

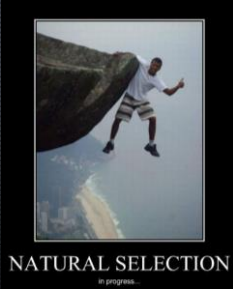
CS 341: ALGORITHMS

Lecture 4: divide & conquer III
Readings: see website

Trevor Brown
<https://student.cs.uwaterloo.ca/~cs341>
trevor.brown@uwaterloo.ca

1

THE SELECTION PROBLEM



NATURAL SELECTION

in progress...

2

THE SELECTION PROBLEM

- Input: An array **A** containing **n distinct** integer values, and an integer **k** between 1 and **n**
- Output: The **k**-th **smallest** integer in **A**
- Minimum** is a special case where $k = 1$
- Median** is a special case where $k = \frac{n}{2}$
- Maximum** is a special case where $k = n$
- Simple algorithm for solving selection?

3

Suppose we choose a **pivot** element **y** in the array **A**, and we **restructure** **A** so that all elements less than **y** precede **y** in **A**, and all elements greater than **y** occur after **y** in **A**. (This is exactly what is done in *Quicksort*, and it takes **linear time**.)

								y								
A	12	4	6	27	23	17	40	9								
	12	4	6	27	23	17	40	9								
Restructure(A, y)	12	4	6	17	9	23	27	40								

Number of elements on each side depend on the value **y**...

								y								
A	12	4	6	27	23	17	40	9								
	12	4	6	27	23	17	40	9								
Restructure(A, y)	4	6	12	27	23	17	40	9								

4

A after $Restructure(A, y)$

1	2	3	...	i_y	...	n	
12	4	6	17	9	23	27	40

A_L A_R
 Number of elements in this range = i_y

Recursive calls

- What's the **k**-th smallest element of **A**?
 - If $k = i_y$ then **y**
 - If $k < i_y$ then the **k**th smallest in A_L
 - If $k > i_y$ then the **(k - i_y)**th smallest in A_R

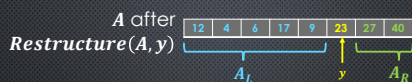
5

```

1 QuickSelect(k, A[1..n])
2   if n = 1 then return A[1] // base case
3
4   y = A[1] // pick an arbitrary pivot
5   (AL, AR, iy) = Restructure(A, y)
6
7   if k == iy return y
8   else if k < iy then return QuickSelect(k, AL)
9   else /* k > iy */ return QuickSelect(k - iy, AR)
10
11 Restructure(A[1..n], y)
12   AL = new array[1..n] // allocate more than enough
13   AR = new array[1..n] // to avoid need for expansion
14   nL = 0, nR = 0
15
16   for i = 1 .. n
17     if A[i] < y then AL[nL++] = A[i]
18     else A[i] > y then AR[nR++] = A[i]
19
20   return (AL, AR, nL+1) // nL+1 is the new index of y
    
```

6

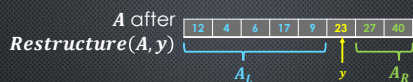
OVERLY OPTIMISTIC ANALYSIS ☺



- If $i_y = \frac{n}{2}$, then we recurse on $\sim \frac{n}{2}$ elements,
- If we could **always** recurse on $\frac{n}{2}$ elements then
 - We would get $T(n) = T(\frac{n}{2}) + \theta(n)$
 - Which would yield $a = 1, b = 2, y = 1, x = \log_2 1 = 0, y > x$ and $T(n) \in \theta(n^y) = \theta(n)$ by the Master theorem.

But we don't always recurse on $\frac{n}{2}$ elements!

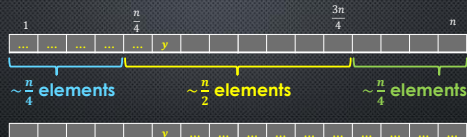
WORST-CASE ANALYSIS



- If we always get $i_y = 1$ and recurse on the right, then
- We get $T(n) = T(n-1) + \theta(n)$
- By the substitution method this is $\theta(n^2)$
- So, sometimes the pivot is good, sometimes it's bad...
- **What about the average case?**

AVERAGE-CASE ANALYSIS

- Definition: we say a pivot y is **good** if $i_y \in (\frac{n}{4}, \frac{3n}{4})$



- For **any good pivot** we recurse on at most $\frac{3n}{4}$ elements
- Probability of an arbitrary pivot being **good?** $p = 1/2$

Reducing the size of the subproblem by at least 1/4

PROOF SKETCH

- Since probability of a good pivot is $1/2$,
- on average, every **two** recursive calls, we will encounter a **good pivot**
- Encountering a good pivot reduces problem size to at most $\frac{3n}{4}$
- So, problem size is reduced to $\frac{3n}{4}$ after **expected linear work**
- **Average case recurrence:** $T(n) = T(\frac{3n}{4}) + \theta(n)$
 - $T(n) \in \theta(n)$

Here is a more rigorous proof of the average-case complexity: We say the algorithm is in **phase j** if the current subarray has size s , where

$$n \left(\frac{3}{4}\right)^{j+1} < s \leq n \left(\frac{3}{4}\right)^j$$

This is just for your notes, in case you want to know how you'd do this analysis formally

Let X_j be a **random variable** that denotes the amount of computation time occurring in phase j . If the pivot is in the middle half of the current subarray, then we transition from phase j to phase $j+1$. This occurs with probability $1/2$, so the expected number of recursive calls in phase j is 2. The computing time for each recursive call is linear in the size of the current subarray, so $E[X_j] \leq 2cn(3/4)^j$ (where $E[\cdot]$ denotes the expectation of a random variable). The total time of the algorithm is given by $X = \sum_{j \geq 0} X_j$. Therefore

$$E[X] = \sum_{j \geq 0} E[X_j] \leq 2cn \sum_{j \geq 0} (3/4)^j = 8cn \in O(n).$$

$$\sum_{k=0}^{\infty} ar^k = \frac{a}{1-r}, \text{ for } |r| < 1.$$

TAKING SELECTION FURTHER

- We just showed:
 - QuickSelect with **average case** runtime in $\theta(n)$
- Next up:
 - Median-of-medians QuickSelect (MOMQuickSelect)
 - **worst case** runtime in $\theta(n)$

Relies on getting a good pivot within $\theta(1)$ recursive calls: **on average**

Must get a good pivot within $\theta(1)$ recursive calls **always**

The algorithm we will see picks a **good pivot** in **every** recursive call

HIGH LEVEL ALGORITHM

- Similar to QuickSelect
 - Choose** a pivot
 - Move smaller elements to the left of the pivot, and larger elements to the right of the pivot
 - Recursively call MOMQuickSelect on one subarray
- Only difference is **how** we choose the pivot
 - Always** want to pick a **good pivot**

13

ALWAYS PICKING A GOOD PIVOT

Example input $A[1..50]$: 11, 38, 6, 21, 20, 17, 14, 9, 7, 5, 8, 34, 49, 47, 28, 18, 44, 31, 46, 48, 27, 4, 2, 50, 23, 45, 3, 13, 43, 22, 10, 32, 35, 41, 24, 1, 30, 12, 15, 26, 16, 19, 36, 33, 37, 39, 25, 40, 29, 42

Group into rows of 5					Find median of each row				
11	38	6	21	20	11	38	6	21	20
17	14	9	7	5	17	14	9	7	5
8	34	49	47	28	8	34	49	47	28
18	44	31	46	48	18	44	31	46	48
27	4	2	50	23	27	4	2	50	23
45	3	13	43	22	45	3	13	43	22
10	32	35	41	24	10	32	35	41	24
1	30	12	15	26	1	30	12	15	26
16	19	36	33	37	16	19	36	33	37
39	25	40	29	42	39	25	40	29	42

Build array of medians: 20, 9, 34, 44, 23, 22, 32, 15, 33, 39

Time complexity?

Recursively find the median of these medians: 23

Recursive problem size?

Time complexity for this step?

HOW GOOD IS THE PIVOT 23?

elements ≤ 23 is at least 3(5). This is at least 3/10ths of our 50-element input, or $3n/10$.

Recall: median of each row

11	38	6	21	20
17	14	9	7	5
8	34	49	47	28
18	44	31	46	48
27	4	2	50	23
45	3	13	43	22
10	32	35	41	24
1	30	12	15	26
16	19	36	33	37
39	25	40	29	42

Imagine sorting each row:

6	11	20	21	38
5	7	9	14	17
1	12	15	26	30
6	11	20	21	38
3	13	22	43	45
2	4	23	27	50
10	24	32	35	41
16	19	33	36	37
8	28	34	47	49
25	29	39	40	42
4	31	44	46	48

Then ordering the medians:

5	7	9	14	17
1	12	15	26	30
6	11	20	21	38
3	13	22	43	45
2	4	23	27	50
10	24	32	35	41
16	19	33	36	37
8	28	34	47	49
25	29	39	40	42
4	31	44	46	48

elements ≥ 23 is at least 3(6). This is $> 3/10$ ths of our 50-element input.

So, after restructuring, pivot 23 must have at least $3n/10$ elements before and after it

This is a good pivot!

We recurse on A_L or A_R , and both have size at most $7n/10$

MOMQuickSelect($k=11, n=14, A$)

```

1 MOMQuickSelect(k, n, A)
2 // base case
3 if n <= 14 then sort(A) and return A[k]
4
5 // divide and conquer to find medians
6 r = (n-5) / 10
7 medians[1..(2*r+1)] = new array
8 for i = 1..(2*r+1)
9   B[1..5] = A[(5*(i-1)+1)..(5*i)]
10  sort(B)
11  medians[i] = B[3]
12
13 y = MOMQuickSelect(r+1, 2*r+1, medians)
14
15 // divide and conquer to find rank k
16 (AL, AR, iy) = Restructure(A, y)
17 if k == iy then return y
18 else if k < iy then return MOMQuickSelect(k, iy-1, AL)
19 else /* k > iy */ then return MOMQuickSelect(k-iy, n-iy, AR)
    
```

11, 38, 6, 21, 20, 17, 14, 9, 7, 5, 8, 34, 49, 47, 28, 18, 44, 31, 46, 48, 27, 4, 2, 50, 23, 45, 3, 13, 43, 22, 10, 32, 35, 41, 49, 47

5, 6, 7, 9, 7, 11, 14, 17, 20, 21, 34, 38, 47, 49

18

MOMQuickSelect($k=11, n=21, A$)

```

1 MOMQuickSelect(k, n, A)
2 // base case
3 if n <= 14 then sort(A) and return A[k]
4
5 // divide and conquer to find medians
6 r = (n-5) / 10
7 medians[1..(2*r+1)] = new array
8 for i = 1..(2*r+1)
9   B[1..5] = A[(5*(i-1)+1)..(5*i)]
10  sort(B)
11  medians[i] = B[3]
12
13 y = MOMQuickSelect(r+1, 2*r+1, medians)
14
15 // divide and conquer to find rank k
16 (AL, AR, iy) = Restructure(A, y)
17 if k == iy then return y
18 else if k < iy then return MOMQuickSelect(k, iy-1, AL)
19 else /* k > iy */ then return MOMQuickSelect(k-iy, n-iy, AR)
    
```

11, 38, 6, 21, 20, 17, 14, 9, 7, 5, 8, 34, 49, 47, 28, 18, 44, 31, 46, 48, 27

$r = \frac{21-5}{10} = 1$ Not considering at most 9 elements

$B = \begin{bmatrix} 11 & 38 & 6 & 21 & 20 \\ 17 & 14 & 9 & 7 & 5 \\ 8 & 34 & 49 & 47 & 28 \\ 18 & 44 & 31 & 46 & 48 \\ 27 & 4 & 2 & 50 & 23 \end{bmatrix}$

$sort(B) = \begin{bmatrix} 6 & 11 & 20 & 21 & 38 \\ 5 & 7 & 9 & 14 & 17 \\ 1 & 12 & 15 & 26 & 30 \\ 3 & 13 & 22 & 43 & 45 \\ 2 & 4 & 23 & 27 & 50 \end{bmatrix}$

medians = 20, 9, 34

$y = MOMQuickSelect(2, 3, [20, 9, 34]) \rightarrow 20$

17

MOMQuickSelect($k=11, n=21, A$)

```

1 MOMQuickSelect(k, n, A)
2 // base case
3 if n <= 14 then sort(A) and return A[k]
4
5 // divide and conquer to find medians
6 r = (n-5) / 10
7 medians[1..(2*r+1)] = new array
8 for i = 1..(2*r+1)
9   B[1..5] = A[(5*(i-1)+1)..(5*i)]
10  sort(B)
11  medians[i] = B[3]
12
13 y = MOMQuickSelect(r+1, 2*r+1, medians)
14
15 // divide and conquer to find rank k
16 (AL, AR, iy) = Restructure(A, y)
17 if k == iy then return y
18 else if k < iy then return MOMQuickSelect(k, iy-1, AL)
19 else /* k > iy */ then return MOMQuickSelect(k-iy, n-iy, AR)
    
```

11, 38, 6, 21, 20, 17, 14, 9, 7, 5, 8, 34, 49, 47, 28, 18, 44, 31, 46, 48, 27

Restructure($A, y=20$) =

$A_L = [11, 6, 17, 14, 9, 7, 5, 8, 18]$

$A_R = [38, 21, 34, 49, 47, 28, 44, 31, 46, 48, 27]$

$i_y = |A_L| + 1 = 10$

$k = 11 > i_y = 10$

$k - i_y = 1$ $n - i_y = 10$

MOMQuickSelect($1, 10, A_R$) = 21

18

Runtime? (unit cost)

```

// base case
if n <= 14 then sort(A) and return A[k]

// divide and conquer to find medians
r = (n-5) / 10
medians[1..(2*r+1)] = new array
for i = 1..(2*r+1)
  B[i..5] = A[(5*(i-1)+1)..(5*i)]
  sort(B)
  medians[i] = B[3]

y = MOMQuickSelect(r+1, 2*r+1, medians)

// divide and conquer to find rank k
(AL, AR, iy) = Restructure(A, y)
if k == iy then return y
else if k < iy then return MOMQuickSelect(k, iy-1, AL)
else /* k > iy */ then return MOMQuickSelect(k-iy, n-iy, AR)
    
```

Rows B ordered by medians

1	7	9	14	17	
2	17	15	26	30	
3	4		27	50	
4	25	29	39	42	
5	18	31	44	46	48

Annotations: $\leq y$ (rows 1-3), $\geq y$ (rows 4-5). $3(r+1)$ elements $\leq y$, $3(r+1)$ elements $\geq y$. So problem size shrinks by at least $3(r+1)$. Observe $n = 10r + 5$.

HOW MUCH DOES THE PROBLEM SHRINK?

- Shrinks by at least $3(r+1)$
- Problem size $\approx n = 10r + 5$
- Subproblem size $\leq n - \text{Shrink} = n - 3(r+1)$
 - $= 10r + 5 - 3r - 3 = 7r + 2$
- Express in terms of n using $r = \lfloor \frac{n-5}{10} \rfloor$
 - Subproblem size $\leq 7 \lfloor \frac{n-5}{10} \rfloor + 2 \leq 7 \frac{n-5}{10} + 2$
 - $= \frac{7n}{10} - 7 \left(\frac{5}{10}\right) + 2 = \frac{7n}{10} - \frac{3}{2} \leq \frac{7n}{10}$

Time complexity

```

// base case
if n <= 14 then sort(A) and return A[k]

// divide and conquer to find medians
r = (n-5) / 10
medians[1..(2*r+1)] = new array
for i = 1..(2*r+1)
  B[i..5] = A[(5*(i-1)+1)..(5*i)]
  sort(B)
  medians[i] = B[3]

y = MOMQuickSelect(r+1, 2*r+1, medians)

// divide and conquer to find rank k
(AL, AR, iy) = Restructure(A, y)
if k == iy then return y
else if k < iy then return MOMQuickSelect(k, iy-1, AL)
else /* k > iy */ then return MOMQuickSelect(k-iy, n-iy, AR)
    
```

$T(n) \in O(n) + T(n/5) + T(7n/10)$ if $n \geq 15$
 $T(n) \in O(1)$ if $n \leq 14$

The key fact is that $1/5 + 7/10 = 9/10 < 1$.
 $T(n) \in O(n) + T(n/5) + T(7n/10)$ if $n \geq 15$
 $T(n) \in O(1)$ if $n \leq 14$

The fact that $T(n) \in \Theta(n)$ can be proven formally using guess-and-check (induction) or informally using the recursion tree method.

- Let $T(n) = c'n + T(\frac{n}{5}) + T(\frac{7n}{10})$ where $c' > 0$
- Want to prove: $T(n) = cn$ for some $c > 0$
- Note c and c' are independent constants
 - c' comes from the work at each level of recursion being $O(n)$
 - c is a positive constant we are trying to show exists
- I.H.: Suppose $\exists c > 0 : T(n') = c'n'$ for $15 \leq n' < n$
- $T(n) = c'n + c \frac{n}{5} + c \frac{7n}{10}$ (by inductive hypoth.)
- $T(n) = cn$** (want this to be true)
- $\Leftrightarrow c'n + c \frac{n}{5} + c \frac{7n}{10} = cn$ (equivalently)
- $\Leftrightarrow c' + c \frac{1}{5} + c \frac{7}{10} = c \Leftrightarrow c = 10c'$ (by algebra)