## CS 498ABD: Algorithms for Big Data

Heavy Hitters / Frequent I law

Lecture 08 September 17, 2020

### Models

### Richer model:

- Want to estimate a function of a vector x ∈ ℝ<sup>n</sup> which is initially assume to be the all 0's vector.
- Each element e<sub>j</sub> of a stream is a tuple (i<sub>j</sub>, Δ<sub>j</sub>) where i<sub>j</sub> ∈ [n] and Δ<sub>i</sub> ∈ ℝ is a real-value: this updates x<sub>ij</sub> to x<sub>ij</sub> + Δ<sub>j</sub>. (Δ<sub>j</sub> can be positive or negative)

### Models

### Richer model:

- Want to estimate a function of a vector x ∈ ℝ<sup>n</sup> which is initially assume to be the all 0's vector.
- Each element e<sub>j</sub> of a stream is a tuple (i<sub>j</sub>, Δ<sub>j</sub>) where i<sub>j</sub> ∈ [n] and Δ<sub>i</sub> ∈ ℝ is a real-value: this updates x<sub>ij</sub> to x<sub>ij</sub> + Δ<sub>j</sub>. (Δ<sub>j</sub> can be positive or negative)
- $\Delta_j > 0$ : cash register model. Special case is  $\Delta_j = 1$ .
- $\Delta_i$  arbitrary: *turnstile* model
- $\Delta_j$  arbitrary but  $x \ge 0$  at all times: *strict turnstile* model
- *Sliding window* model: interested only in the last *W* items (window)

What is  $F_k$  when  $k = \infty$ ?

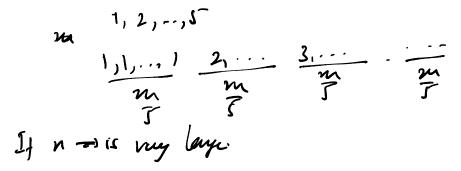
 $(f_{1}, f_{2}, \dots, f_{n})$ 

Fo = max fi

What is  $F_k$  when  $k = \infty$ ? Maximum frequency.

What is  $F_k$  when  $k = \infty$ ? Maximum frequency.

 $F_{\infty}$  very brittle and hard to estimate with low memory. Can show strong lower bounds for very weak relative approximations.



What is  $F_k$  when  $k = \infty$ ? Maximum frequency.

 $F_{\infty}$  very brittle and hard to estimate with low memory. Can show strong lower bounds for very weak relative approximations.

Hence settle for weaker (additive) guarantees.

What is  $F_k$  when  $k = \infty$ ? Maximum frequency.

 $F_{\infty}$  very brittle and hard to estimate with low memory. Can show strong lower bounds for very weak relative approximations.

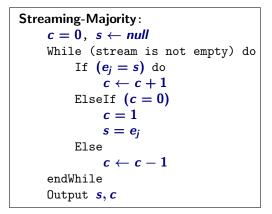
Hence settle for weaker (additive) guarantees.

Heavy Hitters Problem: Find all items i such that  $f_i > m/k$  for some fixed k.

Heavy hitters are very frequent items.

### Majority element problem:

- 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  $f_i > m/2$ ?



Chandra (UIUC)

```
Streaming-Majority:
     c = 0, s \leftarrow null
     While (stream is not empty) do
          If (e_i = s) do
               c \leftarrow c + 1
          ElseIf (c = 0)
               c = 1
               s = e_i
          Else
               c \leftarrow c - 1
     endWhile
     Output s, c
```

**Claim:** If there is a majority element *i* then algorithm outputs s = i and  $c \ge f_i - m/2$ .

```
Streaming-Majority:
     c = 0, s \leftarrow null
     While (stream is not empty) do
          If (e_i = s) do
               c \leftarrow c + 1
          ElseIf (c = 0)
               c = 1
               s = e_i
          Else
               c \leftarrow c - 1
     endWhile
     Output s, c
```

**Claim:** If there is a majority element *i* then algorithm outputs s = i and  $c \ge f_i - m/2$ . **Caveat:** Algorithm may output incorrect element if no majority element. Can verify correctness in a second pass.

Chandra (	UIUC

CS498ABD

C

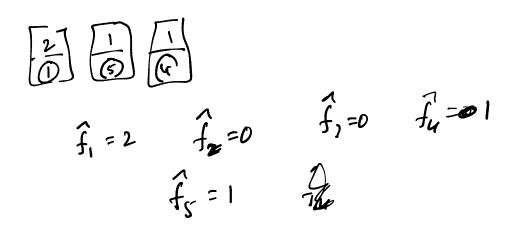
# **Misra-Gries Algorithm**

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

```
MisraGreis(k):
     D is an empty associative array
     While (stream is not empty) do
          e; is current item
          If (e; is in keys(D))
               D[e_i] \leftarrow D[e_i] + 1
          Else if (|keys(A)| < k - 1) then
          D[e_i] \leftarrow 1
         Else
              for each \ell \in keys(D) do
                    D[\ell] \leftarrow D[\ell] - 1
          Remove elements from D whose counter values are 0
endWhile
For each i \in keys(D) set \hat{f}_i = D[i]
For each i \notin keys(D) set f_i = 0
 Chandra (UIUC)
                            CS498ABD
                                          6
                                                                Fall 2020
```

6/1

k=3 1,2,1,4,5,1,2,10,1,3,5,4,--



## Analysis

Space usage O(k).

#### Theorem

For each 
$$i \in [n]$$
:  $f_i - \frac{m}{k+1} \leq \hat{f}_i \leq f_i$ .

#### Corollary

Any item with  $f_i > m/k$  is in D at the end of the algorithm.

A second pass to verify can be used to verify correctness of elements in D.

### Theorem

For each 
$$i \in [n]$$
:  $f_i - \frac{m}{k+1} \leq \hat{f}_i \leq f_i$ .

$$\hat{f}_i \neq \max \{0, f_i - \frac{m}{k+1}\}$$

### Theorem

For each 
$$i \in [n]$$
:  $f_i - \frac{m}{k+1} \leq \hat{f}_i \leq f_i$ .

Easy to see:  $\hat{f}_i \leq f_i$ . Why?

#### Theorem

For each 
$$i \in [n]$$
:  $f_i - \frac{m}{k+1} \leq \hat{f}_i \leq f_i$ .

Easy to see:  $\hat{f}_i \leq f_i$ . Why?

Alternative view of algorithm:

- Maintains counts *C*[*i*] for each *i* (initialized to 0). Only *k* are non-zero at any time.
- When new element *e<sub>j</sub>* comes
  - If  $C[e_j] > 0$  then increment  $C[e_j]$
  - Elself less then k positive counters then set  $C[e_j] = 1$
  - Else decrement all positive counters (exactly **k** of them)

8

• Output  $\hat{f}_i = C[i]$  for each i

Want to show:  $f_i - \hat{f}_i \leq m/(k+1)$ :  $\hat{f}_i \geq f_i - \frac{m}{k+1}$   $\hat{f}_i \neq fi$ 

Want to show:  $f_i - \hat{f}_i \leq m/(k+1)$ :

Suppose we have ℓ occurrences of k counters being decremented.

Want to show:  $f_i - \hat{f}_i \leq m/(k+1)$ :

- Suppose we have ℓ occurrences of k counters being decremented. Then ℓk + ℓ ≤ m which implies ℓ ≤ m/(k + 1).
- Consider  $\alpha = (f_i \hat{f}_i)$  as items are processed. Initially 0. How big can it get?

Want to show:  $f_i - \hat{f}_i \leq m/(k+1)$ :

- Suppose we have ℓ occurrences of k counters being decremented. Then ℓk + ℓ ≤ m which implies ℓ ≤ m/(k + 1).
- Consider  $\alpha = (f_i \hat{f}_i)$  as items are processed. Initially 0. How big can it get?

• If  $e_j = i$  and C[i] is incremented  $\alpha$  stays same

Want to show:  $f_i - \hat{f}_i \leq m/(k+1)$ :

- Suppose we have ℓ occurrences of k counters being decremented. Then ℓk + ℓ ≤ m which implies ℓ ≤ m/(k + 1).
- Consider  $\alpha = (f_i \hat{f}_i)$  as items are processed. Initially 0. How big can it get?
  - If  $e_j = i$  and C[i] is incremented  $\alpha$  stays same
  - If e<sub>j</sub> = i and C[i] is not incremented then α increases by one and k counters decremented — charge to l

Want to show:  $f_i - \hat{f}_i \leq m/(k+1)$ :

- Suppose we have ℓ occurrences of k counters being decremented. Then ℓk + ℓ ≤ m which implies ℓ ≤ m/(k + 1).
- Consider  $\alpha = (f_i \hat{f}_i)$  as items are processed. Initially 0. How big can it get?
  - If  $e_j = i$  and C[i] is incremented  $\alpha$  stays same
  - If e<sub>j</sub> = i and C[i] is not incremented then α increases by one and k counters decremented — charge to l
  - If e<sub>j</sub> ≠ i and α increases by 1 it is because C[i] is decremented
     charge to ℓ

Want to show:  $f_i - \hat{f}_i \leq m/(k+1)$ :

- Suppose we have ℓ occurrences of k counters being decremented. Then ℓk + ℓ ≤ m which implies ℓ ≤ m/(k + 1).
- Consider  $\alpha = (f_i \hat{f}_i)$  as items are processed. Initially 0. How big can it get?
  - If  $e_j = i$  and C[i] is incremented  $\alpha$  stays same
  - If e<sub>j</sub> = i and C[i] is not incremented then α increases by one and k counters decremented — charge to l
  - If e<sub>j</sub> ≠ i and α increases by 1 it is because C[i] is decremented
     charge to ℓ

9

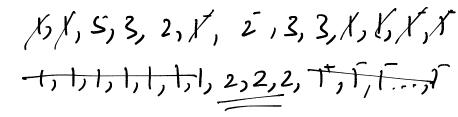
• Hence total number of times  $\alpha$  increases is at most  $\ell$ .

# **Deterministic to Randomized Sketches**

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

Why look for randomized solution?

- Obtain a sketch that allows for deletions
- Additional applications of sketch based solutions
- Will see Count-Min and Count sketches



# **Basic Hashing/Sampling Idea**

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

- Let  $b_1, b_2, \ldots, b_k$  be the k heavy hitters
- Suppose we pick h:[n] 
  ightarrow [ck] for some c>1
- h spreads  $b_1, \ldots, b_k$  among the buckets (k balls into ck bins)
- In ideal situation each bucket can be used to count a separate heavy hitter

