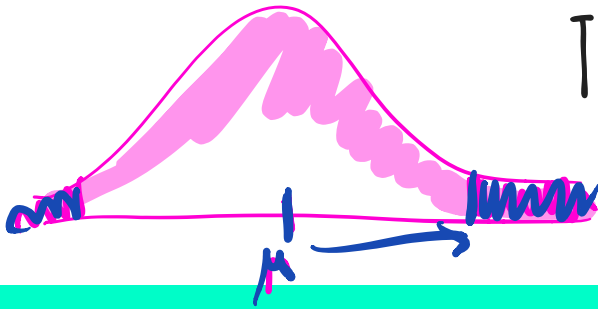


# HEAVY HITTERS TAIL BOUNDS

(continued)



ANNA KARLIN

# PROBLEM

- Input: sequence of  $n$  elements  $x_1, x_2, \dots, x_n$  from a known universe  $U$  (e.g., 8-byte integers).
- Goal: perform a computation on the input, in a single left to right pass where
  - Elements processed in real time
  - Can't store the full data. => minimal storage requirement to maintain working "summary"

# HEAVY HITTERS: KEYS THAT OCCUR MANY TIMES

$x_1$  32,  $x_2$  12,  $x_3$  14, 32, 7,  $x_6$  12, 32, 7, 32, 12,  $x_{11}$  4,

$$f_{32}^6 = 2 \quad f_{32}^7 = 3$$

Applications:

- Determining popular products
- Computing frequent search queries
- Identifying heavy TCP

$x_1, x_2, \dots, x_n$

$$k=100$$

$\forall t \leq n$   $f_x^t$ : # times element  $x$  has appeared in  $x_1, x_2, \dots, x_t$

$$f_x^n \geq 0.01 n$$

$\leq 100$  such elements

Goal: Find/output all elements  $x$  with  $f_x^n \geq \frac{n}{k}$   
 (These elts are "heavy hitters")



Output has size  $O(k)$

much smaller than  $n$   
space impossible

Provably impossible to solve this  
problem exactly with sublinear space

Modified goal: (solve  $\epsilon$ -HH problem)

① If  $f_x^n \geq \frac{n}{k}$ ,  $x$  added to HH list

② If some elt, say  $y$ , added  
to list, then w.p.  $\geq 1 - \delta$

$$f_y^n \geq \frac{n}{k} - \epsilon n$$

Example:

$$k=20 \quad \frac{n}{k} = \frac{n}{20} = 0.05n$$

$$\epsilon = 0.01$$

$$f_y^n \geq 0.05n - 0.01n = 0.04n$$

outputting a list  
of elts we claim  
are HHs.

if an elt  $y$  in list  
w. prob  $\geq \underline{\underline{1-\delta}}$

$$\underline{\underline{f_y^n \geq \frac{n}{k} - \epsilon n}}$$



# COUNT-MIN SKETCH

- Maintain a short summary of the information that still enables answering queries.
- Cousin of the Bloom filter
  - Bloom Filter solves the “membership problem”.
  - We want to extend it to solve a counting problem.

# COUNT-MIN SKETCH

Modified goal: (solve  $\epsilon$ -HH problem)

① If  $f_x^n \geq \frac{n}{k}$ ,  $x$  added to HH list

② If some elt, say  $y$ , added to list, then w.p.  $\geq 1 - \delta$   
 $f_y^n \geq \frac{n}{k} - \epsilon n$

designer specifies  $k, \epsilon, \delta$

$\Rightarrow$

$b, \ell$

Keep 2D array

$\ell$  hash tables, each of size  $b$ .

initialize tables with all 0s

when elt x shows up.

Update (x):  $\forall 1 \leq j \leq l$  increment  $t_j[h_j(x)]$

Count (x): return  $\min_{1 \leq j \leq l} t_j[h_j(x)]$

if  $\text{Count}(x) \geq \frac{n}{k}$ , add x to H+H list.

$$t_j[h_j(x)] \geq f_x^+$$

$$\text{Count}(x) \geq f_x^t$$

	1	2	3	4	5							6
$h_1$	1		4									
$h_2$		1		1	3							
$h_l$												

Example  $l=2$

$t_1$  first row  
 $t_2$  second row

Count(x)  
return 3

$$\min_{1 \leq j \leq 2} t_j[h_j(x)] = 3$$

Suppose

$$h_1(x)=3 \quad h_2(x)=5$$

$$h_1(y)=3 \quad h_2(y)=2$$

$$h_1(z)=1 \quad h_2(z)=4$$

initialize tables with all 0s

when elt  $x$  shows up.

Update ( $x$ ):  $\forall 1 \leq j \leq L$  increment  $t_j[h_j(x)]$

Count ( $x$ ): return  $\min_{1 \leq j \leq L} t_j[h_j(x)]$

if  $\text{Count}(x) \geq \frac{n}{k}$ , add  $x$  to HH list.

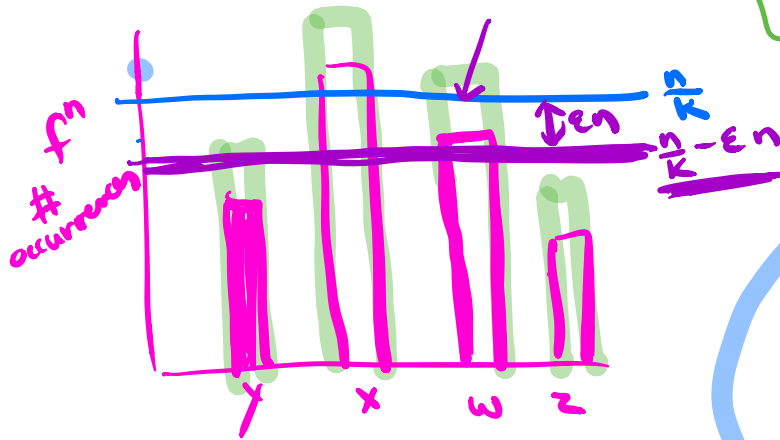
Observations

$$\forall j, \forall x, \forall t$$

$$t_j[h_j(x)] \geq \underline{f_x^t}$$

$\Rightarrow$

$$\text{Count}(x) = \min_{1 \leq j \leq L} t_j[h_j(x)] \geq \underline{f_x^t}$$



$x_1, \dots, x_n$

green show current Count

$\forall y$  that is output  
w.p.  $\geq 1-\delta$   $\underline{f_y^n} \geq \frac{n}{k} - \epsilon$

$$\frac{n}{k} = 0.05n$$

$$\epsilon n = 0.01n$$

$$0.04n$$

b.l.

## Assumptions

① hash functions behave like random maps  
 $h_1, \dots, h_\ell: U \rightarrow \{0, 1, \dots, b-1\}$

$$\forall x \neq y \quad \Pr(h_j(x) = h_j(y)) = \frac{1}{b}$$

② hash fns  $h_1, \dots, h_\ell$   
are indep of each other.

initialize tables with all 0s

when elt  $x$  shows up.

Update ( $x$ ):  $\forall 1 \leq j \leq \ell$  increment  $t_j[h_j(x)]$

Count ( $x$ ): return  $\min_{1 \leq j \leq \ell} t_j[h_j(x)]$

if  $\text{Count}(x) \geq \frac{n}{k}$ , add  $x$  to HH list.

Fix time  $t_j$ ,  $x_1, x_2, \dots, x_t$  have just arrived

$$Z_j^t \triangleq t_j [h_j(x)]$$

$$Z_j^t \geq f_x^t$$

$$Z_j^t = f_x^t + \sum_{y \neq x} f_y^t \underline{w_{xy}}$$

$$\underline{E(Z_j^t)} = f_x^t + \sum_{\substack{\text{distinct } y \neq x \\ y \in \{x_1, \dots, x_t\}}} f_y^t E(w_{xy})$$

$$= f_x^t + \sum_{y \neq x} f_y^t \cdot \frac{1}{b}$$

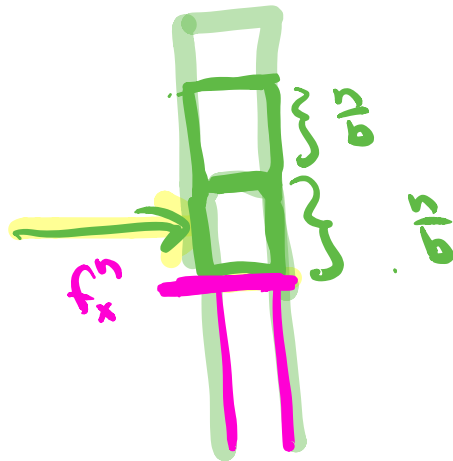
$$= f_x^t + \frac{1}{b} \left[ \sum_{y \neq x} f_y^t \right]$$

$$= \frac{t}{b} - f_x^t \leq t \leq n$$

$$\leq f_x^t + \frac{b}{b}$$

$$E(Z_j^t) - f_x^t \leq \frac{b}{b}$$

$$w_{xy} = \begin{cases} 1 & h_j(x) = h_j(y) \\ 0 & \text{otherwise} \end{cases}$$



$$\Pr(\underline{Z_j^t} - f_x^t \geq 2\frac{n}{b}) = ?$$

By Markov's Inequality,  
if  $X \geq 0$  then  
 $\Pr(X > 2E(X)) \leq \frac{1}{2}$

Let  $X \triangleq \underline{Z_j^t} - f_x^t$

$\Rightarrow X \geq 0$   $\underline{Z_j^t}$  is overestimate of  $f_x^t$

$$\Pr(\underline{Z_j^t} - f_x^t \geq 2E(X)) \leq \frac{1}{2}$$

$\forall j, \forall t \leq n$  (\*)  $\Pr(\underline{Z_j^t} - f_x^t \geq 2\frac{n}{b}) \leq \Pr(\underline{Z_j^t} - f_x^t \geq 2E(X)) \leq \frac{1}{2}$   
 $E(X) \leq \frac{n}{b}$

Putting it all together

$$\Pr(\text{Count}(x) - f_x^n \geq \frac{2n}{b})$$

$$= \Pr(\min(\underline{Z_1^n}, \dots, \underline{Z_e^n}) - f_x^n \geq \frac{2n}{b})$$

$$= \Pr(\underline{Z_1^n} - f_x^n \geq \frac{2n}{b}, \underline{Z_2^n} - f_x^n \geq \frac{2n}{b}, \dots, \underline{Z_e^n} - f_x^n \geq \frac{2n}{b})$$

$$= \prod_{j=1}^e \Pr(\underline{Z_j^n} - f_x^n \geq \frac{2n}{b}) \leq \underbrace{\frac{1}{2} \cdot \frac{1}{2} \cdots \frac{1}{2}}_e = \frac{1}{2^e}$$

↑ independence of hash fns for different tables

$$\Pr(\text{Count}(x) - f_x^n \geq \frac{2n}{b}) \leq \frac{1}{2^e}$$

$k, \epsilon, \delta$ .

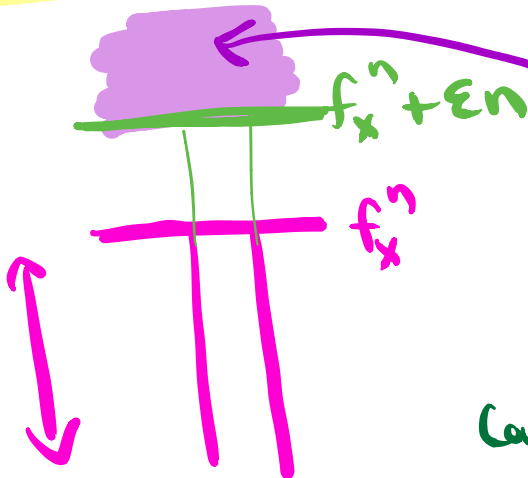
Modified goal: (solve  $\epsilon$ -HH problem)

① If  $f_x^n \geq \frac{n}{k}$ ,  $x$  added to HH list

② If some elt, say  $y$ , added to list, then w.p.  $\geq 1 - \delta$   $\leftarrow$   
 $\nearrow$   $f_y^n \geq \frac{n}{k} - \epsilon n$

$$\Pr(\text{Count}(x) - f_x^n \geq \frac{2n}{b}) \leq \frac{1}{2e}$$

$\epsilon n$        $\delta$



Prob Count(x) ends up in purple  $\leq \delta$

If  $y$  added to HH list,  $\Rightarrow$   
 $\left[ \text{Count}(y) \geq \frac{n}{k} \Rightarrow \begin{aligned} &\text{w.p.} \geq 1 - \delta \\ &f_y^n \geq \frac{n}{k} - \epsilon n \end{aligned} \right]$

$$\epsilon x = \frac{2n}{b}$$

$$\Rightarrow b = \frac{2n}{\epsilon} (*)$$

$$\delta = \frac{1}{2e}$$

$$\Rightarrow 2^L = \frac{1}{\delta}$$

$$\Rightarrow L = \log_2\left(\frac{1}{\delta}\right) (*)$$



# COUNT-MIN SKETCH

- Elegant small space data structure.  $O(\underline{b \cdot l} + k)$
- Space used is independent of  $n$ .
- Is implemented in several real systems.
  - AT&T used in network switches to analyze network traffic.
  - Google uses a version on top of Map Reduce parallel processing infrastructure and in log analysis.
- Huge literature on sketching and streaming algorithms (algorithms like Distinct Elements, Heavy Hitters and many many other very cool algorithms).

## Hash functions

Say hashing  $n$  32 bit integers into table of size  $b$ .  
Pick prime number  $p > \min(n, 2^{32})$

$$\mathcal{H} = \left\{ h_{e,g}(x) = [(ex + g) \bmod p] \bmod b \mid \begin{array}{l} 1 \leq e \leq p-1 \\ 0 \leq g \leq p-1 \end{array} \right\}$$

$x \in 0, \dots, b-1$

family of hash fns:  $(p-1) \cdot p$  # of fns in family.

If  $h$  is chosen uniformly at random from  $\mathcal{H}$

$$\forall x \neq y \quad \Pr(h(x) = h(y)) \leq \frac{2}{b}.$$

2 different hash fns

each selected uniformly at random  
and independently from  $\mathcal{H}$

		1	2	3	4	5			6
e.g. 1	1								
e.g. 2	2								
e.g. 2	2								

# 6.1 TAIL BOUNDS

MOST SLIDES BY JOSHUA FAN AND ALEX TSUN

# AGENDA

- MARKOV'S INEQUALITY
- ~~CHEBYSHEV'S INEQUALITY~~
- ~~THE LAW OF LARGE NUMBERS~~



# MARKOV'S INEQUALITY (INTUITION)

The score distribution of an exam is modelled by a rv  $X$  with range  $\Omega_X \subseteq [0, 110]$  (for extra credit).

If the average was  $E[X] = 50$ , at most what percentage of the class could have gotten 100 (or higher)?

at most 50% of class  
could have gotten score  
of 100 or higher.

Pf by  $\rightarrow \leftarrow$   
suppose

> 50% got score  $\geq 100$

$$E(X) = \sum_{\text{scores} \geq 100} \text{score} \Pr(X = \text{score}) > 50$$

$\underbrace{\hspace{10em}}_{> 50\%}$

$$\geq 100 \cdot \frac{1}{2} > 50$$

$\leq$

a) 100%

$\Rightarrow$  b) 50%

c) 25%

d) no bound



# MARKOV'S INEQUALITY (INTUITION)

The score distribution of an exam is modelled by a rv  $X$  with range  $\Omega_X \subseteq [0, 110]$  (for extra credit).

If the average was  $E[X] = 50$ , at most what percentage of the class could have gotten 100 (or higher)?

$$\frac{1}{2}$$

If the average was  $E[X] = 25$ , at most what percentage of the class could have gotten 100 (or higher)?

at most  $\frac{1}{4}$  can get score  $\geq 100$

- $\leq$
- a) 100%
  - b) 50%
  - c) 25%
  - d) no bound



# MARKOV'S INEQUALITY (INTUITION)

The score distribution of an exam is modelled by a rv  $X$  with range  $\Omega_X \subseteq [0, 110]$  (for extra credit).

$\Omega_X$  all #s  $\leq 110$

If the average was  $E[X] = 50$ , at most what percentage of the class could have gotten 100 (or higher)?

$\frac{1}{2}$

If the average was  $E[X] = 25$ , at most what percentage of the class could have gotten 100 (or higher)?

$\frac{1}{4}$

What if you could get a negative score?

Any r.v. with  $E[X] = 50$

$$0.01(-4900) + 0.99 \cdot 100 = 50$$

- $\leq$
- a) 100%
  - b) 50%
  - c) 25%

d) no bound

0.01	-4900	0.99	100
		0.99999	

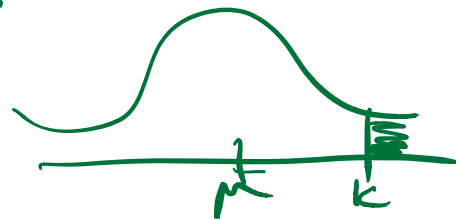


# MARKOV'S INEQUALITY

Markov's Inequality: Let  $X \geq 0$  be a **nonnegative** random variable (discrete or continuous), and let  $k > 0$ . Then,

bound on upper tail  
↓

$$P(X \geq k) \leq \frac{E[X]}{k}$$



$$E(X) = \sum_{\substack{x | x < k \\ x \in \Omega_X}} x \Pr(X=x) + \sum_{\substack{x | x \geq k \\ x \in \Omega_X}} x \Pr(X=x)$$

$\geq 0$

$$\geq \sum_{x | x \geq k} x \Pr(X=x) \geq k \sum_{x | x \geq k} \Pr(X=x) = k \Pr(X \geq k)$$

$$x \geq k$$

$$E(X) \geq k \Pr(X \geq k) \\ \equiv \Pr(X \geq k) \leq \frac{E(X)}{k}$$

## MARKOV'S INEQUALITY

**Markov's Inequality:** Let  $X \geq 0$  be a **nonnegative** random variable (discrete or continuous), and let  $k > 0$ . Then,

$$P(X \geq k) \leq \frac{E[X]}{k}$$

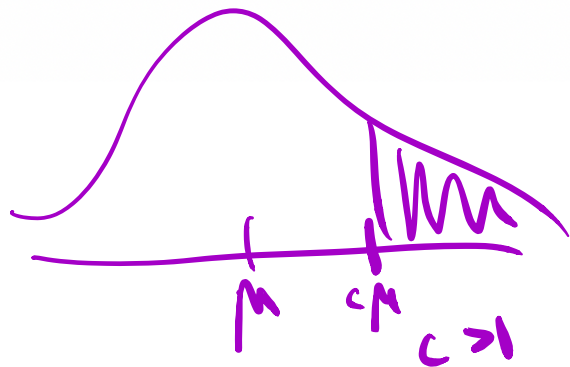
$$\Rightarrow \Pr(X \geq cE(X)) \leq \frac{E(X)}{cE(X)} = \frac{1}{c}$$

Alternatively,

Setting  $k = cE(X)$

$$P(X \geq cE[X]) \leq \frac{1}{c}$$

$$c > 1$$



$$c=2$$

$$\Pr(X \geq 2E(X)) \leq \frac{1}{2} \\ \text{for } X \geq 0$$

# MARKOV'S INEQUALITY (PROOF)



**Markov's Inequality:** Let  $X \geq 0$  be a **nonnegative** rv and let  $k > 0$ . Then,

$$P(X \geq k) \leq \frac{E[X]}{k}$$

**Proof (Markov):**

$$\begin{aligned} X \geq 0 & \quad E[X] = \int_0^{\infty} x f_X(x) dx = \int_0^k x f_X(x) dx + \int_k^{\infty} x f_X(x) dx \\ & \quad \geq \int_k^{\infty} x f_X(x) dx \geq \int_k^{\infty} k f_X(x) dx = k \int_k^{\infty} f_X(x) dx = k P(X \geq k) \end{aligned}$$

Rearranging gives

$$P(X \geq k) \leq \frac{E[X]}{k}$$

