# CS 473: Algorithms

Ruta Mehta

University of Illinois, Urbana-Champaign

Spring 2021

# CS 473: Algorithms, Spring 2021

# **Streaming Algorithms**

Lecture 12 March 16, 2021

Most slides are courtesy Prof. Chekuri

# Streaming Algorithms

A topic that is both very old, and very current!

#### Dawn of CS..

Data was stored on tapes, and amount of RAM was very small.

- Too much data, too little space.
- Store only summary or sketch of data.

# Streaming Algorithms

A topic that is both very old, and very current!

#### Dawn of CS..

Data was stored on tapes, and amount of RAM was very small.

- Too much data, too little space.
- Store only summary or sketch of data.

#### Now..

Terabytes of memory, Gigabytes of RAM.

- Data streams: Humongous amount of data (sometimes never ending)!
- Can go over it at most once, and sometimes not even that!
- Store only summary: sub-linear space-time algorithms.

An internet router sees a stream of packets, and may want to know,

- which connection is using the most packets
- how many different connections
- median of the file sizes transferred since mid-night
- $\bullet$  which connections are using more than 0.1% of the bandwidth.

Computing aggregative information about data streams.

#### Outline

Computation with data streams.

#### Heavy-hitters

- Majority element (by R. Boyer and J.S. Moore)
- ε-heavy hitters deterministic
- Approximate counting

Counting using hashing – Count-min Sketch (Cormode-Muthukrishnan'05)

Variant of Bloom filters.

#### Data Streams

A stream of data elements,  $S = a_1, a_2, \ldots$ 

Say  $a_t$  arrive at time t. Let us assume that  $a_t$ 's are numbers for this lecture.

#### Data Streams

A stream of data elements,  $S = a_1, a_2, \ldots$ 

Say  $a_t$  arrive at time t. Let us assume that  $a_t$ 's are numbers for this lecture.

Denote 
$$a_{[1..t]} = \langle a_1, a_2, \ldots, a_t \rangle$$
.

Given some function we want to compute it continually, while using limited space.

• at any time t we should be able to query the function value on the stream seen so far, i.e.,  $a_{[1..t]}$ .

$$S = 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32, ...$$

#### Computing Sum

$$F(a_{[1..t]}) = \sum_{i=1}^t a_i$$

Outputs are: 3, 4, 21, 25, 16, 48, 149, 152, -570, ...

$$S = 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32, ...$$

#### Computing Sum

$$F(a_{[1..t]}) = \sum_{i=1}^t a_i$$

Outputs are: 3, 4, 21, 25, 16, 48, 149, 152, -570, ...

Keep a counter, and keep adding to it.

After T rounds, the number can be at most  $T2^b$ .  $O(b + \log T)$  space.

$$S = 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32, ...$$

#### Computing max

$$F(a_{[1..t]}) = \max_{i=1}^t a_i$$

Outputs are: 3, 3, 17, 17, 17, 32, 101, 101, ...

Just need to store **b** bits.

$$S = 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32, ...$$

#### Computing max

$$F(a_{[1..t]}) = \max_{i=1}^t a_i$$

Outputs are: 3, 3, 17, 17, 17, 32, 101, 101, ...

Just need to store **b** bits.

Median?

$$S = 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32, ...$$

#### Computing max

$$F(a_{[1..t]}) = \max_{i=1}^t a_i$$

Outputs are: 3, 3, 17, 17, 17, 32, 101, 101, ...

Just need to store **b** bits.

Median? A lot more tricky

$$S = 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32, ...$$

#### Computing max

$$F(a_{[1..t]}) = \max_{i=1}^t a_i$$

Outputs are: 3, 3, 17, 17, 17, 32, 101, 101, ...

Just need to store **b** bits.

Median? A lot more tricky

# distinct elements?

$$S = 3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32, ...$$

#### Computing max

$$F(a_{[1..t]}) = \max_{i=1}^t a_i$$

Outputs are: 3, 3, 17, 17, 17, 32, 101, 101, ...

Just need to store **b** bits.

Median? A lot more tricky

# distinct elements? also tricky!

Ruta (UIUC) CS473 8 Spring 2021 8 / 3:

# Streaming Algorithms: Framework

```
\begin{split} &\langle \mathsf{Initialize\ summary\ information} \rangle \\ &\mathsf{While\ stream}\ \textit{\textbf{S}}\ \mathsf{is\ not\ done} \\ & \quad \textit{\textbf{x}} \leftarrow \mathsf{next\ element\ in}\ \textit{\textbf{S}} \\ & \quad \langle \mathsf{Do\ something\ with}\ \textit{\textbf{x}}\ \mathsf{and\ update\ summary\ information} \rangle \\ & \quad \langle \mathsf{Output\ something\ if\ needed} \rangle \end{split}
```

# Streaming Algorithms: Framework

```
(Initialize summary information)
```

While stream **S** is not done

 $x \leftarrow$  next element in S

⟨Do something with *x* and update summary information⟩ ⟨Output something if needed⟩

Return (summary)

Despite of restrictions, we can compute interesting functions if we can tolerate some error.

# Streaming Algorithms: One-sided Error

#### No false negative

Anything that needs to be considered/counted should be counted.

#### There may be false positive

We may over count. That is we may consider/count something that shouldn't have been counted.

### Part I

# Heavy Hitters

Find the element that occur strictly more than half the time, if any.

Note that at most one such element!

Find the element that occur strictly more than half the time, if any.

Note that at most one such element!

$$E, D, B, D, D_5, D, B, B, B, B, B, B_{11}, E, E, E, E, E_{16}$$

- At time **5**, it is **D**.
- At time **11**, it is **B**
- At time 16, none!

Finding a Majority Element

#### Treasure hunt

Once upon a time...

Finding a Majority Element

#### Treasure hunt

Once upon a time... there was a treasure hidden in a cave that different gangs were after. Only one-on-one fight is the unsaid rule (wild west style). Thus, if two members from different gangs face each other, then they shoot each other and both die.

Finding a Majority Element

#### Treasure hunt

Once upon a time... there was a treasure hidden in a cave that different gangs were after. Only one-on-one fight is the unsaid rule (wild west style). Thus, if two members from different gangs face each other, then they shoot each other and both die.

Which gang will get the treasure?

Suppose more than half the bandits are part of gang ALGO, then?

Finding a Majority Element

#### Treasure hunt

Once upon a time... there was a treasure hidden in a cave that different gangs were after. Only one-on-one fight is the unsaid rule (wild west style). Thus, if two members from different gangs face each other, then they shoot each other and both die.

Which gang will get the treasure?

Suppose more than half the bandits are part of gang ALGO, then?

Gang ALGO will get the treasure for sure!

Find the element that accrue strictly more than half the time, if any.

### R. Boyer and J. S. Moore Algorithm

Initialize: mem=∅ and counter=0

Find the element that accrue strictly more than half the time, if any.

### R. Boyer and J. S. Moore Algorithm

```
Initialize: mem=0 and counter=0
```

```
When element a_t arrives if (counter -- 0)
```

if (counter 
$$== 0$$
)

set mem= $a_t$  and counter=1

Find the element that accrue strictly more than half the time, if any.

### R. Boyer and J. S. Moore Algorithm

Initialize: mem= $\emptyset$  and counter=0

```
When element a_t arrives if (counter == 0) set mem=a_t and counter=1 else if (a_t == mem) then counter++
```

Find the element that accrue strictly more than half the time, if any.

### R. Boyer and J. S. Moore Algorithm

```
Initialize: mem=\emptyset and counter=0

When element a_t arrives

if (counter == 0)

set mem=a_t and counter=1

else if (a_t == mem) then counter++

else counter-- (discard a_t and a copy of mem)

Return mem.
```

Find the element that accrue strictly more than half the time, if any.

### R. Boyer and J. S. Moore Algorithm

```
Initialize: mem=\emptyset and counter=0

When element a_t arrives

if (counter == 0)

set mem=a_t and counter=1

else if (a_t == mem) then counter++

else counter-- (discard a_t and a copy of mem)

Return mem.
```

Even if no majority element, something is returned – False positive.

### R. Boyer and J. S. Moore Algorithm

Initialize: mem= $\emptyset$  and counter=0

```
When element a_t arrives if (counter == 0) set mem=a_t and counter=1 else if (a_t == mem) then counter++ else counter-- (discard a_t and a copy of mem) Return mem.
```

$$E, D, B, D, D_5, D, B, B, B, B, B, B_{11}, E, E, E, E, E_{16}$$

a <sub>t</sub>	Е	D	В	D	D	D	В	В	В	В	В	• • •
mem	Е	Е	В	В	D	D	D	D	В	В	В	• • •
counter	1	0	1	0	1	2	1	0	1	2	3	

Correctness, if majority element

#### Lemma

If there is a majority element, the algorithm will output it.

#### Proof.

 Decreasing counter is like throwing away a copy of element in mem.

Correctness, if majority element

#### Lemma

If there is a majority element, the algorithm will output it.

#### Proof.

- Decreasing counter is like throwing away a copy of element in mem.
- We do this every time  $a_t$  is different than mem, and there are less than half such  $a_t$ .

Correctness, if majority element

#### Lemma

If there is a majority element, the algorithm will output it.

#### Proof.

- Decreasing counter is like throwing away a copy of element in mem.
- We do this every time  $a_t$  is different than mem, and there are less than half such  $a_t$ .
- Sometimes mem may not contain the majority element.

Correctness, if majority element

#### Lemma

If there is a majority element, the algorithm will output it.

#### Proof.

- Decreasing counter is like throwing away a copy of element in mem.
- We do this every time  $a_t$  is different than mem, and there are less than half such  $a_t$ .
- Sometimes mem may not contain the majority element.
   However, even if we are throwing away the majority element every time, since they are more than half all can't be thrown.

Correctness, if majority element

#### Lemma

If there is a majority element, the algorithm will output it.

#### Proof.

- Decreasing counter is like throwing away a copy of element in mem.
- We do this every time  $a_t$  is different than mem, and there are less than half such  $a_t$ .
- Sometimes mem may not contain the majority element.
   However, even if we are throwing away the majority element every time, since they are more than half all can't be thrown.

In fact at any time t, mem contains majority element of sub-stream  $a_{[1..t]}$ , if any.

# Part II

Heavy Hitters

#### Definition

Given a stream  $S = a_1, a_2, ...$ , define count of element e at any time t to be

$$count_t(e) = |\{i \leq t \mid a_i = e\}|$$

e is called  $\epsilon$ -heavy hitter at time t if count<sub>t</sub> $(e) > \epsilon t$ .

### Definition

Given a stream  $S = a_1, a_2, ...$ , define count of element e at any time t to be

$$count_t(e) = |\{i \leq t \mid a_i = e\}|$$

e is called  $\epsilon$ -heavy hitter at time t if count $_t(e) > \epsilon t$ .

### Goal:

Maintain a structure containing all the  $\epsilon$ -heavy hitters so far.

At any point there are at most  $1/\epsilon$  such elements.

#### Definition

Given a stream  $S = a_1, a_2, ...$ , define count of element e at any time t to be

$$count_t(e) = |\{i \le t \mid a_i = e\}|$$

e is called  $\epsilon$ -heavy hitter at time t if count $_t(e) > \epsilon t$ .

#### Goal:

Maintain a structure containing all the  $\epsilon$ -heavy hitters so far. At any point there are at most  $1/\epsilon$  such elements.

## Crucial Note: false positive are OK, but no false negative

We are NOT allowed to miss any heavy-hitters, but we could store non-heavy-hitters.

Ruta (UIUC) CS473 18 Spring 2021 18 / 33

If  $\epsilon = 1/2$  then the majority element!

Ruta (UIUC) CS473 19 Spring 2021 19 / 33

If  $\epsilon = 1/2$  then the majority element!

$$E, D, B, D, D_5, D, B, A, B, B, B_{11}, E, E, E, E, E_{16}$$

### 1/3-heavy hitters

- At time **5**, it is **D**.
- At time 11, both B and D.
- At time 15,

If  $\epsilon = 1/2$  then the majority element!

$$E, D, B, D, D_5, D, B, A, B, B, B_{11}, E, E, E, E, E_{16}$$

#### 1/3-heavy hitters

- At time **5**, it is **D**.
- At time 11, both B and D.
- At time 15, none!
- At time 16,

If  $\epsilon = 1/2$  then the majority element!

$$E, D, B, D, D_5, D, B, A, B, B, B_{11}, E, E, E, E, E_{16}$$

### 1/3-heavy hitters

- At time **5**, it is **D**.
- At time 11, both B and D.
- At time 15, none!
- At time **16**, it is **E**.

If  $\epsilon = 1/2$  then the majority element!

$$E, D, B, D, D_5, D, B, A, B, B, B_{11}, E, E, E, E, E_{16}$$

### 1/3-heavy hitters

- At time **5**, it is **D**.
- At time 11, both B and D.
- At time 15, none!
- At time 16, it is *E*.

As time passes, the set of heavy hitters may change completely.

If 
$$\epsilon=1/2$$
 then the majority element! Set  $k=\lceil 1/\epsilon \rceil-1$ . (if  $\epsilon=1/2$  then  $k=1$ )

### Algorithm

Keep an array T[1, ..., k] to hold elements Keep an array C[1, ..., k] to hold their counters

Initialize: C[j] = 0 and  $T[j] = \emptyset$  for all i.

If 
$$\epsilon=1/2$$
 then the majority element! Set  $k=\lceil 1/\epsilon \rceil-1$ . (if  $\epsilon=1/2$  then  $k=1$ )

### Algorithm

```
Keep an array T[1,\ldots,k] to hold elements
Keep an array C[1,\ldots,k] to hold their counters
Initialize: C[j]=0 and T[j]=\emptyset for all i.
When element a_t arrives,
If (a_t==T[j] for some j\leq k), then C[j]++.
```

If 
$$\epsilon=1/2$$
 then the majority element! Set  $k=\lceil 1/\epsilon \rceil-1$ . (if  $\epsilon=1/2$  then  $k=1$ )

### Algorithm

```
Keep an array T[1, \ldots, k] to hold elements
Keep an array C[1, \ldots, k] to hold their counters
Initialize: C[j] = 0 and T[j] = \emptyset for all i.
When element a_t arrives,
If (a_t == T[j] for some j \leq k), then C[j] + +.
Else if (C[j] == 0 for some j \leq k), then
```

If 
$$\epsilon=1/2$$
 then the majority element! Set  $k=\lceil 1/\epsilon \rceil-1$ . (if  $\epsilon=1/2$  then  $k=1$ )

## Algorithm

```
Keep an array T[1,\ldots,k] to hold elements
Keep an array C[1,\ldots,k] to hold their counters
Initialize: C[j]=0 and T[j]=\emptyset for all i.
When element a_t arrives,
If (a_t==T[j] for some j\leq k), then C[j]++.
Else if (C[j]==0 for some j\leq k), then
Set T[j]\leftarrow a_t and C[j]\leftarrow 1.
```

```
If \epsilon = 1/2 then the majority element!
Set k = \lceil 1/\epsilon \rceil - 1. (if \epsilon = 1/2 then k = 1)
```

### Algorithm

```
Keep an array T[1, \ldots, k] to hold elements
Keep an array C[1, \ldots, k] to hold their counters
Initialize: C[i] = 0 and T[i] = \emptyset for all i.
When element a_t arrives,
   If (a_t == T[j] \text{ for some } j < k), then C[j] + +.
    Else if (C[i] == 0 for some i < k, then
     Set T[i] \leftarrow a_t and C[i] \leftarrow 1.
   Else do C[j] — for all j. (discard a_t and a copy of all T[j])
```

```
If \epsilon = 1/2 then the majority element!
Set k = \lceil 1/\epsilon \rceil - 1. (if \epsilon = 1/2 then k = 1)
```

### Algorithm

```
Keep an array T[1, \ldots, k] to hold elements
Keep an array C[1, \ldots, k] to hold their counters
Initialize: C[i] = 0 and T[i] = \emptyset for all i.
When element a_t arrives,
   If (a_t == T[j] \text{ for some } j < k), then C[j] + +.
    Else if (C[i] == 0 for some i < k, then
     Set T[i] \leftarrow a_t and C[i] \leftarrow 1.
   Else do C[j] — for all j. (discard a_t and a copy of all T[j])
```

Same as the *Majority algorithm* for  $\epsilon = 1/2$ .

#### Algorithm Analysis

At any time t, our estimates are:

$$est_t(e) = C[j]$$
 if  $e = T[j]$   
= 0 otherwise

#### Lemma

Estimates satisfy:  $est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$ 

#### Algorithm Analysis

At any time t, our estimates are:

$$est_t(e) = C[j]$$
 if  $e = T[j]$   
= 0 otherwise

#### Lemma

Estimates satisfy: 
$$est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$$

For each element, count is maintained up to  $\epsilon t$  error!

#### Algorithm Analysis

At any time t, our estimates are:

$$est_t(e) = C[j]$$
 if  $e = T[j]$   
= 0 otherwise

#### Lemma

Estimates satisfy: 
$$est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$$

For each element, count is maintained up to  $\epsilon t$  error!

If e is not an  $\epsilon$ -heavy hitter then  $\operatorname{count}_t(e) \leq \epsilon t$ , and hence  $\operatorname{est}_t(e) = 0$  is correct up to  $\epsilon t$  error.

#### Algorithm Analysis

At any time t, our estimates are:

$$est_t(e) = C[j]$$
 if  $e = T[j]$   
= 0 otherwise

#### Lemma

Estimates satisfy:  $est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$ 

### Corollary

For any time t, T contains all the  $\epsilon$ -heavy hitters in  $a_{[1..t]}$ .

### Proof.

If e is a heavy hitter at time t then  $count_t(e) > \epsilon t$ .

#### Algorithm Analysis

At any time t, our estimates are:

$$est_t(e) = C[j]$$
 if  $e = T[j]$   
= 0 otherwise

#### Lemma

Estimates satisfy:  $est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$ 

### Corollary

For any time t, T contains all the  $\epsilon$ -heavy hitters in  $a_{[1..t]}$ .

### Proof.

If e is a heavy hitter at time t then  $\operatorname{count}_t(e) > \epsilon t$ . Using the lemma,

$$\operatorname{est}_t(e) \geq \operatorname{count}_t(e) - \epsilon t$$

#### Algorithm Analysis

At any time t, our estimates are:

$$est_t(e) = C[j]$$
 if  $e = T[j]$   
= 0 otherwise

#### Lemma

Estimates satisfy:  $est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$ 

### Corollary

For any time t, T contains all the  $\epsilon$ -heavy hitters in  $a_{[1..t]}$ .

#### Proof.

If e is a heavy hitter at time t then  $\operatorname{count}_t(e) > \epsilon t$ . Using the lemma,

$$\operatorname{est}_t(e) \geq \operatorname{count}_t(e) - \epsilon t > 0$$

#### Algorithm Analysis

At any time t, our estimates are:

$$est_t(e) = C[j]$$
 if  $e = T[j]$   
= 0 otherwise

#### Lemma

Estimates satisfy:  $est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$ 

#### Proof.

Counter for e increases only when we see e,  $\therefore$  est<sub>t</sub> $(e) \leq \text{count}_t(e)$ .

#### Algorithm Analysis

At any time t, our estimates are:

$$est_t(e) = C[j]$$
 if  $e = T[j]$   
= 0 otherwise

#### Lemma

Estimates satisfy:  $est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$ 

#### Proof.

Counter for e increases only when we see e,  $\therefore$  est<sub>t</sub> $(e) \le \text{count}_t(e)$ . We want count<sub>t</sub> $(e) - \text{est}_t(e) \le \epsilon t$ . It decreases by one,

ullet when we decrease all k counters, and see an element outside T

#### Algorithm Analysis

At any time t, our estimates are:

$$est_t(e) = C[j]$$
 if  $e = T[j]$   
= 0 otherwise

#### Lemma

Estimates satisfy:  $est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$ 

#### Proof.

Counter for e increases only when we see e,  $\therefore$  est<sub>t</sub> $(e) \le \text{count}_t(e)$ . We want count<sub>t</sub> $(e) - \text{est}_t(e) \le \epsilon t$ . It decreases by one,

- ullet when we decrease all k counters, and see an element outside T
- this is like discarding k + 1 elements.
- up to time t, we have only t elements to discard

#### Algorithm Analysis

At any time t, our estimates are:

$$est_t(e) = C[j]$$
 if  $e = T[j]$   
= 0 otherwise

#### Lemma

Estimates satisfy:  $est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$ 

#### Proof.

Counter for e increases only when we see e,  $\therefore$  est<sub>t</sub>(e)  $\leq$  count<sub>t</sub>(e).

- We want  $\operatorname{count}_t(e) \operatorname{est}_t(e) \le \epsilon t$ . It decreases by one,
  - ullet when we decrease all k counters, and see an element outside T
  - this is like discarding k + 1 elements.
  - up to time t, we have only t elements to discard

So at most  $t/(k+1) < t\epsilon$  such decreases.

Space usage

Set 
$$k = \lceil 1/\epsilon \rceil - 1$$
. (if  $\epsilon = 1/2$  then  $k = 1$ )

## Algorithm

Keep an array T[1, ..., k] to hold elements Keep an array C[1, ..., k] to hold their counters

÷

Maintains  $O(1/\epsilon)$  counters and elements.

Space usage

Set 
$$k = \lceil 1/\epsilon \rceil - 1$$
. (if  $\epsilon = 1/2$  then  $k = 1$ )

### Algorithm

Keep an array T[1, ..., k] to hold elements Keep an array C[1, ..., k] to hold their counters

:

Maintains  $O(1/\epsilon)$  counters and elements.  $O(\log t)$  for each counter.  $O(\Sigma)$  for each element, where  $\Sigma$  is the description of largest element.

Space usage

Set 
$$k = \lceil 1/\epsilon \rceil - 1$$
. (if  $\epsilon = 1/2$  then  $k = 1$ )

### Algorithm

Keep an array T[1, ..., k] to hold elements Keep an array C[1, ..., k] to hold their counters

:

Maintains  $O(1/\epsilon)$  counters and elements.

 $O(\log t)$  for each counter.  $O(\Sigma)$  for each element, where  $\Sigma$  is the description of largest element.

Total:  $O(1/\epsilon(\log t + \Sigma))$ .

Recall: maintains counts for all elements up to  $\epsilon t$  error.

# Part III

# Use of Hash Functions

### **Problem Statement:**

At any time t, estimate the number of times every element appeared so far.

### **Problem Statement:**

At any time t, estimate the number of times every element appeared so far.

If error up to  $\epsilon t$  is OK, then we can use  $\epsilon$ -heavy hitter algorithm.

#### **Problem Statement:**

At any time t, estimate the number of times every element appeared so far.

If error up to  $\epsilon t$  is OK, then we can use  $\epsilon$ -heavy hitter algorithm.

It takes  $O(1/\epsilon(\log t + \Sigma))$  space.

#### **Problem Statement:**

At any time t, estimate the number of times every element appeared so far.

If error up to  $\epsilon t$  is OK, then we can use  $\epsilon$ -heavy hitter algorithm.

It takes  $O(1/\epsilon(\log t + \Sigma))$  space.

Can we do better?

Yes – Bloom filter like idea

### Recall: Bloom Filter

#### Storage for inserts and lookups

Sample hash functions  $h_1, \ldots, h_d$  independently and uniformly at random from some family  $\mathcal{H}$ .

```
egin{aligned} \mathsf{Insert}(e) \ \mathsf{For} \ i = 1...d \ \mathsf{Set} \ T_i[h_i(e)] \leftarrow 1 \end{aligned}
```

```
Lookup(e)

For i = 1...d

If (T_i[h_i(e)] == 0) then return "No"

Return "Yes"
```

If e inserted, then Lookup(e) will always return "Yes".

### Recall: Bloom Filter

#### Storage for inserts and lookups

Sample hash functions  $h_1, \ldots, h_d$  independently and uniformly at random from some family  $\mathcal{H}$ .

```
egin{aligned} \mathsf{Insert}(e) \ \mathsf{For} \ i = 1...d \ \mathsf{Set} \ \mathcal{T}_i[h_i(e)] \leftarrow 1 \end{aligned}
```

```
Lookup(e)
For i = 1...d
If (T_i[h_i(e)] == 0) then return "No"
Return "Yes"
```

If e inserted, then Lookup(e) will always return "Yes".

- e not inserted, but still it can return "Yes" with very low probability.
  - Due to some e's being inserted with  $h_i(e') = h_i(e)$ .
  - If  $\Pr_{h_i \sim \mathcal{H}}[e \text{ not inserted and } T_i[h_i(e)] = 1] \leq \alpha$ , then combined error probability would be at most  $\alpha^d$ .

### Count Min-Sketch

By G. Cormode and S. M. Muthukrishnan'05

Keep d arrays  $C_1, ..., C_d$ , each to hold m counters.

By G. Cormode and S. M. Muthukrishnan'05

Keep d arrays  $C_1, ..., C_d$ , each to hold m counters.

 $\mathcal{H}$ : 2-universal family of hash functions mapping U to  $\{0, \ldots, m-1\}$ . Sample  $h_1, \ldots, h_d$  independently and uniformly at random from  $\mathcal{H}$ .

Ruta (UIUC) CS473 28 Spring 2021 28 / 33

By G. Cormode and S. M. Muthukrishnan'05

Keep d arrays  $C_1, ..., C_d$ , each to hold m counters.

 $\mathcal{H}$ : 2-universal family of hash functions mapping U to  $\{0,\ldots,m-1\}$ . Sample  $h_1,\ldots,h_d$  independently and uniformly at random from  $\mathcal{H}$ .

```
CMInsert(e)
For i = 1...d
Do C_i[h_i(e)] + +
```

By G. Cormode and S. M. Muthukrishnan'05

Keep d arrays  $C_1, ..., C_d$ , each to hold m counters.

 $\mathcal{H}$ : 2-universal family of hash functions mapping U to  $\{0,\ldots,m-1\}$ . Sample  $h_1,\ldots,h_d$  independently and uniformly at random from  $\mathcal{H}$ .

CMInsert(e)  
For 
$$i = 1...d$$
  
Do  $C_i[h_i(e)] + +$ 

```
\begin{aligned} &\mathsf{CMEstimate}(e) \\ & est \leftarrow \infty \\ &\mathsf{For} \ i = 1...d \\ & est \leftarrow \min\{est, \mathit{C}_i[h_i(e)]\} \\ & \mathsf{Return} \ est \end{aligned}
```

By G. Cormode and S. M. Muthukrishnan'05

Keep d arrays  $C_1, ..., C_d$ , each to hold m counters.

 $\mathcal{H}$ : 2-universal family of hash functions mapping U to  $\{0,\ldots,m-1\}$ . Sample  $h_1,\ldots,h_d$  independently and uniformly at random from  $\mathcal{H}$ .

```
CMInsert(e)
For i = 1...d
Do C_i[h_i(e)] + +
```

```
 \begin{aligned} \mathsf{CMEstimate}(e) \\ est &\leftarrow \infty \\ \mathsf{For} \ i &= 1...d \\ est &\leftarrow \mathsf{min}\{est, \textit{C}_i[\textit{h}_i(e)]\} \\ \mathsf{Return} \ est \end{aligned}
```

As element  $a_t$  arrives at time t, call CMInsert $(a_t)$ .

To get count of e at any time t, call CMEstimate(e).

Ruta (UIUC) CS473 28 Spring 2021 28 / 33

By G. Cormode and S. M. Muthukrishnan'05

```
CMInsert(e)

For i = 1...d

Do C_i[h_i(e)] + +
```

```
CMEstimate(e)
est \leftarrow \infty
For i = 1...d
est \leftarrow \min\{est, C_i[h_i(e)]\}
Return est
```

At time t, let  $est_t(e) = CMEstimate(e) = min_{i=1}^d C_i[h_i(e)]$ .

Ruta (UIUC) CS473 29 Spring 2021 29 / 33

By G. Cormode and S. M. Muthukrishnan'05

```
CMInsert(e)

For i = 1...d

Do C_i[h_i(e)] + +
```

```
 \begin{aligned} \mathsf{CMEstimate}(e) \\ est &\leftarrow \infty \\ \mathsf{For} \ i = 1...d \\ est &\leftarrow \mathsf{min}\{est, \mathit{C}_i[\mathit{h}_i(e)]\} \\ \mathsf{Return} \ est \end{aligned}
```

At time t, let  $\operatorname{est}_t(e) = \operatorname{CMEstimate}(e) = \min_{i=1}^d C_i[h_i(e)]$ . Observation:  $\operatorname{est}_t(e) \geq \operatorname{count}_t(e)$ .

Ruta (UIUC) CS473 29 Spring 2021 29 / 33

By G. Cormode and S. M. Muthukrishnan'05

```
CMInsert(e)
For i = 1...d
Do C_i[h_i(e)] + +
```

```
\begin{aligned} &\mathsf{CMEstimate}(e) \\ & \textit{est} \leftarrow \infty \\ & \mathsf{For} \ i = 1...d \\ & \textit{est} \leftarrow \mathsf{min}\{\textit{est}, \textit{C}_i[\textit{h}_i(e)]\} \\ & \mathsf{Return} \ \textit{est} \end{aligned}
```

```
At time t, let \operatorname{est}_t(e) = \operatorname{CMEstimate}(e) = \min_{i=1}^d C_i[h_i(e)]. Observation: \operatorname{est}_t(e) \geq \operatorname{count}_t(e).
```

Question: How big  $(est_t(e) - count_t(e))$  can be?

By G. Cormode and S. M. Muthukrishnan'05

```
CMInsert(e)
For i = 1...d
Do C_i[h_i(e)] + +
```

```
\begin{aligned} &\mathsf{CMEstimate}(e) \\ & \textit{est} \leftarrow \infty \\ & \mathsf{For} \ i = 1...d \\ & \textit{est} \leftarrow \mathsf{min}\{\textit{est}, \textit{C}_i[\textit{h}_i(e)]\} \\ & \mathsf{Return} \ \textit{est} \end{aligned}
```

```
At time t, let \operatorname{est}_t(e) = \operatorname{CMEstimate}(e) = \min_{i=1}^d C_i[h_i(e)]. Observation: \operatorname{est}_t(e) \geq \operatorname{count}_t(e).
```

Question: How big  $(est_t(e) - count_t(e))$  can be?

**Recall:** Any  $e, y \in U$ , if  $e \neq y$  then  $\Pr[h_i(y) = h_i(e)] = \frac{1}{m} \forall i$ .

Ruta (UIUC) CS473 29 Spring 2021 29 / 33

By G. Cormode and S. M. Muthukrishnan'05

Let  $f'_x = \operatorname{est}_t(x)$  and  $f_x = \operatorname{count}_t(x)$ .

Fix an element e. We want to bound  $(f'_e - f_e)$ .

Observations:

By G. Cormode and S. M. Muthukrishnan'05

Let 
$$f'_x = \operatorname{est}_t(x)$$
 and  $f_x = \operatorname{count}_t(x)$ .

Fix an element e. We want to bound  $(f'_e - f_e)$ .

#### Observations:

Define indicator variable  $X_{i,x} = [h_i(x) = h_i(e)]$ .

$$E[X_{i,x}] = Pr[h_i(x) = h_i(e)] = 1/m$$

By G. Cormode and S. M. Muthukrishnan'05

Let 
$$f'_x = \operatorname{est}_t(x)$$
 and  $f_x = \operatorname{count}_t(x)$ .

Fix an element e. We want to bound  $(f'_e - f_e)$ .

#### Observations:

Define indicator variable  $X_{i,x} = [h_i(x) = h_i(e)]$ .

$$E[X_{i,x}] = Pr[h_i(x) = h_i(e)] = 1/m$$

Let  $X_i := \sum_{x \neq e} X_{i,x} f_x$  be the total over counting at  $C_i[h_i(e)]$ .

$$C_i[h_i(e)] = X_i + f_e$$

By G. Cormode and S. M. Muthukrishnan'05

Let 
$$f'_x = \operatorname{est}_t(x)$$
 and  $f_x = \operatorname{count}_t(x)$ .

Fix an element e. We want to bound  $(f'_e - f_e)$ .

#### Observations:

Define indicator variable  $X_{i,x} = [h_i(x) = h_i(e)]$ .

$$\mathsf{E}[X_{i,x}] = \mathsf{Pr}[h_i(x) = h_i(e)] = 1/m$$

Let  $X_i := \sum_{x \neq e} X_{i,x} f_x$  be the total over counting at  $C_i[h_i(e)]$ .

$$C_i[h_i(e)] = X_i + f_e$$

and since at most t elements have arrived so far,

$$\mathsf{E}[X_i] = \sum_{x \neq e} \mathsf{E}[X_{i,x}] \, f_x = \frac{1}{m} \sum_{x \neq e} f_x \le$$

By G. Cormode and S. M. Muthukrishnan'05

Let 
$$f'_x = \operatorname{est}_t(x)$$
 and  $f_x = \operatorname{count}_t(x)$ .

Fix an element e. We want to bound  $(f'_e - f_e)$ .

#### Observations:

Define indicator variable  $X_{i,x} = [h_i(x) = h_i(e)]$ .

$$\mathsf{E}[X_{i,x}] = \mathsf{Pr}[h_i(x) = h_i(e)] = 1/m$$

Let  $X_i := \sum_{x \neq e} X_{i,x} f_x$  be the total over counting at  $C_i[h_i(e)]$ .

$$C_i[h_i(e)] = X_i + f_e$$

and since at most t elements have arrived so far,

$$\mathsf{E}[X_i] = \sum_{\mathsf{x} \neq e} \mathsf{E}[X_{i,\mathsf{x}}] \, f_\mathsf{x} = \frac{1}{m} \sum_{\mathsf{x} \neq e} f_\mathsf{x} \le \frac{t}{m}$$

By G. Cormode and S. M. Muthukrishnan'05

We have  $C_i[h_i(e)] = X_i + f_e$  and  $E[X_i] \leq \frac{t}{m}$ .

By G. Cormode and S. M. Muthukrishnan'05

We have 
$$C_i[h_i(e)] = X_i + f_e$$
 and  $E[X_i] \leq \frac{t}{m}$ .

Then, for  $\epsilon > 0$ 

$$\Pr[C_i[h_i(e)] - f_e \ge \epsilon t] = \Pr[X_i \ge \epsilon t]$$

[definition]

By G. Cormode and S. M. Muthukrishnan'05

We have 
$$C_i[h_i(e)] = X_i + f_e$$
 and  $\mathbf{E}[X_i] \leq \frac{t}{m}$ .

Then, for  $\epsilon > 0$ 

$$\Pr[C_i[h_i(e)] - f_e \ge \epsilon t] = \Pr[X_i \ge \epsilon t]$$
 [definition]   
  $\le \frac{\mathbb{E}[X_i]}{\epsilon t}$  [Markov's inequality]

By G. Cormode and S. M. Muthukrishnan'05

We have 
$$C_i[h_i(e)] = X_i + f_e$$
 and  $\mathbf{E}[X_i] \leq \frac{t}{m}$ .

Then, for  $\epsilon > 0$ 

$$\begin{array}{ll} \Pr[C_i[h_i(e)] - f_e \geq \epsilon t] &=& \Pr[X_i \geq \epsilon t] & \text{[definition]} \\ &\leq \frac{\underline{t}[X_i]}{\epsilon t} & \text{[Markov's inequality]} \\ &\leq \frac{t/m}{\epsilon t} = \frac{1}{m\epsilon} & \text{[derived above]} \end{array}$$

By G. Cormode and S. M. Muthukrishnan'05

We have 
$$C_i[h_i(e)] = X_i + f_e$$
 and  $E[X_i] \leq \frac{t}{m}$ .

Then, for  $\epsilon > 0$ 

$$\begin{array}{ll} \Pr[C_i[h_i(e)] - f_e \geq \epsilon t] &=& \Pr[X_i \geq \epsilon t] & \text{[definition]} \\ &\leq & \frac{\underline{\epsilon}[X_i]}{\epsilon t} & \text{[Markov's inequality]} \\ &\leq & \frac{t/m}{\epsilon t} = \frac{1}{m\epsilon} & \text{[derived above]} \end{array}$$

Recall: 
$$f'_e = \operatorname{est}_t(e) = \min_{i=1}^d C_i[h_i(e)].$$

By G. Cormode and S. M. Muthukrishnan'05

We have 
$$C_i[h_i(e)] = X_i + f_e$$
 and  $E[X_i] \leq \frac{t}{m}$ .

Then, for  $\epsilon > 0$ 

$$\begin{array}{ll} \Pr[C_i[h_i(e)] - f_e \geq \epsilon t] &=& \Pr[X_i \geq \epsilon t] & \text{[definition]} \\ &\leq & \frac{\underline{e}[X_i]}{\epsilon t} & \text{[Markov's inequality]} \\ &\leq & \frac{t/m}{\epsilon t} = \frac{1}{m\epsilon} & \text{[derived above]} \end{array}$$

Recall: 
$$f'_e = \operatorname{est}_t(e) = \min_{i=1}^d C_i[h_i(e)].$$

$$\Pr[f'_e - f_e \ge \epsilon t] =$$

By G. Cormode and S. M. Muthukrishnan'05

We have 
$$C_i[h_i(e)] = X_i + f_e$$
 and  $E[X_i] \leq \frac{t}{m}$ .

Then, for  $\epsilon > 0$ 

$$\begin{array}{ll} \Pr[\textit{C}_i[\textit{h}_i(e)] - \textit{f}_e \geq \epsilon t] &=& \Pr[\textit{X}_i \geq \epsilon t] & \text{[definition]} \\ &\leq& \frac{\underline{\mathsf{E}}[\textit{X}_i]}{\epsilon t} & \text{[Markov's inequality]} \\ &\leq& \frac{t/m}{\epsilon t} = \frac{1}{m\epsilon} & \text{[derived above]} \end{array}$$

Recall: 
$$f'_e = \operatorname{est}_t(e) = \min_{i=1}^d C_i[h_i(e)].$$

$$\Pr[f'_e - f_e \ge \epsilon t] = \Pr[C_i[h_i(e)] - f_e \ge \epsilon t \text{ for all } i]$$

By G. Cormode and S. M. Muthukrishnan'05

We have 
$$C_i[h_i(e)] = X_i + f_e$$
 and  $E[X_i] \leq \frac{t}{m}$ .

Then, for  $\epsilon > 0$ 

$$\begin{array}{ll} \Pr[\textit{C}_i[\textit{h}_i(e)] - \textit{f}_e \geq \epsilon t] & = \Pr[\textit{X}_i \geq \epsilon t] & \text{[definition]} \\ \leq \frac{\underline{\textbf{f}}[\textit{X}_i]}{\epsilon t} & \text{[Markov's inequality]} \\ \leq \frac{t/m}{\epsilon t} = \frac{1}{m\epsilon} & \text{[derived above]} \end{array}$$

Recall: 
$$f'_e = \operatorname{est}_t(e) = \min_{i=1}^d C_i[h_i(e)].$$

$$\Pr[f'_e - f_e \ge \epsilon t] = \Pr[C_i[h_i(e)] - f_e \ge \epsilon t \text{ for all } i]$$
  
= 
$$\Pr[X_i \ge \epsilon t \text{ for all } i]$$

By G. Cormode and S. M. Muthukrishnan'05

We have 
$$C_i[h_i(e)] = X_i + f_e$$
 and  $E[X_i] \leq \frac{t}{m}$ .

Then, for  $\epsilon > 0$ 

$$\begin{array}{ll} \Pr[C_i[h_i(e)] - f_e \geq \epsilon t] &=& \Pr[X_i \geq \epsilon t] & \text{[definition]} \\ &\leq & \frac{\underline{\epsilon}[X_i]}{\epsilon t} & \text{[Markov's inequality]} \\ &\leq & \frac{t/m}{\epsilon t} = \frac{1}{m\epsilon} & \text{[derived above]} \end{array}$$

Recall: 
$$f'_e = \operatorname{est}_t(e) = \min_{i=1}^d C_i[h_i(e)].$$

$$\Pr[f'_e - f_e \ge \epsilon t] = \Pr[C_i[h_i(e)] - f_e \ge \epsilon t \text{ for all } i]$$

$$= \Pr[X_i \ge \epsilon t \text{ for all } i]$$

$$= \Pi^d_{i=1} \Pr[X_i \ge \epsilon t] \text{ [independence of } h_i \text{'s]}$$

By G. Cormode and S. M. Muthukrishnan'05

We have 
$$C_i[h_i(e)] = X_i + f_e$$
 and  $\mathbf{E}[X_i] \leq \frac{t}{m}$ .

Then, for  $\epsilon > 0$ 

$$\begin{array}{ll} \Pr[C_i[h_i(e)] - f_e \geq \epsilon t] &=& \Pr[X_i \geq \epsilon t] & \text{[definition]} \\ &\leq & \frac{\underline{t}[X_i]}{\epsilon t} & \text{[Markov's inequality]} \\ &\leq & \frac{t/m}{\epsilon t} = \frac{1}{m\epsilon} & \text{[derived above]} \end{array}$$

Recall: 
$$f'_e = \operatorname{est}_t(e) = \min_{i=1}^d C_i[h_i(e)].$$

$$\Pr[f'_e - f_e \ge \epsilon t] = \Pr[C_i[h_i(e)] - f_e \ge \epsilon t \text{ for all } i]$$

$$= \Pr[X_i \ge \epsilon t \text{ for all } i]$$

$$= \prod_{i=1}^d \Pr[X_i \ge \epsilon t] \text{ [independence of } h_i \text{'s]}$$

$$\le \left(\frac{1}{\epsilon m}\right)^d \text{ [derived above]}$$

By G. Cormode and S. M. Muthukrishnan'05

$$\mathsf{Pr}[\mathsf{est}_t(e) - \mathsf{count}_t(e) \geq \epsilon t] \leq \left( rac{1}{\epsilon m} 
ight)^d$$

By G. Cormode and S. M. Muthukrishnan'05

$$\Pr[\operatorname{est}_t(e) - \operatorname{count}_t(e) \geq \epsilon t] \leq \left(\frac{1}{\epsilon m}\right)^d \leq \delta$$

By G. Cormode and S. M. Muthukrishnan'05

$$\Pr[\operatorname{est}_t(e) - \operatorname{count}_t(e) \geq \epsilon t] \leq \left(rac{1}{\epsilon m}
ight)^d \leq \delta$$

Set  $m = \lceil 2/\epsilon \rceil$  and  $d = \lceil \lg 1/\delta \rceil$ .

By G. Cormode and S. M. Muthukrishnan'05

$$\Pr[\operatorname{est}_t(e) - \operatorname{count}_t(e) \ge \epsilon t] \le \left(\frac{1}{\epsilon m}\right)^d \le \delta$$

Set 
$$m = \lceil 2/\epsilon \rceil$$
 and  $d = \lceil \lg 1/\delta \rceil$ .

Space:

By G. Cormode and S. M. Muthukrishnan'05

$$\Pr[\operatorname{est}_t(e) - \operatorname{count}_t(e) \ge \epsilon t] \le \left(\frac{1}{\epsilon m}\right)^d \le \delta$$

Set  $m = \lceil 2/\epsilon \rceil$  and  $d = \lceil \lg 1/\delta \rceil$ .

Space: m \* d counters each of size  $\lg(t) = O(\frac{1}{\epsilon} \lg \frac{1}{\delta} \lg t)$  bits.

By G. Cormode and S. M. Muthukrishnan'05

#### Lemma

Given  $\epsilon, \delta > 0$ , we can estimate count<sub>t</sub>(e), at any time t for any element e, up to  $\epsilon t$  error with probability at least  $(1 - \delta)$  using  $O(\frac{1}{\epsilon} \lg \frac{1}{\delta})$  many counters.