b a} = n^{\log_{1/p} 1} = n^0$

Case 3: $\Theta(f(n))$

52

Divide and conquer strategy



Split data in half and recurse on two halves

Assume it works! How do we get the answer to the entire problem?

- Often have to do a bit of extra work
- Be careful about solutions that could span/combine the two halves

53

Data structures



What is a data structure?

Way of storing data that **facilitates particular operations**

54

Data structures



What are some of the data structures that you've seen?

55

Data structures review



List

1. What operations do they support?

Ordered Set

2. What are they good at?

Priority Queue

3. How can we implement them? (Are there variations?)

Unordered Set

4. What are the runtimes for the operations? (Do variations matter?)

56

Lists

get/set at index

append (add at the end of the list)

remove

add/insert



57

Ordered Set

insert

remove

contains

next/prev (successor/predecessor)



58

Priority Queue

insert

remove

min/max



59

Unordered Set

insert

remove

contains



60