# $\left\{ \text{ Csc 80030 } \mid \text{ Lecture 4} \right\}$

## PROBABILISTIC ANALYSIS & RANDOMIZED ALGORITHMS

Hoeteck Wee · hoeteck@cs.qc.edu

## part 0

- ► HOMEWORK: HW2 is out; review HW1
- ► TODAY: randomized median finding & balls n' bins

PART 1 randomized median finding

## Median Finding

#### MEDIAN FINDING Problem

Input: a set *S* of *n* values from some totally ordered universe Goal: output the median element *m* of *S* 

- ► WHAT'S KNOWN: "easier" than sorting there is a deterministic linear-time algorithm.
- ► TODAY: a randomized linear-time algorithm based on random sampling and the fact that we can sort a set of size  $O(n^{3/4})$  in O(n) time.

#### RANDMEDIAN Algorithm

Input: A list S of n distinct values

- 1. Sample  $\ell$  from *S* such that rank<sub>*S*</sub>( $\ell$ )/ $n \in [\frac{1}{2} 2\delta, \frac{1}{2}]$ ;
- 2. Sample *u* from *S* such that rank<sub>*S*</sub>(*u*)/ $n \in [\frac{1}{2}, \frac{1}{2} + 2\delta]$ ;
- 3. By comparing with each value in *S*, compute  $C = \{y \in S \mid \ell \le y \le u\};$
- 4. Output the  $(\frac{1}{2}n \operatorname{rank}_{S}(\ell) + 1)$ 'th smallest element in the sorted set *C*.
- CORRECTNESS:  $\operatorname{rank}_{S}(\operatorname{output}) = \operatorname{rank}_{C}(\operatorname{output}) + (\operatorname{rank}_{S}(\ell) 1).$
- FACT:  $|C| \le 4\delta n + 1$ , can sort in  $O(\delta n \log n)$  time; set  $\delta = 1/n^{1/4}$ .
- GOAL: implement steps 1, 2 in O(n) time.

#### RANDMEDIAN Subroutine

- 1. Sample  $\ell$  from *S* such that rank<sub>*S*</sub>( $\ell$ )/ $n \in [\frac{1}{2} 2\delta, \frac{1}{2}];$ 
  - 1.1 Pick a set R of  $n^{3/4}$  elements in S with replacement.
  - 1.2 Output the element x in R whose rank is  $(\frac{1}{2} \delta)|R|$ .
- INTUITION:  $\operatorname{rank}_{S}(x)/n \approx \operatorname{rank}_{R}(x)/|R| = \frac{1}{2} \delta$
- CLAIM:  $\Pr_R[\operatorname{rank}_S(x)/n \notin [\frac{1}{2} 2\delta, \frac{1}{2}]$  is tiny.
- SUB-CLAIM 1:  $\Pr_R[\operatorname{rank}_S(x)/n > \frac{1}{2}]$  is tiny.
  - ► rank<sub>S</sub>(x)/n >  $\frac{1}{2} \Rightarrow #$ {elements in R > median(S)}  $\geq (\frac{1}{2} + \delta)|R|$
  - X: #{elements in R > median(S)}, so  $X \sim B(|R|, \frac{1}{2})$
  - Show  $\Pr[X \ge (\frac{1}{2} + \delta)|R|]$  is tiny.

### Analysis of RANDMEDIAN

- SUB-CLAIM 1:  $\Pr_R[\operatorname{rank}_S(x)/n > \frac{1}{2}]$  is tiny.
  - $X \sim B(|R|, \frac{1}{2})$ , will show  $\Pr[X \ge (\frac{1}{2} + \delta)|R|]$  is tiny.

PROOF:

- $E[X] = \frac{1}{2}|R|$  and  $Var[X] = \frac{1}{4}|R|$ .
- Applying Chebyshev's,

$$\begin{split} \Pr\bigl[X \geq (\frac{1}{2} + \delta) |R|\bigr] &\leq & \Pr\bigl[|X - \mathbf{E}[X]| \geq \delta |R|\bigr] \\ &\leq & \frac{\operatorname{Var}[X]}{(\delta |R|)^2} = \frac{1}{4n^{1/4}} \end{split}$$

### Analysis of RANDMEDIAN

- SUB-CLAIM 2:  $\Pr_R[\operatorname{rank}_S(x)/n < \frac{1}{2} 2\delta]$  is tiny.
  - ► rank<sub>S</sub>(x)/n <  $\frac{1}{2} 2\delta \Rightarrow$ #{elements in *R* with rank<sub>S</sub> <  $(\frac{1}{2} - 2\delta)n$ } ≥  $(\frac{1}{2} - \delta)|R|$
  - ► *Y*: #{elements in *R* with rank<sub>*S*</sub> <  $(\frac{1}{2} 2\delta)n$ }, so  $Y \sim B(|R|, \frac{1}{2} 2\delta)$
  - Show  $\Pr[Y \ge (\frac{1}{2} \delta)|R|]$  is tiny.

▶ PROOF:

- ►  $E[Y] = (\frac{1}{2} 2\delta)|R|$  and  $Var[Y] = (\frac{1}{4} 4\delta^2)|R| \le \frac{1}{4}|R|$ .
- Applying Chebyshev's,

$$\begin{aligned} \Pr[Y \ge (\frac{1}{2} - \delta)|R|] &\leq & \Pr[|Y - \mathrm{E}[Y]| \ge \delta|R|] \\ &\leq & \frac{\mathrm{Var}[Y]}{(\delta|R|)^2} \le \frac{1}{4n^{1/4}} \end{aligned}$$

## Analysis of RANDMEDIAN

#### **RANDMEDIAN** Algorithm

Input: A list S of n distinct values

- 1. Sample  $\ell$  from *S* such that rank<sub>*S*</sub>( $\ell$ )/ $n \in [\frac{1}{2} 2\delta, \frac{1}{2}]$ ;
- 2. Sample *u* from *S* such that rank<sub>*S*</sub>(*u*)/ $n \in [\frac{1}{2}, \frac{1}{2} + 2\delta]$ ;

3. ...

- $\Pr\left[\operatorname{rank}_{S}(\ell)/n \notin \left[\frac{1}{2} 2\delta, \frac{1}{2}\right]\right] \leq \frac{1}{2n^{1/4}}$
- similarly,  $\Pr[\operatorname{rank}_{S}(u)/n \notin [\frac{1}{2}, \frac{1}{2} + 2\delta]] \leq \frac{1}{2n^{1/4}}$
- Pr[RANDMEDIAN outputs incorrect answer]  $\leq \frac{1}{n^{1/4}}$

PART 2 | balls 'n bins

QUESTION. What is the probability that amongst 30 people in a room, two share the same birthday?

MODEL. Everyone's birthday is independently and uniformly chosen at random amongst 365 days.

ANALYSIS. Pr[all birthdays are distinct] is  $(1 - \frac{1}{365}) \cdot (1 - \frac{2}{365}) \cdot (1 - \frac{3}{365}) \cdots (1 - \frac{29}{365}) \approx 0.2937$ MORE GENERALLY... For *m* people and *n* "birthdays", it's

$$(1 - \frac{1}{n}) \cdot (1 - \frac{2}{n}) \cdot (1 - \frac{3}{n}) \cdots (1 - \frac{m-1}{n})$$
  
$$\approx \prod_{j=1}^{m-1} e^{-j/n} = e^{-m(m-1)/2n} \approx e^{-m^2/2n}$$

 $\Rightarrow$  constant prob of "collision" whenever  $m \gtrsim \sqrt{2n \ln 2}$ 

- ▶ *m* balls thrown into *n* bins
  - location of each ball independent and random
- ► Example: job scheduling
  - balls = tasks, bins = processors
- Quantities of interest
  - average load = expected number of balls in each bin
  - maximum load = number of balls in fullest bin
  - number of empty bins (= number of idle processors)

### Average/Maximum Load

•  $L_i$  be r.v. for # balls in Bin i

• 
$$L_i \sim B(m, \frac{1}{n})$$
, so  $\mathbb{E}[L_i] = \frac{m}{n}$ ,  $\operatorname{Var}[L_i] = \frac{m}{n}(1 - \frac{1}{n})$ 

#### Chernoff Bound

Let  $X_1, \ldots, X_n$  be *independent*  $\{0, 1\}$ -r.v.'s. Let  $X = X_1 + \cdots + X_n$  and  $\mu = \mathbb{E}[X]$ . Then, for all  $\delta > 0$ ,  $\Pr[X \ge (1 + \delta)\mu] \le e^{-\frac{\mu\delta^2}{2+\delta}}$ 

• APPLICATION. bounding  $Pr[L_i \ge 2 \ln n + 1]$  for m = n

• set 
$$\mu = 1, \delta = 2 \ln n$$
, so  $\frac{\mu \delta^2}{2+\delta} \ge 2 \ln n$   
 $\Rightarrow \Pr[L_i \ge 2 \ln n+1] \le e^{-2 \ln n} = \frac{1}{n^2}$ 

- By union bound,  $\Pr[\bigvee_{i=1}^{n} (L_i \ge 2 \ln n + 1)] \le \frac{1}{n}$
- Hence,  $\Pr[\text{maximum load} \le 2 \ln n + 1] \ge 1 \frac{1}{n}$ .
- e.g. n = 1 million, max load is at most 30 w.h.p.

▶ BETTER ANALYSIS.

$$\begin{aligned} \Pr[L_i \ge k] &= \Pr[\exists \text{ subset of } k \text{ balls all of which fall into bin } i] \\ &\leq \binom{n}{k} \cdot (1/n)^k \\ &\leq (ne/k)^k \cdot (1/n)^k = (e/k)^k \\ &\leq 1/n^2 \quad \text{ for } k \ge \frac{3\ln n}{\ln \ln n} \end{aligned}$$

BETTER BOUND.

- ▶ obtain a bound of  $O(\frac{\log n}{\log \log n})$  instead of  $O(\log n)$  for the maximum load.
- e.g. n = 1 million, max load is at most 16 w.h.p.

## **Empty Bins**

- Let *X* be random variable for # empty bins.
- Let  $X_i$  be r.v. indicating whether Bin *i* is empty.
- $\Pr[X_i = 1] = (1 \frac{1}{n})^m$  and  $\mathbb{E}[X] = n(1 \frac{1}{n})^m$ .
- ► NOTE.  $X_i$  and  $X_j$  are *not* independent, e.g.  $\Pr[X_i = 1 \land X_j = 1] = (1 - \frac{2}{n})^m \neq \Pr[X_i = 1] \cdot \Pr[X_j = 1]$

#### Empty Bins: Variance

• Recall  $\operatorname{Var}[X] = \operatorname{E}[X^2] - \operatorname{E}[X]^2$ 

- ►  $E[X^2] = E[(X_1 + \dots + X_n)^2] = \sum_{i=1}^n E[X_i^2] + \sum_{i \neq j} E[X_i X_j]$
- If  $X_i \in \{0, 1\}$ , then  $E[X_i^2] = E[X_i]$
- Computing  $E[X_iX_j]$

•  $E[X_iX_j] = Pr[X_iX_j = 1] = Pr[X_i = 1 \land X_j = 1] = (1 - \frac{2}{n})^m$ 

► Computing Var[X]

► 
$$E[X^2] = n(1 - \frac{1}{n})^m + n(n-1)(1 - \frac{2}{n})^m$$
  
►  $Var[X] = n(1 - \frac{1}{n})^m + n(n-1)(1 - \frac{2}{n})^m - n^2(1 - \frac{1}{n})^{2m}$ 

THE END | HW2 due next week