

Lecture #5 . 1/26

Eli Goldin

Last time:

stream $x_1, \dots, x_m \in [N]$
 $n = \#IP = 2^{32}$

Goal: estimate # distinct.

Algo: Bottom-k

using $h: [N] \rightarrow [Q]$.

$k = \Theta(1/\epsilon^2)$ enough

for $(1 \pm \epsilon)$ approx $(d \gg 1/\epsilon^2)$.

Add remarks:

1) $\rightarrow [Q]$

$h: [N] \rightarrow \{0, 1, \dots, M-1\}$.

$\rightarrow \{0, 1/M, 2/M, \dots, \frac{M-1}{M}\}$.

M large enough.

$h(x)$ vs $h(y)$.

tiny issues $h(x)$ can collide
 $h(y)$

on different items.

$$M = n^3 \quad \mathbb{H}[\text{collisions}] \leq \frac{n^2}{M} \leq \frac{1}{n}$$

2) fully random?
possible with "limited randomness"

Def: \mathbb{H} is k -wise independent if
if $\forall a_1, a_2, \dots, a_k \in \{0, 1, \dots, \frac{M-1}{n}\}$
 $\forall x_1, \dots, x_k \in [n]$ distinct.

$$\Pr[h(x_i) = a_i \forall i=1..k] = \frac{1}{M^k}$$

$$\Pr[h(x) = a \wedge h(y) = a]$$

$$\mathbb{H} \text{ univ. if } \Pr[h(x) = h(y)] = \sum_{a \in \{0, \dots, \frac{M-1}{n}\}} \Pr[h(x) = h(y) = a]$$

$$= \sum_a \frac{1}{M^2} = \frac{M}{M^2} = \frac{1}{M}$$

Rem: if \mathbb{H} is 2-wise indep \Rightarrow universal.

Rem: $\forall k \geq 2$, can construct k -wise indep \mathbb{H} of size $\lg |\mathbb{H}| \leq O(k \cdot \lg n)$.

Thm [Thorup '13]: for Boolean- k also enough to use 2-wise indep h.f.

Problem: Stream $x_1, x_2, \dots, x_m \in [U]$.

Most frequent?

Def: freq. vector $f_x =$ how many times x appeared in str.

$$f \in \mathbb{R}_+^n$$

$$\|f\|_1 = \sum_x f_x = m = \text{len. stream.}$$

$\rightarrow x$ s.t. f_x is maximal.

Simple sol: store vector f explicitly
 \rightarrow space $O(n)$.

Thm: - can't do exactly.

- can't do 2-approx, using random.

Bad cases: all items appear
once
except 1' item app 2x.

Related Problem: - report items which are
sufficiently frequent

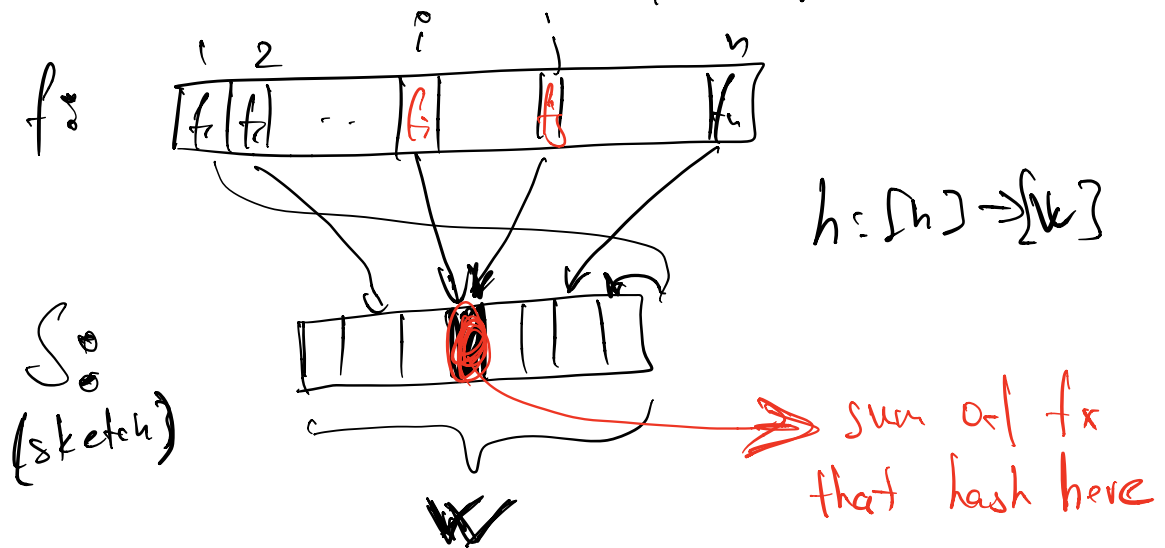
Def: $\phi \in (0, 1)$, item $x \in [n]$ is ϕ -heavy hitter

$$\text{if } f_x \geq \phi \cdot \sum_{y \in [n]} f_y = \phi \cdot m.$$

New goal: find all items that are ϕ -FH. OR space depends on ϕ .

Algorithm CountMin:

Basic idea: to get some estimate \hat{f}_x for f_x .



Bucket - Algor

-init: $S[1..k] = 0$

- when we see x in stream $\{f_x := f_{x+1}\}$
 $S[h(x)] := S[h(x)] + 1$.

→ @ end: to get $\hat{f}_x = S[h(x)]$.

Obs: $\hat{f}_x \geq f_x$.

ϕ -HH:
 $f_x \geq \phi \cdot m$.

Claim: $\Pr_h[\hat{f}_x > f_x + \epsilon \phi m] \leq 0.1$,
 as long as $w \geq \Omega\left(\frac{1}{\epsilon \phi}\right)$.

Proof: $\mathbb{E}_h[\hat{f}_x] = \mathbb{E}_h\left[\sum_{y \in [m]} \mathbb{1}[h(x)=h(y)] \cdot f_y\right]$

$$= \mathbb{E}\left[f_x + \sum_{y \neq x} \mathbb{1}[h(x)=h(y)] \cdot f_y\right]$$

$$= f_x + \sum_{y \neq x} f_y \cdot \underbrace{\Pr[h(x)=h(y)]}_{\frac{1}{w}}$$

$$\leq f_x + \frac{\sum_{y \neq x} f_y}{w} = f_x + \frac{m}{w}$$

$$\mathbb{E}[\underbrace{\hat{f}_x - f_x}_1] \leq \frac{m}{w}$$

↳ positive r.v.

⇒ By Markov's Ineq:

$$\Pr \left[\hat{f}_x - f_x \geq \underbrace{10 \cdot \frac{m}{\epsilon}}_w \right] \leq 0.1.$$

$$= \epsilon \cdot \phi \cdot m$$

$$\Rightarrow \text{set } w = 10/\epsilon\phi.$$



The ϕ -HH algorithm:

- for each $x \in [n]$, compute \hat{f}_x from S ,

- if $\hat{f}_x \geq \phi \cdot m \Rightarrow$ report as ϕ -HH.

Corollary: - if $f_x \geq \phi m \Rightarrow$ algo reports it as HH.

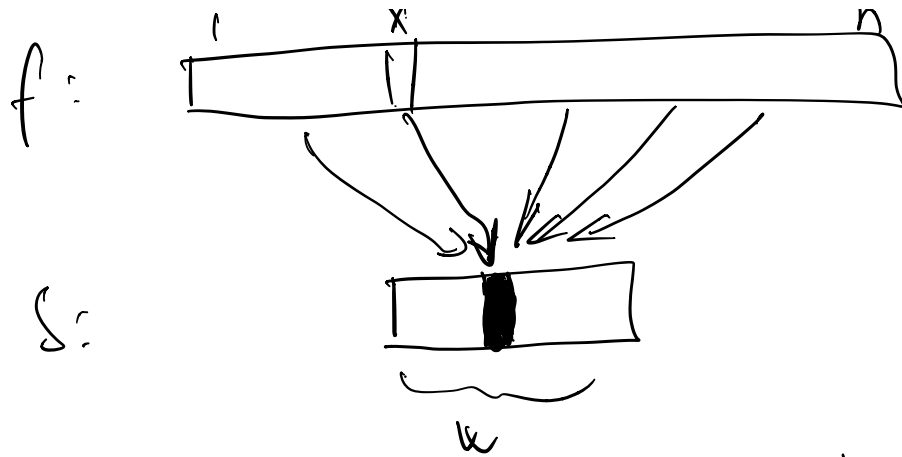
- if $f_x < \phi m - \epsilon \phi m = (1-\epsilon)\phi m$,

then not reported as HH.

↳ $\hat{f}_x < \phi m$ with $\geq 90\%$.

Issue:

↳ some non-HH x 's are reported as HH's.



In fact all y 's s.t. $h(y) = h(x)$
~~are~~ will be considered heavy.

Full Count Min algo

Repeat "basic idea" for $L = O(\lg n)$ times.

Algo:

- inits $S_i [1..w]$, $h_i : [n] \rightarrow w$

$i = 1..L$

- @ iter x in stream:

$$S_i [h_i(x)] := S_i [h_i(x)] + 1$$

$i = 1..L$

- estimator: $\hat{f}_x := \min_{i=1..L} S_i [h_i(x)]$

L indep experiments

Claim: for $w = \frac{10}{\epsilon \phi}$,

$$L = O(\lg n)$$

$$\left[\vec{f}_x - f_x \leq \epsilon \phi_m, \text{ for all } x \in [n] \right]$$

with probability $\geq 1 - \frac{1}{n}$.

pf: by Claim above, \forall fixed $i \in [L]$:

$$\Pr_{h_i} [S_i[h_i(x)] \Rightarrow f_x > \epsilon \phi_m] < \frac{1}{10}$$

for fixed $x \in [n] \rightarrow \Pr_{h_1, \dots, h_L} [\vec{f}_x - f_x > \epsilon \phi_m]$

$$= \prod_{i=1}^L \Pr_{h_i} [S_i[h_i(x)] - f_x > \epsilon \phi_m]$$

$$\leq \left(\frac{1}{10}\right)^L = e^{-L \cdot \ln 10} \stackrel{\rightarrow = \ln n^2}{\approx} \frac{1}{n^2}$$

$$\text{for } L = \frac{2 \cdot \ln n}{\ln 10} = O(\lg n).$$

$$\Pr [\exists x \text{ s.t. } \vec{f}_x - f_x > \epsilon \phi_m] \leq \quad \leftarrow \text{union bound}$$

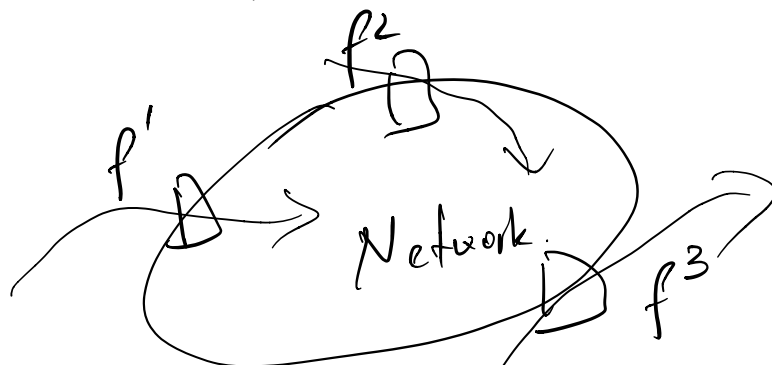
$$\leq \sum_{x \in [n]} \Pr [\hat{f}_x - f_x > \epsilon \phi^m]$$

$$\leq n \cdot \frac{1}{n^2} = \frac{1}{n}. \quad \square$$

Space: $O(L \cdot w)$ words (count up to m)

$\Rightarrow O\left(\frac{\lg n}{\epsilon \phi}\right)$ words for ϕ -HH
(up to $1 \pm \epsilon$ approx. ^{*})

Obs: consider situation where:

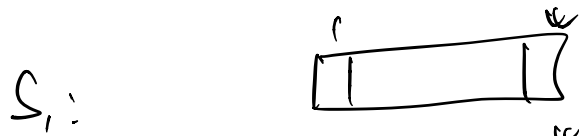
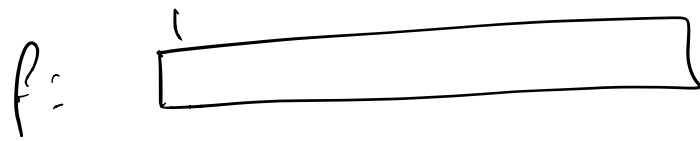


$$f = f^1 + f^2 + f^3 \quad \rightarrow \quad f_x = f_x^1 + f_x^2 + f_x^3 \quad \forall x \in [n].$$

How to compute sketch that able to estimate \hat{f}_x for all $x \in [n]$.

Solution: each router computes
 CM sketch of its freq. vect.
 (using same hash funct.),
 then sum it all up.

CM sketch $(x) \in \mathbb{R}^{2 \cdot w}$



$S =$ $A \cdot f$
 A matrix of size
 $2w \times n$.

$A_{(ij),x} = 1$ iff $h_i(x) = j$
 $i \in [2], j \in [w], x \in [n]$.

$$Af_1 + Af_2 + Af_3 = A(f_1 + f_2 + f_3)$$

—
sketch
@ router 1

sketch of
freq. vector $f_1 + f_2 + f_3$.

Cost Min is a linear sketch.