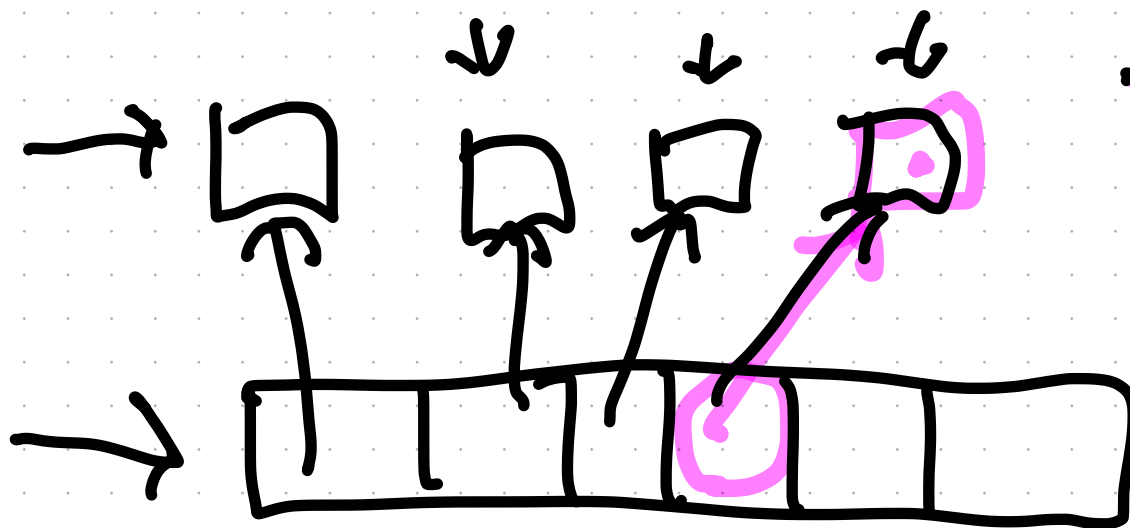


Big Data Algorithms

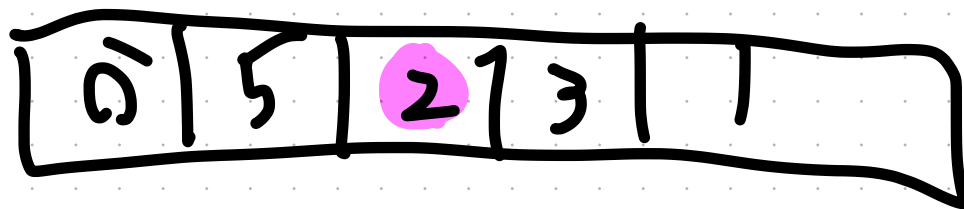
Week 3

start at 8:05

— Better analysis of hashing



— Sublinear
algorithms
/ streaming



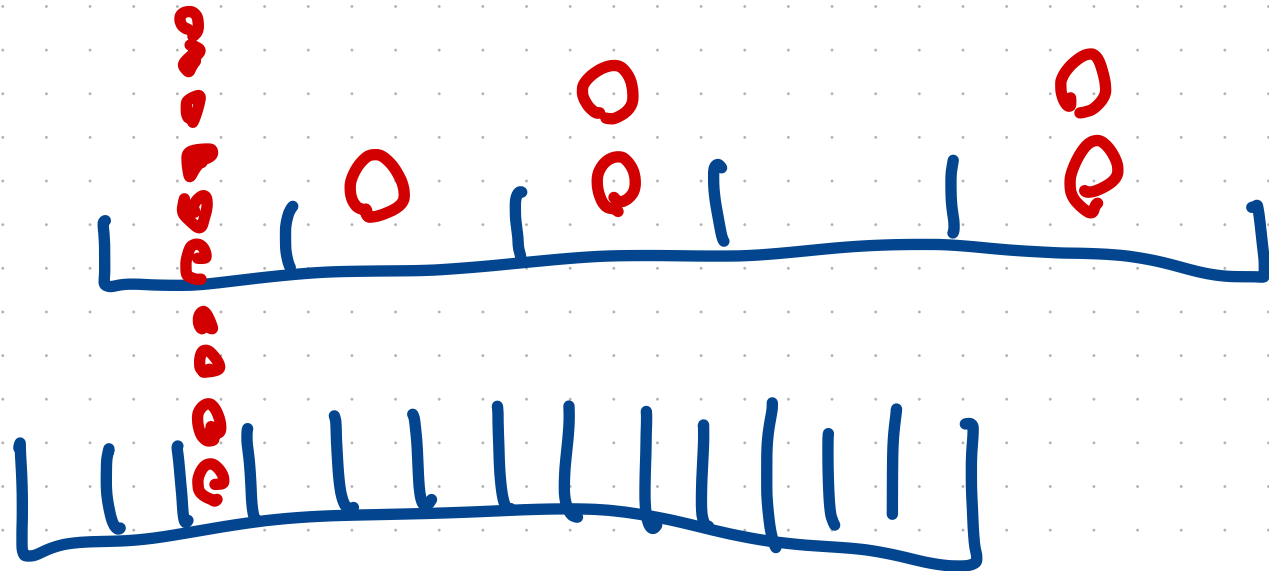
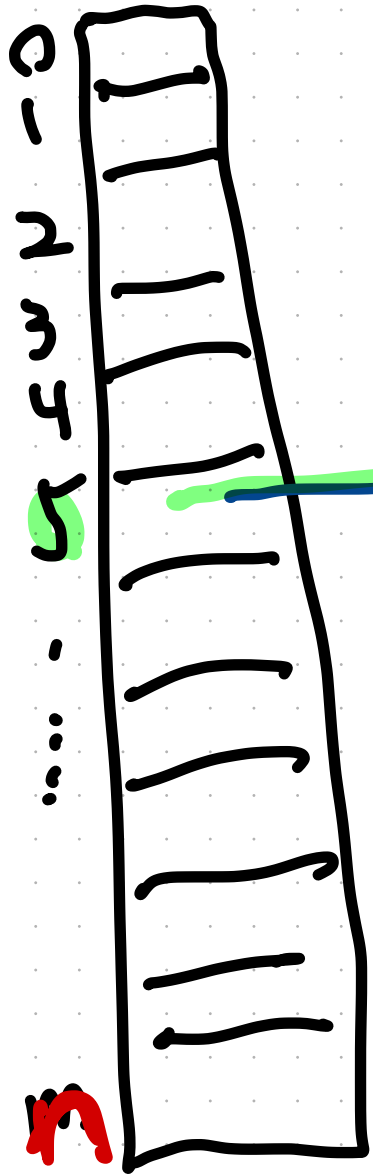
$$\text{hash}(x) \% m$$

$$\text{hash}(\underline{10}) \% m = 5$$

Insert n items

"Balls in Bins"

Given n balls, throw them
randomly into n bins



N Balls N bins

① What is the average size bin?

② What is the expected size of a bin?

- ③ Can you say: Bin i has size at most x with prob $\leq \rho$

- ④ Can you say: All bins have size at most x with prob $\leq \rho$

Average

b_i = size of bin i

$$\text{Average Bin size} = \frac{1}{N} \underbrace{\sum_{i=1}^N b_i}_N = \frac{1}{N} \cdot N = 1$$

Expected Value

$$E[X] = \sum_i i \Pr[X=i]$$

↑
random variable

$$E[b_i] = \sum_{j=0}^N j \Pr[b_i=j]$$

what is this

$$\Pr[b_i=0] = \left(1 - \frac{1}{N}\right)^N \approx \frac{1}{e}$$

↑
fact from last time

$$\Pr[b_i=j] = \binom{N}{j} \left(1 - \frac{1}{N}\right)^{N-j} \left(\frac{1}{N}\right)^j$$

Linearity of Expectation

$$E[X+Y] = E[X] + E[Y]$$

$$E[X \cdot Y] \text{ Not true that always } = E[X \cdot Y]$$

$$\begin{aligned} E[b_i] &= E\left[\sum_{j=1}^N b_{i,j}\right] \\ &= \sum_{j=1}^N E[b_{i,j}] \\ &= \sum_{j=1}^N \frac{1}{N} = N \cdot \frac{1}{N} = 1 \end{aligned}$$

$$b_{i,j} = \begin{cases} 1 & \text{if } j^{\text{th}} \text{ ball is in } b_i \\ 0 & \text{otherwise} \end{cases}$$

"Indicator random variable"

$$\begin{aligned} E[b_{i,j}] &= \sum_k k \Pr[b_{i,j} = k] \\ &= \Pr[b_{i,j} = 1] \\ &= \frac{1}{N} \end{aligned}$$

Markov's Inequality

100 students

Average score is 4

Claim: At most 50 students have
a grade of 8 or higher

Assumption: Scores are not negative


$$\Pr[X > a] \leq \frac{E[X]}{a}$$

wanted: bin i has size $\geq x$ with pr $\leq \frac{1}{x}$

$$\Pr[b_i > x] \leq \frac{E[b_i]}{x} = \frac{1}{x}$$

bin i has size ≥ 100 with probability $\leq \frac{1}{100}$
 $\leq 1\%$ 99%

Chernoff Bounds

They bound $\sum_{i=1}^K X_i$ where the X_i are independent

 rand var

Defn. X, Y are ind if $\Pr[X=i \text{ and } Y=j] = \Pr[X=i] \cdot \Pr[Y=j]$

$b_{i,\bar{j}}$ = Ball \bar{j} is in
bin i

$b_{i,\bar{j}}$ and $b_{i,\bar{j}'}$ $\bar{j}' \neq \bar{j}$ Ind

$b_{i,\bar{j}}$ and $b_{i',\bar{j}}$ $i' \neq i$

$\left(\frac{1}{\omega}\right)^n b_i$ b_j $i \neq j$

Break until

9:10

Chernoff Bounds

They bound $\sum_{i=1}^K X_i$ where the X_i are independent

$$\Pr[X \leq (1-\delta)E[X]] \leq \frac{1}{e^{\delta^2 E[X]/2}} \quad 0 \leq \delta \leq 1$$

$$\Pr[X \geq (1+\delta)E[X]] \leq \frac{1}{e^{\delta^2 E[X]/(2+\delta)}} \quad 0 \leq \delta$$

$$\Pr[b_i \geq (1+\delta) \underbrace{E[b_i]}_1] \leq e^{-\delta^2 \underbrace{E[b_i]}_1 / (2+\delta)}$$

$$\Pr[b_i \geq (1+\delta)] \leq e^{-\delta^2 / (2+\delta)}$$

What is $\Pr[b_i \geq 100] \leq e^{-99^2/(2+99)}$
 $\delta = 99$
 $= e^{-\frac{9801}{11}}$

1. Գործարարական գործունեության
 զարգացումը
 2021

With High Probability

Event e happens "with polynomially high probability" means it happens with prob $1 - O(\frac{1}{n^c})$ for some c

Intuition: Chernoff $e^{-\delta}$
Set $\delta = c \ln N \rightarrow e^{-c \ln N} = N^{-c} = \frac{1}{N^c}$

What is the chance $b_i \geq 1 + d \ln n$?

$$\Pr[b_i \geq (1+\delta)] \leq e^{-\delta^2/(2+\delta)}$$

$$\Pr[b_i \geq 1 + d \ln n] \leq e^{-(d \ln n)/(2 + d \ln n)}$$

$$\approx \frac{1}{n^{d-1}}$$

$$\Pr[b_i \geq 1 + 3 \ln n] \leq \frac{1}{n^2}$$

"bin i is $\leq 1 + 3 \ln n$ w.h.p."

"bin i is $O(\log n)$ w.h.p."

Union Bound

$$P[X=i \text{ and } Y=j] \leq P[X=i] + P[Y=j]$$

$$P\left[\begin{array}{c} 4 \text{ coins tossed} \\ \text{all heads} \end{array}\right] \leq 4 \cdot P[\text{one heads}]$$

$$= 4 \cdot \frac{1}{2}$$

$$Pr[b_i \geq 1 + 3 \ln n] \leq \frac{1}{n^2} = 2$$

$$Pr\left[\bigwedge_{i=1}^n \text{All in } 1..N \quad b_i \geq 1 + 3 \ln N\right] \leq n \cdot \frac{1}{N^2} = \frac{1}{N}$$

Approximate Median Finding

10000

~~X~~ 3, 99, 32, 78, 2, 4, 6, 8

Average: 25.8 \rightarrow 1136

Median: 6 \rightarrow 8

ϵ - Approx Median

x is an ϵ -approx median
if it is between the

$\left(\frac{1-\epsilon}{2}\right)N^{\text{th}}$ and $\left(\frac{1+\epsilon}{2}\right)N^{\text{th}}$ largest item

Stupid alg: Pick a random item x

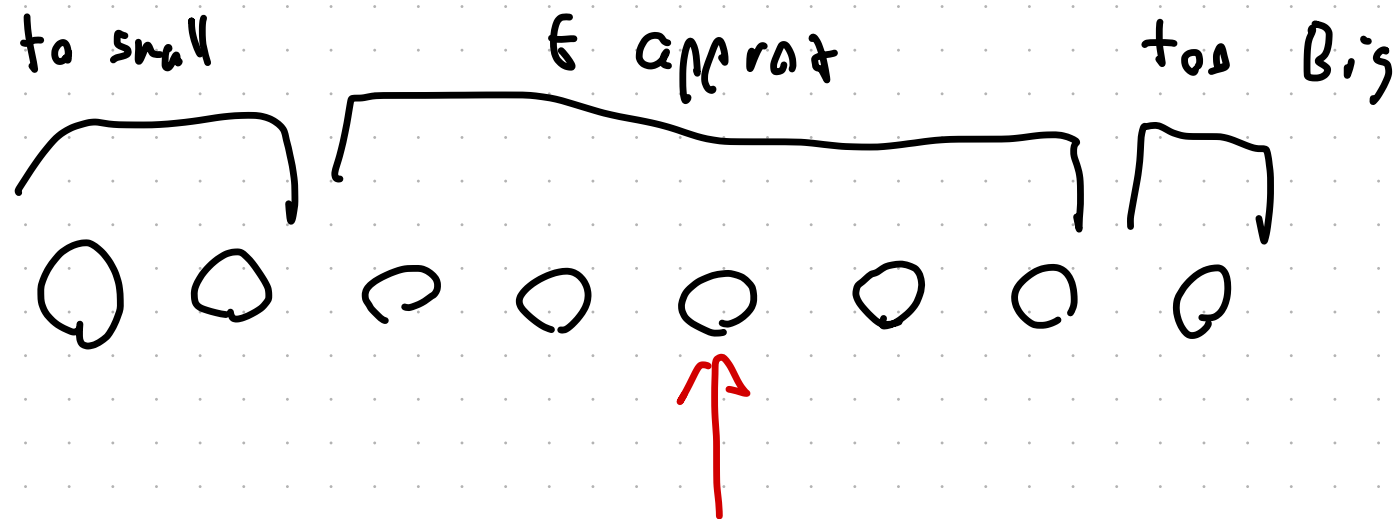
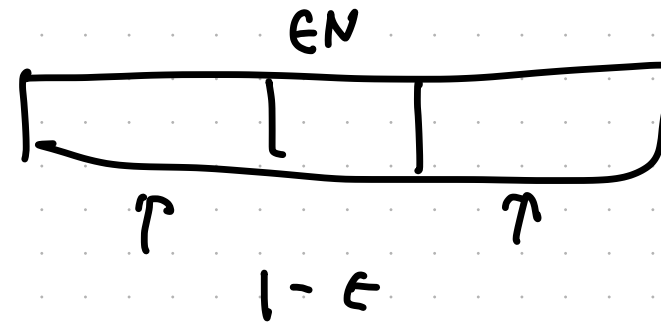
$\rightarrow x$ is a 50%-median with prob 50%

Better algorithm: Pick k random elements
return the median of these elements

X_i = i th sample is not a ϵ -approx median
Bad sample

$$E[X_i] = 1 - \epsilon$$

X = Total number of
Bad samples



$\kappa=2$

$\kappa=1$

