#### Introduction to Algorithms 6.046J/18.401J



#### **LECTURE 6** Order Statistics

- Randomized divide and conquer
- Analysis of expected time
- Worst-case linear-time order statistics
- Analysis

#### **Prof. Erik Demaine**

September 28, 2005

Copyright © 2001-5 by Erik D. Demaine and Charles E. Leiserson



### **Order statistics**

Select the *i*th smallest of *n* elements (the element with *rank i*).

- *i* = 1: *minimum*;
- *i* = *n*: *maximum*;
- $i = \lfloor (n+1)/2 \rfloor$  or  $\lceil (n+1)/2 \rceil$ : *median*.

*Naive algorithm*: Sort and index *i*th element. Worst-case running time =  $\Theta(n \lg n) + \Theta(1)$ =  $\Theta(n \lg n)$ ,

using merge sort or heapsort (not quicksort).



#### Randomized divide-andconquer algorithm

RAND-SELECT(A, p, q, i)  $\triangleright$  ith smallest of A[p... *Q* if p = q then return A[p] $r \leftarrow \text{RAND-PARTITION}(A, p, q)$  $\triangleright k = \operatorname{rank}(A[r])$  $k \leftarrow r - p + 1$ if i = k then return A[r]if i < kthen return RAND-SELECT(A, p, r-1, i) else reti k n RAND-SELECT(A, r + 1, q, i - k)



September 28, 2005 Copyright © 2001-5 by Erik D. Demaine and Charles E. Leiserson

L6.3



#### Select the i = 7th smallest:

#### Partition:





# Intuition for analysis

(All our analyses today assume that all elements are distinct.)

Lucky:  $T(n) = T(9n/10) + \Theta(n)$  $= \Theta(n)$ 

Unlucky:  $T(n) = T(n-1) + \Theta(n)$  $= \Theta(n^2)$   $n^{\log_{10/9}1} = n^0 = 1$ CASE 3

arithmetic series

#### Worse than sorting!



# Analysis of expected time

The analysis follows that of randomized quicksort, but it's a little different.

Let T(n) = the random variable for the running time of RAND-SELECT on an input of size n, assuming random numbers are independent.

For k = 0, 1, ..., n-1, define the *indicator random variable* 

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



# **Analysis (continued)**

To obtain an upper bound, assume that the *i*th element always falls in the larger side of the partition:

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

$$= \sum_{k=0}^{n-1} X_k (T(\max\{k, n-k-1\}) + \Theta(n)).$$



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

Take expectations of both sides.



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

Linearity of expectation.



$$\begin{split} E[T(n)] &= E\bigg[\sum_{k=0}^{n-1} X_k \big( T(\max\{k, n-k-1\}) + \Theta(n) \big) \bigg] \\ &= \sum_{k=0}^{n-1} E\big[ X_k \big( T(\max\{k, n-k-1\}) + \Theta(n) \big) \big] \\ &= \sum_{k=0}^{n-1} E\big[ X_k \big] \cdot E\big[ T(\max\{k, n-k-1\}) + \Theta(n) \big] \end{split}$$

# Independence of $X_k$ from other random choices.



$$\begin{split} E[T(n)] &= E\left[\sum_{k=0}^{n-1} X_k \left( T(\max\{k, n-k-1\}) + \Theta(n) \right) \right] \\ &= \sum_{k=0}^{n-1} E\left[ X_k \left( T(\max\{k, n-k-1\}) + \Theta(n) \right) \right] \\ &= \sum_{k=0}^{n-1} E\left[ X_k \right] \cdot E\left[ T(\max\{k, n-k-1\}) + \Theta(n) \right] \\ &= \frac{1}{n} \sum_{k=0}^{n-1} E\left[ T(\max\{k, n-k-1\}) \right] + \frac{1}{n} \sum_{k=0}^{n-1} \Theta(n) \end{split}$$

#### Linearity of expectation; $E[X_k] = 1/n$ .



$$E[T(n)] = E\left[\sum_{k=0}^{n-1} X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)\right]$$
  
=  $\sum_{k=0}^{n-1} E[X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)]$   
=  $\sum_{k=0}^{n-1} E[X_k] \cdot E[T(\max\{k, n-k-1\}) + \Theta(n)]$   
=  $\frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\})] + \frac{1}{n} \sum_{k=0}^{n-1} \Theta(n)$   
 $\leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \Theta(n)$  Upper terms appear twice.



Use fact:

#### Hairy recurrence

(But not quite as hairy as the quicksort one.)

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

**Prove:**  $E[T(n)] \leq cn$  for constant c > 0.

• The constant *c* can be chosen large enough so that  $E[T(n)] \leq cn$  for the base cases.

$$\sum_{k=\lfloor n/2 \rfloor}^{n-1} k \le \frac{3}{8}n^2 \quad \text{(exercise).}$$



 $E[T(n)] \leq \frac{2}{n} \sum_{k=|n/2|}^{n-1} ck + \Theta(n)$ 

#### Substitute inductive hypothesis.



$$E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} ck + \Theta(n)$$
$$\leq \frac{2c}{n} \left(\frac{3}{8}n^2\right) + \Theta(n)$$

Use fact.



$$E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} ck + \Theta(n)$$
$$\leq \frac{2c}{n} \left(\frac{3}{8}n^2\right) + \Theta(n)$$
$$= cn - \left(\frac{cn}{4} - \Theta(n)\right)$$

Express as *desired – residual*.



$$E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} ck + \Theta(n)$$
$$\leq \frac{2c}{n} \left(\frac{3}{8}n^2\right) + \Theta(n)$$
$$= cn - \left(\frac{cn}{4} - \Theta(n)\right)$$
$$\leq cn,$$

# if *c* is chosen large enough so that cn/4 dominates the $\Theta(n)$ .



# Summary of randomized order-statistic selection

- Works fast: linear expected time.
- Excellent algorithm in practice.
- But, the worst case is *very* bad:  $\Theta(n^2)$ .
- *Q*. Is there an algorithm that runs in linear time in the worst case?
- *A.* Yes, due to Blum, Floyd, Pratt, Rivest, and Tarjan [1973].

#### **IDEA:** Generate a good pivot recursively.



# Worst-case linear-time order statistics

#### Select(i, n)

- 1. Divide the *n* elements into groups of 5. Find the median of each 5-element group by rote.
- 2. Recursively SELECT the median *x* of the  $\lfloor n/5 \rfloor$  group medians to be the pivot.
- 3. Partition around the pivot x. Let  $k = \operatorname{rank}(x)$ .
- **4.** if i = k then return x
  - elseif i < k

then recursively SELECT the *i*th smallest element in the lower part else recursively SELECT the (i-k)th smallest element in the upper part Same as RAND-SELECT









1. Divide the *n* elements into groups of 5.





lesser 1. Divide the *n* elements into groups of 5. Find the median of each 5-element group by rote.



L6.22

greater





lesser

greater

- 1. Divide the *n* elements into groups of 5. Find the median of each 5-element group by rote.
- 2. Recursively SELECT the median *x* of the  $\lfloor n/5 \rfloor$  group medians to be the pivot.







At least half the group medians are  $\leq x$ , which is at least  $\lfloor n/5 \rfloor / 2 \rfloor = \lfloor n/10 \rfloor$  group medians.





#### Analysis (Assume all elements are distinct.)



At least half the group medians are  $\leq x$ , which is at least  $\lfloor n/5 \rfloor / 2 \rfloor = \lfloor n/10 \rfloor$  group medians. • Therefore, at least  $3 \lfloor n/10 \rfloor$  elements are  $\leq x$ .



L6.25

greater

lesser



#### Analysis (Assume all elements are distinct.)



At least half the group medians are  $\leq x$ , which is at least  $\lfloor n/5 \rfloor / 2 \rfloor = \lfloor n/10 \rfloor$  group medians.

- Therefore, at least  $3 \lfloor n/10 \rfloor$  elements are  $\leq x$ .
- Similarly, at least  $3\lfloor n/10 \rfloor$  elements are  $\geq x$ .

September 28, 2005 Copyright © 2001-5 by Erik D. Demaine and Charles E. Leiserson

L6.26

greater

lesser



## Minor simplification

- For  $n \ge 50$ , we have  $3\lfloor n/10 \rfloor \ge n/4$ .
- Therefore, for  $n \ge 50$  the recursive call to SELECT in Step 4 is executed recursively on  $\le 3n/4$  elements.
- Thus, the recurrence for running time can assume that Step 4 takes time *T*(3*n*/4) in the worst case.
- For n < 50, we know that the worst-case time is  $T(n) = \Theta(1)$ .



# **Developing the recurrence**

# T(n) SELECT(i, n)

 $\Theta(n) \left\{ \begin{array}{l} 1. \text{ Divide the } n \text{ elements into groups of 5. Find} \\ \text{the median of each 5-element group by rote.} \end{array} \right.$  $T(n/5) \begin{cases} 2. \text{ Recursively SELECT the median } x \text{ of the } \lfloor n/5 \rfloor \\ \text{group medians to be the pivot.} \end{cases}$  $\Theta(n) \qquad 3. \text{ Partition around the pivot } x. \text{ Let } k = \text{rank}(x). \end{cases}$ 

4. if i = k then return x else if i < k

*T*(3*n*/4) ≺

then recursively SELECT the *i*th smallest element in the low smallest element in the lower part else recursively SELECT the (i-k)th smallest element in the upper part



# Solving the recurrence

$$T(n) = T\left(\frac{1}{5}n\right) + T\left(\frac{3}{4}n\right) + \Theta(n)$$

**Substitution:**  $T(n) \le cn$ 

$$T(n) \leq \frac{1}{5}cn + \frac{3}{4}cn + \Theta(n)$$
  
=  $\frac{19}{20}cn + \Theta(n)$   
=  $cn - \left(\frac{1}{20}cn - \Theta(n)\right)$   
 $\leq cn$ ,

if *c* is chosen large enough to handle both the  $\Theta(n)$  and the initial conditions.



#### Conclusions

- Since the work at each level of recursion is a constant fraction (19/20) smaller, the work per level is a geometric series dominated by the linear work at the root.
- In practice, this algorithm runs slowly, because the constant in front of *n* is large.
- The randomized algorithm is far more practical.

**Exercise:** Why not divide into groups of 3?