## Lecture 10 **Computing medians and quicksort**

Chinmay Nirkhe | CSE 421 Spring 2025





# Previously in CSE 421...

## Multiplication

- Matrix multiplication
  - $O(n^{2.87})$  time algorithm for  $n \times n$  matrices
  - Strassen's divide and conquer algorithm
- Integer multiplication
  - $O(n^{1.58})$  time algorithm for multiplying *n*-bit numbers
  - Karatsuba's divide and conquer algorithm
- **Polynomial multiplication** 
  - $O(n \log n)$  time algorithm for multiplying degree *n* polynomials
  - Convert to evaluation basis via Fast Fourier transform for quick evaluation



#### Median

- Input: Input list  $X = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$  for n odd.
- Output: The median element i.e.  $y_{(n+1)/2}$  when Y = sort(X).

- An upper bound for the runtime is  $O(n \log n)$  from sorting + selecting.
- Can we do better? Could we achieve O(n)?

#### Median

- Consider a divide and conquer algorithm for median
- What would the recurrence relation have to be for T(n) = O(n)?
- Case 1: T(n) = 2T(n/2) + O(1)
  - Challenge is to split the problem X into two halves with O(1) compute
  - And to "stitch" the solutions to the two subproblems together in O(1) compute
- Case 2: T(n) = T(n/2) + O(n)
  - With O(n) time, we can make a constant number of passes through the list X
  - After constant number of passes, we need to find a sublist X' of size n/2 which must contain the median
  - Then we recurse on the sublist X'

#### Selection

- Let's define a more general problem called "Selection"
  - Input: list  $(X, k) \in \mathbb{R}^n \times [n]$ .
  - Output: The k-th element  $y_k$  when  $\vec{y} = \operatorname{sort}(\vec{x})$ .
- Generalizes the median problem

#### **Selection** Find the 6th element



## 03365223

#### **Selection** Find the 6th element



#### Selection Find the 6th element



#### Selection

- **Recursive algorithm** Selection(X, k):
  - Randomly sample *j* from [*n*]. Call  $x_j$  the "**pivot**".
  - Filter X into  $X_L$ ,  $X_E$ , and  $X_R$  based on if  $x_i < x_j$ ,  $x_i = x_j$ , or  $x_i > x_j$ .
  - If  $|X_I| \ge k$ , recursively output Selection $(X_I, k)$ .
  - Else if,  $|X_L| + |X_E| \ge k$ , output  $x_i$ .
  - Else, recursively output Selection( $X_R, k |X_L| |X_E|$ ).

## **Runtime analysis**

the input size from n to n/b for b > 1

• 
$$T(n) = T(n/b) + cn \implies T(n) = \frac{c}{1 - 1/c}$$

- However, each call may not reduce the size from n to n/b
- Depends on how close the randomly chosen  $x_i$  is to the middle
  - If pivot  $x_i$  was the largest element, then  $|X_L| = n 1$ ,  $|X_E| = 1$ , and  $|X_R| = 0$ .
  - Decreases instance size from *n* to n-1.
  - Fortunately, the probability this occurs is 1/n.

• In order to apply the master theorem, we would need to argue that each recursive call was reducing

$$-n$$

## **Runtime analysis**

- Amortized analysis:
  - If pivot  $x_j$  is the  $\ell$ -th element, then the next problem is of size  $\leq \max\{\ell, n \ell\}.$
  - With probability  $\geq 1/2$ , pivot  $x_j$  is the  $\ell$ -th element for  $\ell \in \{n/4, \dots, 3n/4\}$ .
  - The expected compute in reducing from *n*-sized instance to a 3n/4-sized instance is O(n).
- Total **expected** runtime: T(n) = T(3)



$$3n/4) + O(n) \implies T(n) = O(n).$$

## **Runtime analysis**

- Amortized analysis:
  - If pivot  $x_i$  is the  $\ell$ -th element, then the next problem is of size  $\leq \max\{\ell, n \ell\}$ .
  - With probability  $\geq 1/2$ , pivot  $x_i$  is the  $\ell$ -th element for  $\ell \in \{n/4, \dots, 3n/4\}$ .
  - The expected compute in reducing from *n*-sized instance to a 3n/4-sized instance is O(n).
    - $\geq 1/2$  probability, shrinks in 1 reduction.
    - $\geq 1/4$  probability, shrinks in 2 reductions.
    - ...  $\geq 1/2^{j}$  probability, shrinks in *j* reductions ...
    - Expected compute is  $\leq O(n) \cdot (\frac{1}{2} + \frac{1}{4} \cdot 2 + \frac{1}{8})$
- Total expected runtime:  $T(n) = T(3n/4) + O(n) \implies T(n) = O(n)$ .

$$\cdot 3 + \ldots) = O(n) \cdot 2$$

#### Derandomization

- The worst case runtime is  $O(n^2)$ .
  - Only happens with  $2^{-\Omega(n\log n)}$  probability.
- But, is there an algorithm that didn't require randomness?
- If we could guarantee that the pivot  $x_j$  was in the middle half, then each recursion would decrease in size by 3/4.
- Blum-Pratt-Floyd-Rivest-Tarjan (1973): Calculate a pivot in the middle 4n/10 in time O(n).

• Express the *n* elements as a  $5 \times (n/5)$  matrix of elements



- Express the *n* elements as a  $5 \times (n/5)$  matrix of elements
- Calculate the medians of each of the columns:  $Y = (y_1, y_2, ..., y_{n/5})$



the 5 clements in the

column

- Express the *n* elements as a  $5 \times (n/5)$  matrix of elements
- Calculate the medians of each of the columns:  $Y = (y_1, y_2, ..., y_{n/5})$



- Express the *n* elements as a  $5 \times (n/5)$  matrix of elements
- Calculate the medians of each of the columns:  $Y = (y_1, y_2, ..., y_{n/5})$
- Choose the pivot as the median of the medians:
   *p* ← median(*Y*)



#### **Pivot selection algorithm Runtime analysis**

- Express the *n* elements as a  $5 \times (n/5)$  matrix of elements
- Calculate the medians of O(1) per col.  $\in$ Total O(n). each of the columns:  $Y = (y_1, y_2, \dots, y_{n/5})$
- Choose the pivot as the median of the medians:  $p \leftarrow \operatorname{median}(Y) \quad T(n/s) \quad recursively$

Total time:  $T(n) = T(\frac{n}{2}) + O(n) \implies T(n) = O(n)$ 





• There are  $\geq n/10$  columns such that  $y_i \geq p$ .



- There are  $\geq n/10$  columns such that  $y_i \geq p$ .
- In each such column, there are 3 elements  $\geq y_i$ .



- There are  $\geq n/10$  columns such that  $y_j \geq p$ .
- In each such column, there are 3 elements  $\geq y_i$ .
- Therefore, there are  $\geq 3n/10$ elements  $\geq p$ .
- Similarly, there are  $\geq 3n/10$ elements  $\leq p$ . • So, p is in



## Median/Selection algorithm

- Input:  $(X, k) \in \mathbb{R}^n \times [n]$
- **Output:** the k-th item in the list X
- Algorithm:
  - Calculate  $p \leftarrow \text{median-of-medians}(X)$  in a
  - Filter X into  $X_L$ ,  $X_E$ , and  $X_R$  based on p
  - If  $|X_L| \ge k$ , recurse Selection $(X_L, k)$ 
    - Else if  $|X_L| + |X_E| \ge k$ , return p
    - Else, return Selection( $X_R, k |X_L| |$

$$Total: T(n) = T(\frac{7}{10}n) + T(\frac{n}{5}) + O(n)$$
  

$$\Rightarrow T(n) = O(n)$$
  
Pset problem on how to analyze this  
generalization of Master theorem  

$$na 5 \times (n/5) \text{ division.}$$
  

$$recursive T(\frac{n}{5}) + O(n)$$
  

$$\int recursive T(\frac{7}{10}n)$$
  

$$-|X_E|).$$

## Quicksort algorithm

- Sorting algorithm Quicksort(X):
  - Pick a pivot p (either randomized or with median-of-medians)
  - Filter X into  $X_L$ ,  $X_E$ ,  $X_R$  by comparing elements with p
  - Concatenate  $Sort(X_I), X_E, Sort(X_I)$ Computing expected runt variable

• The algorithm we just analyzed, "Quickselect", can be generalized to sorting

$$X_R$$
).

#### **Quicksort algorithm Runtime analysis**

- Runtime depends on pivot selection
- Median-of-means:
  - $T(n) \le T(\alpha n) + T(n \alpha n) + O(n)$  for  $\alpha \in [0.3, 0.7]$
  - T(n) = O(n) by analysis you will solve on your pset
- Choose random element:
  - Worst case:  $O(n^2)$  time
  - Amortized:  $O(n \log n)$  (next!)

#### Quicksort algorithm Runtime analysis for random choice of pivot

- Observations:
  - The runtime of Quicksort is proportional to the number of comparisons
  - The algorithm only compares two elements if is the pivot
- Let  $Y = (y_1, \dots, y_n)$  be the sorted version of the input.

• Let 
$$p_{ij} = \mathbf{Pr} \left[ y_i \text{ and } y_j \text{ are compared} \right]$$

• Claim: 
$$p_{ij} \leq \frac{2}{j-i+1}$$
 when  $i < j$ .

Expected number of comparisons:  

$$\sum_{i < j} P_{ij} \leq 2 \sum_{i < l} \sum_{j = i+l}^{n} \frac{1}{j - i + l}$$

$$= 2 \sum_{i = l}^{n} \sum_{i = l} \sum_{k = l}^{n-i+l} \frac{1}{k+l}$$
one  

$$= 2 \sum_{i = l}^{n} \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n-i+l}$$

$$\leq \log(n-i+l) + 1$$

$$\leq \log(n) + 1$$

$$\leq 2 n \log n + 2n$$
Runtime of quicksort = O(n log n).



### **Proof of claim**

• Claim: 
$$p_{ij} \leq \frac{2}{j-i+1}$$
 when  $i < j$ .

• **Proof**:

- $y_i \leq y_j$  and  $y_i$  and  $y_j$  are compared at most once
  - Comparisons only occur when one of them is the pivot
  - Case 1:  $y_i, y_i \in X_E$  and we never recurse on  $X_E$
  - Case 2:  $y_i \in X_E, y_i \in X_R$  and we never compare between  $X_L, X_E$ , and  $X_R$
  - Case 3:  $y_i \in X_L, y_j \in X_E$  and we never compare between  $X_L, X_E$ , and  $X_R$

- If and when  $y_i$  and  $y_j$  are compared during sort(X') then  $y_i, y_{i+1}, y_{i+2}, \dots, y_j \in X'$ 
  - Can be formally proven via induction  $\bullet$
  - So  $|X'| \ge j i + 1$ .
  - Probability that either  $y_i$  or  $y_j$  is chosen as pivot is ullet $\leq \frac{2}{i-i+1}.$



# Sorting in the real world

#### Quicksort

- Fast almost always, especially for in-memory sorting.
- Works well with caches due to good locality of reference.
- In practice,
  - Don't filter  $X_L, X_E$ , and  $X_R$ . Use in-place swaps.
  - When *n* is small, insertion sorting is a better base case.
  - Pick pivot randomly for small n, median of 3 random values for medium n, and median-ofmedians on 9 elements for large n
  - Never actually run the median-of-medians pivot finding routine

## Sorting in the real world

- Mergesort
  - Used when data is expressed as a linked list and RAM access to entries in the middle of the list is non-existent
  - Sorting over a dataset that cannot be stored in memory
  - Uses O(n) extra space when sorting arrays over Quicksort

## Sorting in the real world

#### Insertion sort

- Best when data is almost sorted already
- $O(n^2)$  when far from sorted
- Heap sort memory efficient choice
- Bucket sort distribution aware sorting
- Etc...