

# Count-min & Count-Sketch 26

Goal: <sup>turnstile</sup> given  $\hat{\cdot}$  stream of items  
Find "Heavy Hitter":  
Most frequent items.

$X \in \mathbb{R}^n$ ,  $x_i = \#$  times  $i$  appears  
Find  $\{(i, x_i) \mid x_i \text{ "large"}\}$

$\hat{X} \in \mathbb{R}^n$  with  $\|X - \hat{X}\|_\infty < \boxed{\text{bound}}$

How small can error be?

Simplest:  $\epsilon \cdot \|X\|_1 \leftarrow$  "Heavy Hitters"  
Doable w/  $O\left(\frac{\log n}{\epsilon}\right)$  words

$l_\infty / l_1$   
guarantee

Fancier:  $\frac{1}{K} \|X - X_K\|_1$   
where  $X_K \in \mathbb{R}^n$  has  $K$  largest entries of  $X$

Not  $\epsilon \cdot \#$  items seen, but  
 $\epsilon \cdot \#$  less frequent items seen.

Still  $O(K \log n)$  words  
Count-min sketch, Comale-Muthukrishnan

$l_\infty / l_2$   
guarantee

Even fancier:  $\frac{1}{\sqrt{K}} \cdot \|X - X_K\|_2$

$O(K \log n)$ , Count-Sketch,  
Charikar, Chen, Farach-Colton

HW will show:  $\frac{1}{\sqrt{K}} \|X - X_K\|_2 \leq \frac{1}{K} \|X - X_K\|_1$

How do these behave?

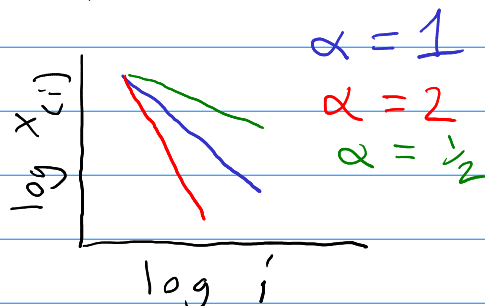
Zipf's law:  $i^{\text{th}}$  most common word in English appears  $\sim \frac{1}{i}$  times

Power law:  $i^{\text{th}}$  largest  $X_j \sim \frac{1}{i^\alpha}$  for some  $\alpha > 0$ .

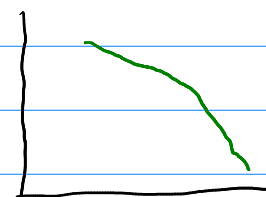
Seen in URLs in web  
Frequencies in music  
Population in cities

[approximately] when it comes to largest entries]

Log log plot: line



actual data, often



Similar behavior: lognormal, —

$$\|X\|_1 = \sum_{i=1}^n \frac{1}{i^\alpha} \sim \begin{cases} 1 & \text{if } \alpha > 1 \\ \log n & \text{if } \alpha = 1 \\ n^{1-\alpha} & \text{if } \alpha < 1 \end{cases}$$

$$\|X\|_2 = \sqrt{\sum_{i=1}^n \frac{1}{i^{2\alpha}}} \sim \begin{cases} 1 & \text{if } \alpha > \frac{1}{2} \\ \sqrt{\log n} & \text{if } \alpha = \frac{1}{2} \\ n^{\frac{1}{2}-\alpha} & \text{if } \alpha < \frac{1}{2} \end{cases}$$

What effect does  $\|x - x_k\|$  have?

$$\|x - x_k\|_1 = \sum_{i=k+1}^n \frac{1}{i^\alpha} = \begin{cases} k^{1-\alpha} & \text{if } \alpha > 1 \\ \log \frac{n}{k} & \text{if } \alpha = 1 \\ n^{1-\alpha} & \text{if } \alpha < 1 \end{cases}$$

Suppose  $\alpha > 1$ .

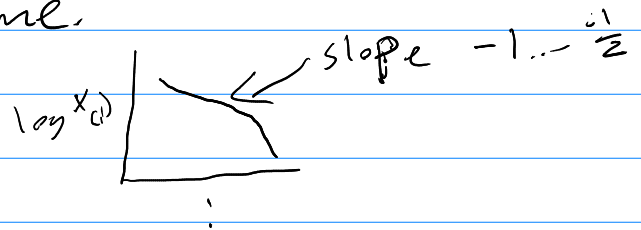
$$\frac{1}{2} < \alpha < 1$$

Heavy Hitters:  $\frac{1}{k} \|x\|_1 \sim \frac{1}{k} \frac{n^{1-\alpha}}{k}$

Count-min:  $\frac{1}{k} \|x - x_k\|_1 \sim \frac{1}{k^\alpha} \frac{n^{1-\alpha}}{k}$

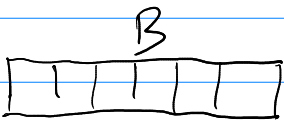
Count-sketch:  $\frac{1}{\sqrt{k}} \|x - x_k\|_2 \sim \frac{1}{k^\alpha} \frac{1}{k^\alpha}$

Count-sketch much better if  $\frac{1}{2} < \alpha < 1$ .  
Most common regime.



How they work



Hash to  $y$ :   $B = o(k)$

$h: [n] \rightarrow [B]$

How to deal with collisions?

$$y_j = \sum_{\substack{u \in [n] \\ h(u) = j}} x_u$$

Estimate  $\hat{x}_u$  as  $y_{h(u)}$

In strict turnstile ( $x_u \geq 0 \forall u$ ):  
 $\hat{x}_u \geq y_{h(u)}$

In general,

$$\begin{aligned} E[\hat{x}_u - x_u] &= E\left[\sum_{u' \neq u} x_{u'} \cdot \mathbb{1}_{h(u') = h(u)}\right] \\ &= \sum_{u' \neq u} x_{u'} \cdot \Pr[h(u') = h(u)] \end{aligned}$$

$$\leq \|x\|_1 \cdot \frac{1}{B} \quad \text{w/ Pairwise independence}$$

So if  $B = 4k$ ,  $|\hat{x}_u - x_u| \leq \frac{\|x\|_1}{k}$  w. p.  $\frac{3}{4}$

$$|\hat{x}_u - x_u| \leq \frac{\|x\|_1}{k} \quad \text{w.p. } \frac{3}{4} \quad \text{for fixed } u.$$

Want  $\|\hat{x}_u - x_u\| \leq \dots$ , i.e. all  $u$ .

Repeat  $R$  times:

$h_1: y^{(1)} = \text{[rectangle with } \beta \text{ above]}$   
 $h_2: y^{(2)} = \text{[rectangle]}$   
 $\vdots$   
 $h_R: y^{(R)} = \text{[rectangle with } \beta \text{ above]}$

Each  $\hat{x}_u^{(r)}$  is close to  $x_u$  w/  $\frac{3}{4}$  prob.  
 Different  $r$  are independent.  
 How to combine?

Strict turnstile:  $\hat{x}_u^{(r)} \geq x_u \quad \forall r$

Set  $\hat{x}_u = \min_r \hat{x}_u^{(r)}$

works if any of  $R$  work. Failure prob  $4^{-R}$   
 $R = \log n \Rightarrow$  all  $n$   $\hat{x}_u$  work.

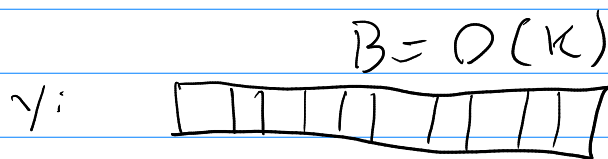
Each  $r$  uses  $O(k)$  space  
 $\Rightarrow O(k \log n)$  total

to solve heavy hitters

# Extensions

- What about nonstrict turnstile?  $\rightarrow$ 
  - still  $|\hat{x}_u^{(r)} - x_u| \leq \frac{\|x\|_1}{K}$  w.p.  $\frac{3}{4}$ .
  - but  $\min \hat{x}_u^{(r)}$  is bad.  
Replace min w/ median.  
Same asymptotics, worse constant.

- bound w/  $\frac{\|x - x_K\|_1}{K}$



Split  $x$  into  $x_K$ : largest  $K$  terms  
 $x - x_K$ : rest

Most likely:

-  $h(u)$  will not collide w/ any of top  $K$   
[Pr.  $1 - \frac{K}{B}$ ]

- error from rest  $\leq \frac{\|x - x_K\|_1}{K}$  [as before]

- bound w/  $\frac{\|x - x_K\|_2}{\sqrt{K}}$

problem: what if all  $x_u = 1$ ?  $\rightarrow$

$$\forall j, \quad y_j \approx \frac{n}{B} \Rightarrow \hat{x}_u \approx \frac{n}{B} = \Theta\left(\frac{\|x\|_2}{K}\right)$$
$$|\hat{x}_u - 1| \gg \frac{\|x\|_2}{\sqrt{K}} = \sqrt{\frac{n}{K}}$$

Change to

$$Y_i^{(r)} = \sum_{u: h(u)=i} x_u \cdot S_r(u)$$

for  $S_r: [n] \rightarrow \{-1, 1\}$

$$\hat{x}_u^{(r)} = Y_{h(u)}^{(r)} \cdot S_r(u)$$

$$\hat{x}_u = \text{median}_r \hat{x}_u^{(r)}$$

Analysis:

$$E[(\hat{x}_u^{(r)} - x_u)^2] = E\left[\left(\sum_{u' \neq u} 1_{h(u')=h(u)} \cdot x_{u'} \cdot S(u) \cdot S(u')\right)^2\right]$$

w/ Pairwise indep.

of  $S$ , cross terms cancel!

$$= \sum_{u' \neq u} x_{u'}^2 \cdot \Pr[h(u')=h(u)]$$

$$= \frac{\|x\|_2^2}{B}$$

Hence  $|\hat{x}_u^{(r)} - x_u| \leq \frac{\|x\|_2}{\sqrt{B}}$  w.p.  $\frac{3}{4}$ .

Rest, same as before.