# Algorithms for Big Data (FALL 25)

Lecture 7

HEAVY HITTERS: MISRA-GRIES, COUNTMIN AND COUNTSKETCH

ALI VAKILIAN (vakilian@vt.edu)





#### Frequent Items Problem ( $F_{\infty}$ -Moment)

**Recall:** What is  $F_{\infty}$ ?

- $F_{\infty}$  is very brittle and hard to estimate with low memory.
- Even strong lower bounds even for very weak relative approximations.

Hence, we need to settle for weaker (additive) guarantees.

**Heavy Hitters Problem:** Find all items i such that  $f_i > \frac{m}{k}$  for some fixed k. Heavy hitters are **very** frequent items.

## Finding Majority Element (interview question)



**Offline:** given an array/list A of m integers, is there an element that occurs more than m/2 times in A?

**Streaming:** is there an i such that  $f_i > m/2$ ?

### **Boyer-Moore Voting**

**Lemma.** If there exists a majority item i, the algorithm outputs s = i and  $c \ge f_i - \frac{m}{2}$ .

Why it works?

```
Majority (in streams):
   let c \leftarrow 0, s \leftarrow null
   foreach item e_i in the stream do:
         if e_i = s then
            c \leftarrow c + 1
         else if c = 0
            c \leftarrow 1 and s \leftarrow e_i
         else
            c \leftarrow c - 1
   return c and s
```

### **Boyer-Moore Voting**

**Lemma.** If there exists a majority item i, the algorithm outputs s = i and  $c \ge f_i - \frac{m}{2}$ .

Why it works?

What if no majority item exists? How to verify?

```
Majority (in streams):
   let c \leftarrow 0, s \leftarrow null
   foreach item e_i in the stream do:
         if e_i = s then
            c \leftarrow c + 1
         else if c = 0
            c \leftarrow 1 and s \leftarrow e_i
         else
            c \leftarrow c - 1
   return c and s
```

#### Extension to k Heavy Hitters

**Offline:** given an array/list A of m integers, is there an element that occurs more than m/k times in A?

**Streaming:** is there an i such that  $f_i > m/k$ ?

Idea. Extending Boyer-More Voting algorithm to this more general setting.

Space: O(k)

**Theorem.** For each  $i \in [n]$ ,  $f_i - \frac{m}{k+1} \le \hat{f_i} \le f_i$ 

 $\implies$  any item with  $f_i > \frac{m}{k}$  is in D.

```
Misra-Gries (k):
   let D be an empty array of size k
   foreach item e_i in the stream do
         if e_i \in D then
           D[e_i] \leftarrow D[e_i] + 1
         else if D has less than k items
            add e_i to D and set D[e_i] \leftarrow 1
         else
           foreach \ell \in D do
              D[\ell] \leftarrow D[\ell] - 1 (if 0, remove)
   foreach \ell \in D, set \hat{f}_{\ell} \leftarrow D[\ell] (zero for rest)
```

**Theorem.** For each  $i \in [n]$ ,  $f_i - \frac{m}{k+1} \le \hat{f_i} \le f_i$ . **Proof.** 

•  $\hat{f}_i \leq f_i$  is easy.

**Theorem.** For each  $i \in [n]$ ,  $f_i - \frac{m}{k+1} \le \hat{f_i} \le f_i$ .

Proof.

#### Alternative view of algorithm.

- maintain count C[i] for each i (initialized to 0).  $\leq k$  are nonzero at anytime.
- when next element  $e_i$  arrives:
  - if  $C[e_j] > 0$  then increment  $C[e_j]$
  - else if < k positive counters, then **set**  $C[e_j] = 1$
  - lacktriangle else, **decrement all** positive counters (exactly k of them)
- output  $\hat{f}_i = C[i]$  for each i

**Theorem.** For each  $i \in [n]$ ,  $f_i - \frac{m}{k+1} \le \hat{f_i} \le f_i$ .

Goal. 
$$f_i - \hat{f}_i \leq \frac{m}{k+1}$$

• Suppose **decrement all** occur  $\ell$  times, then  $\ell k + \ell \leq m \Longrightarrow \ell \leq \frac{m}{k+1}$ 

"each decrement all remove k previously added items and involves an insertion causing this operation. (each deals with k+1 distinct elements)"

Define  $\alpha = f_i - \hat{f_i}$ . It is initially zero (as both are equal to zero).

How big can it get?

**Theorem.** For each  $i \in [n]$ ,  $f_i - \frac{m}{k+1} \le \hat{f_i} \le f_i$ .

Define  $\alpha = f_i - \hat{f_i}$ . It is initially zero (as both are equal to zero).

#### How big can it get?

- If  $e_i = i$  and C[i] is incremented, then  $\alpha$  stays the same.
- If  $e_j = i$  and C[i] is not incremented, then  $\alpha$  increases by one and k counters decremented (charge to one of the  $\ell$  events).
- If  $e_j \neq i$  and C[i] is decremented, then  $\alpha$  increases by one. This only happens in **decrement all** scenario (again charge to one of the  $\ell$  events).

So, 
$$\alpha \le \ell \le m/(k+1)$$

#### Wrap-up: Deterministic vs Randomized

- Cannot improve O(k) space if one wants additive error of at most m/k.
- Somewhat rare to have a deterministic algorithm that is near-optimal.

#### Why may we still look for randomized solutions?

- Supporting deletions
- Extra properties of sketch-based solutions

### Basic Hashing Idea

**Heavy Hitters Problem:** Find all items i such that  $f_i \ge m/k$ .

- Let  $b_1, ..., b_k$  be the k heavy hitters (at most k)
- Suppose we pick a hash function  $h: [n] \to [ck]$  for some c > 1
- h maps the heavy hitters into different buckets (k balls into ck bins)
- Then, ideally, we would like to use the count of items in each bucket as an estimate for the frequency of one heavy hitters.

Repeating this idea with independent hashes improves the estimate

#### CountMin Sketch [Cormode-Muthukrishnan]

- k pairwise independent hash functions  $h_1, ..., h_k$ ; each  $[n] \to [w]$
- Equivalently, a table C with k rows and w columns.
- Store one counter per entry in the table, which keep the aggregate frequency of items mapped to the entry by the corresponding hash function.  $C[\ell, s]$  is the counter for bucket s in hash function  $h_{\ell}$ .
- Let  $f \in \mathbb{R}^n$  be the final frequency vector. For  $\ell \in [k]$ ,  $s \in [w]$ ,

$$C[\ell, s] = \sum_{i:h_{\ell}(i)=s} f_i$$

- For every  $\ell \in [k]$ ,  $C[\ell, h_{\ell}(i)]$  is an over-estimate of  $f_i$ .
- $\circ$  We have k such estimate, how good is the quality of best of them?

### CountMin Sketch in Streaming

• Each of k estimates for  $f_i$  is overcounting its frequency

• Picking the minimum such estimate is reasonable.

#### **CountMin Sketch (stream):**

**let**  $h_1, ..., h_k$  be pairwise independent hash functions from  $[n] \rightarrow [w]$ 

foreach item  $e_t = (i_t, \Delta_t)$  in the stream do: for  $\ell = 1$  to k do:  $C[\ell, h_{\ell}(i_t)] \leftarrow C[\ell, h_{\ell}(i_t)] + \Delta_t$ 

//frequency estimates

**foreach**  $i \in [n]$ , set  $\tilde{f}_i = \min_{\ell \in [k]} C[\ell, h_{\ell}(i)]$ 

#### CountMin Sketch: Main Property

Theorem. Consider strict turnstile streaming (i.e., always  $f \ge 0$ ). Let  $k = \Omega(\log \frac{1}{\delta})$  and  $w > \frac{2}{\varepsilon}$ . Then, for any fixed  $i \in [n]$ ,  $f_i \le \tilde{f}_i$ , and

$$\Pr[\tilde{f}_i \geq f_i + \varepsilon ||f||_1] \leq \delta$$

- Unlike Misra-Gries, CountMin overestimates.
- Items are not stored (can be recovered via queries).
- Handles deletion (works in strict turnstile model)
- Space complexity:  $O(\frac{\log \frac{1}{\delta}}{\varepsilon} \cdot \log m)$  bits

#### CountMin: Analysis

- Consider an item i and fix a row  $\ell$ .
- Define  $Z_{\ell} = C[\ell, h_{\ell}(i)]$  the value of counter in row  $\ell$  that i is hashed to.

$$\mathbb{E}[Z_{\ell}] = f_i + \sum_{j \neq i} \Pr[h_{\ell}(j) = h_{\ell}(i)] \cdot f_j$$

$$= f_i + \sum_{j \neq i} \frac{1}{w} \cdot f_j \qquad \text{// pairwise independence of } h_{\ell}$$

$$\leq f_i + \varepsilon ||f||_1 / 2 \qquad \text{// } w > 2/\varepsilon$$

Applying Markov,  $\Pr[Z_{\ell} - f_i \ge \varepsilon ||f||_1] \le 1/2$ 

Since k hash functions are independent,

$$\Pr[\min_{\ell \in [k]} Z_{\ell} \ge f_i + \varepsilon ||f||_1] \le \frac{1}{2^k} \le \delta \quad // k = \Omega(\log \frac{1}{\delta})$$

#### CountMin Sketch

## **Space Complexity:** $O(\frac{1}{\varepsilon} \log n \log m)$ bits

Theorem. Consider strict turnstile streaming (i.e., always  $f \ge 0$ ). Let  $k = \Omega(\log \frac{1}{\delta})$  and  $w > \frac{2}{\epsilon}$ . Then, for any fixed  $i \in [n]$ ,  $f_i \le \tilde{f}_i$ , and

$$\Pr[\tilde{f}_i \geq f_i + \varepsilon ||f||_1] \leq \delta$$

• Setting  $\delta = 1/n^2$ , a CountMin with  $O(\log n)$  rows and  $O(1/\varepsilon)$  columns, for every  $i \in [n]$ ,

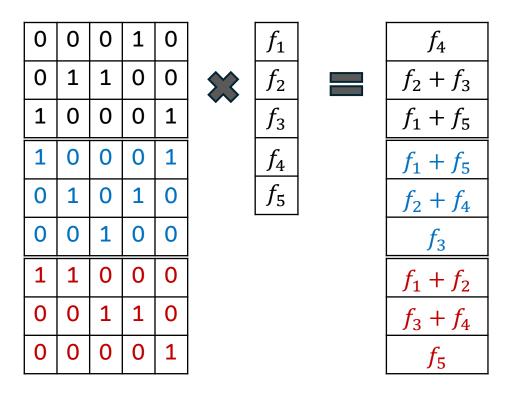
$$\Pr[\tilde{f}_i > f_i + \varepsilon || f ||_1] \le 1/n^2$$

• By union bound over all n items, with probability  $\geq 1 - 1/n$ , for all  $i \in [n]$ 

$$\tilde{f}_i \le f_i + \varepsilon ||f||_1$$

#### CountMin is a Linear Sketch

0	0	0	0	0	*	$f_1$ $f_2$	$f_4$ $f_2 + f_3$
to pucket 3	to bucket <b>2</b>	to bucket <b>2</b>	to bucket 1 O	to pucket <b>3</b>		$f_3$ $f_4$ $f_5$	$f_1 + f_5$
em 1 is hashed to bucket 3	em 2 is hashed to bucket 2	em <b>3</b> is hashed to bucket <b>2</b>	em <b>4</b> is hashed to bucket <b>1</b>	em <b>5</b> is hashed to bucket <b>3</b>			



hash function  $h_i$  as a Matrix-Vector Multiplication  $\Pi_{w \times m} \ m{f}_{m imes 1}$ 

CountMin as a Matrix-Vector Multiplication  $\Pi_{(k\cdot w) imes m} \, m{f}_{m imes 1}$ 

#### CountSketch

- Simialr to CountMin, keeps track of a table of  $k \times w$  counters
- Inspired by AMS sketch, assign u.a.r signs  $\{-1, +1\}$  to items
- Counters can get even negative

#### CountSketch (stream):

**let**  $h_1, ..., h_k$  be pairwise independent hash functions from  $[n] \rightarrow [w]$ 

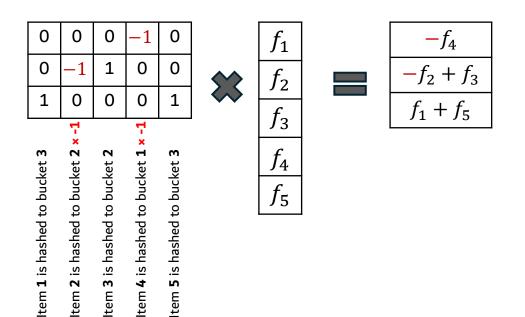
**let**  $g_1, ..., g_k$  be pairwise independent hash functions from  $[n] \rightarrow \{-1, +1\}$ 

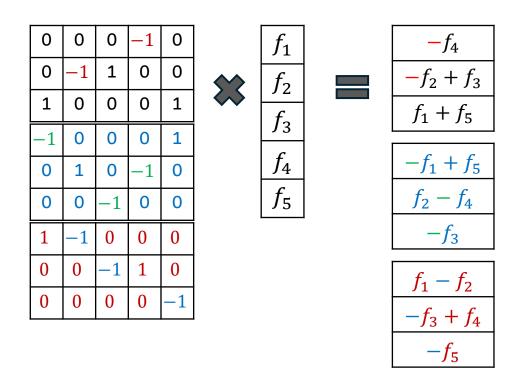
foreach item  $e_t = (i_t, \Delta_t)$  in the stream do: for  $\ell = 1$  to k do:  $\mathcal{C}[\ell, h_\ell(i_t)] \leftarrow \mathcal{C}[\ell, h_\ell(i_t)] + g_\ell(i_t) \cdot \Delta_t$ 

//frequency estimates

**foreach**  $i \in [n]$ , set  $\tilde{f}_i = \underset{\ell \in [k]}{\operatorname{median}} \{g_{\ell}(i) \cdot C[\ell, h_{\ell}(i)]\}$ 

#### Why CountSkecth is a Linear Sketch?





#### CountSketch: Main Property

Theorem. Consider strict turnstile streaming (i.e., always  $f \ge 0$ ). Let  $k = \Omega(\log \frac{1}{\delta})$  and  $w > \frac{3}{\varepsilon^2}$ . Then, for any fixed  $i \in [n]$ ,  $\mathbb{E}[\tilde{f}_i] = f_i$ , and

$$\Pr[\left|\tilde{f}_i - f_i\right| \ge \varepsilon ||f||_2] \le \delta$$

#### **Comparison to CountMin**

- Error is w.r.t.  $||f||_2$  instead of  $||f||_1$ . Note  $||f||_2 \le ||f||_1$ , and in some cases  $||f||_2 \ll ||f||_1$
- Space complexity:  $O(\frac{1}{\varepsilon^2} \cdot \log n)$  bits

### CountSketch Analysis

- Consider an item i and fix a row  $\ell$ .
- Define  $Z_{\ell} = g_{\ell}(i) \mathcal{C}[\ell, h_{\ell}(i)]$  the value of counter in row  $\ell$  that i is hashed to.

For  $j \in [n]$  let  $Y_j$  be the indicator r.v. that is 1 if  $h_{\ell}(i) = h_{\ell}(j)$ ; i.e., i and j collide in  $h_{\ell}$ 

$$\mathbb{E}ig[Y_jig] = \mathbb{E}ig[Y_j^2ig] = 1/w$$
 from pairwise independence of  $h_\ell$ 

$$Z_{\ell} = g_{\ell}(i)C[\ell, h_{\ell}(i)] = g_{\ell}(i)f_i + \sum_{j \neq i} g_{\ell}(i)f_jY_j$$

$$\mathbb{E}[Z_{\ell}] = f_i + \sum_{j \neq i} \mathbb{E}[g_{\ell}(i)g_{\ell}(j)Y_j] \cdot f_j$$

$$= f_i \qquad // \text{ pairwise independence of } g_{\ell}$$

Since 
$$\mathbb{E}[g_{\ell}(i)g_{\ell}(j)Y_j] = \mathbb{E}[g_{\ell}(i)g_{\ell}(j)]\mathbb{E}[Y_j] = 0$$

#### CountSketch Analysis: Variance

• Define  $Z_{\ell} = g_{\ell}(i) \mathcal{C}[\ell, h_{\ell}(i)]$  the value of counter in row  $\ell$  that i is hashed to. For  $j \in [n]$  let  $Y_i$  be the indicator r.v. that is 1 if  $h_{\ell}(i) = h_{\ell}(j)$ ; i.e., i and j collide in  $h_{\ell}$  $\mathbb{E}[Y_i] = \mathbb{E}[Y_i^2] = 1/w$  from pairwise independence of  $h_\ell$  $=\mathbb{E}[(Z_{\ell}-f_i)^2]$  $= \mathbb{E} \left| \left( \sum_{j \neq i} g_{\ell}(i) g_{\ell}(j) Y_{j} f_{j} \right)^{2} \right|$  $= \mathbb{E} \left[ \sum_{j \neq i} g_{\ell}(i)^{2} g_{\ell}(j)^{2} Y_{i}^{2} f_{i}^{2} + \sum_{i,i' \neq i} g_{\ell}(i)^{2} g_{\ell}(j) g_{\ell}(j') Y_{i} Y_{i'} f_{i} f_{i'} \right]$  $=\sum_{i\neq i}f_i^2\mathbb{E}[Y_i^2]$  $\leq \|f\|_2^2/w$ 

Using Chebyshev, 
$$\Pr[|Z_{\ell} - f_i| \ge \varepsilon ||f||_2] \le \frac{\operatorname{Var}(Z_{-\ell})}{\varepsilon^2 ||f||_2^2} \le \frac{1}{\varepsilon^2 w} \le 1/3$$

#### Countsketch: Concentration

Using Chebyshev, 
$$\Pr[|Z_{\ell} - f_i| \ge \varepsilon ||f||_2] \le \frac{\operatorname{Var}(Z_{-\ell})}{\varepsilon^2 ||f||_2^2} \le \frac{1}{\varepsilon^2 w} \le 1/3$$

Then, by Chernoff bound,

$$\Pr[|\text{median}\{Z_1, \dots, Z_k\} - f_i| \ge \varepsilon ||f||_2] \le e^{-\Omega(k)} \le \delta$$