CS 473: Algorithms

Chandra Chekuri Ruta Mehta

University of Illinois, Urbana-Champaign

Fall 2016

Streaming Algorithms

Lecture 12 October 5, 2016

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 $\mathbf{a_t}$ arrive at time \mathbf{t} . Let us assume that $\mathbf{a_t}$'s are numbers for this lecture.

Data Streams

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

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

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

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

$$\mathsf{F}(\mathsf{a}_{[1..t]}) = \textstyle \sum_{\mathsf{i}=1}^{\mathsf{t}} \mathsf{a}_{\mathsf{i}}$$

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

$$\mathsf{F}(\mathsf{a}_{[1..t]}) = \textstyle \sum_{\mathsf{i}=1}^{\mathsf{t}} \mathsf{a}_{\mathsf{i}}$$

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

$$\mathsf{F}(\mathsf{a}_{[1..t]}) = \mathsf{max}_{\mathsf{i}=1}^\mathsf{t} \, \mathsf{a}_{\mathsf{i}}$$

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

$$\mathsf{F}(\mathsf{a}_{[1..t]}) = \mathsf{max}_{\mathsf{i}=1}^\mathsf{t} \, \mathsf{a}_{\mathsf{i}}$$

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

$$\mathsf{F}(\mathsf{a}_{[1..t]}) = \mathsf{max}_{\mathsf{i}=1}^\mathsf{t} \, \mathsf{a}_{\mathsf{i}}$$

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

$$\mathsf{F}(\mathsf{a}_{[1..t]}) = \mathsf{max}_{\mathsf{i}=1}^\mathsf{t} \, \mathsf{a}_{\mathsf{i}}$$

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

$$\mathsf{F}(\mathsf{a}_{[1..t]}) = \mathsf{max}_{\mathsf{i}=1}^\mathsf{t} \, \mathsf{a}_{\mathsf{i}}$$

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!

Streaming Algorithms: Framework

```
⟨Initialize summary information⟩
While stream S is not done
    x ← next element in S
    ⟨Do something with x and update summary information⟩
    ⟨Output something if needed⟩
Return ⟨summary⟩
```

Streaming Algorithms: Framework

```
\langleInitialize summary information\rangle
```

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!

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

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

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

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 \mathbf{a_t} arrives
```

```
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=∅ and counter=0
```

```
When element \mathbf{a_t} arrives if (counter == 0) set mem=\mathbf{a_t} and counter=1 else if (\mathbf{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 \mathbf{a_t} arrives

if (counter == 0)

set mem=\mathbf{a_t} and counter=1

else if (\mathbf{a_t} == mem) then counter++

else counter-- (discard \mathbf{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 \mathbf{a_t} arrives

if (counter == 0)

set mem=\mathbf{a_t} and counter=1

else if (\mathbf{a_t} == mem) then counter++

else counter-- (discard \mathbf{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

R. Boyer and J. S. Moore Algorithm

```
When element \mathbf{a_t} arrives if (counter == 0)
```

Initialize: mem= \emptyset and counter=0

R. Boyer and J. S. Moore Algorithm

Initialize: mem= \emptyset and counter=0

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

R. Boyer and J. S. Moore Algorithm

Initialize: mem= \emptyset and counter=0

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

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

a _t	Е	D	В	D	D	D	В	В	В	В	В	• • •
mem	Е	Ε	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.
- Even if we are throwing away the majority element every time, since they are more than half all cannot be thrown away.



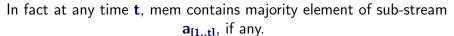
Correctness, if majority element

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.
- Even if we are throwing away the majority element every time, since they are more than half all cannot be thrown away.



15 / 32

Correctness, if majority element

Gang war interpretation!

Every element is a gang member. When we have two members from different gangs, they shoot each other.

Correctness, if majority element

Gang war interpretation!

Every element is a gang member. When we have two members from different gangs, they shoot each other.

If there is a gang with more than n/2 members, that will be the only one whose members will survive!

ϵ -Heavy Hitters

Definition

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

$$\mathsf{count}_t(e) = |\{i \le t \mid a_i = e\}|$$

It is called ϵ -heavy hitter at time **t** if count_t(**e**) $> \epsilon$ **t**.

Definition

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

$$\mathsf{count}_t(e) = |\{i \le t \mid a_i = e\}|$$

It is called ϵ -heavy hitter at time **t** if count_t(**e**) > ϵ **t**.

Goal:

Maintain a structure containing all the ϵ -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

$$\mathsf{count}_t(e) = |\{i \le t \mid a_i = e\}|$$

It is called ϵ -heavy hitter at time **t** if count_t(**e**) $> \epsilon$ **t**.

Goal:

Maintain a structure containing all the ϵ -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.

ϵ -Heavy Hitters: Example

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

ϵ -Heavy Hitters: Example

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

$$\mathsf{E}, \mathsf{D}, \mathsf{B}, \mathsf{D}, \mathsf{D}_5, \mathsf{D}, \mathsf{B}, \mathsf{B}, \mathsf{B}, \mathsf{B}, \mathsf{B}, \mathsf{B}_{11}, \mathsf{E}, \mathsf{E}, \mathsf{E}, \mathsf{E}, \mathsf{E}, \mathsf{E}_{16}$$

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 \mathbf{k}=\lceil 1/\epsilon \rceil-1. (if \epsilon=1/2 then \mathbf{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 \mathbf{k}=\lceil 1/\epsilon \rceil-1. (if \epsilon=1/2 then \mathbf{k}=1)
```

```
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.
When element a_t arrives,
If (a_t == T[j] for some j < k), then C[j] + +.
```

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

```
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 \mathbf{k}=\lceil 1/\epsilon \rceil-1. (if \epsilon=1/2 then \mathbf{k}=1)
```

```
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 \mathbf{k}=\lceil 1/\epsilon \rceil-1. (if \epsilon=1/2 then \mathbf{k}=1)
```

```
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.
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 \mathbf{k}=\lceil 1/\epsilon \rceil-1. (if \epsilon=1/2 then \mathbf{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.
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 ϵt 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 ϵt error!

If e is not an ϵ -heavy hitter then $\operatorname{count}_t(e) \leq \epsilon t$, and hence $\operatorname{est}_t(e) = 0$ is correct up to ϵt 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 ϵ -heavy hitters in $a_{[1..t]}$.

Proof.

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

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 ϵ -heavy hitters in $a_{[1..t]}$.

Proof.

If **e** is a heavy hitter at time **t** then $count_t(e) > \epsilon t$. Using the lemma,

$$\mathsf{est}_{\mathsf{t}}(\mathbf{e}) \geq \mathsf{count}_{\mathsf{t}}(\mathbf{e}) - \epsilon \mathbf{t}$$

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 ϵ -heavy hitters in $a_{[1..t]}$.

Proof.

If **e** is a heavy hitter at time **t** then $count_t(e) > \epsilon t$. Using the lemma,

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



Algorithm Analysis

At any time t, our estimates are:

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

Lemma

Estimates satisfy: $0 \le count_t(e) - est_t(e) \le \frac{t}{k+1} \le \epsilon t$

Proof.

Counter for **e** increases only when we see **e**, \therefore est_t(**e**) \leq 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:
$$0 \le count_t(e) - est_t(e) \le \frac{t}{k+1} \le \epsilon t$$

Proof.

Counter for **e** increases only when we see **e**, \therefore est_t(**e**) \leq count_t(**e**). count_t(**e**) — est_t(**e**) increases by one,

• 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:
$$0 \le count_t(e) - est_t(e) \le \frac{t}{k+1} \le \epsilon t$$

Proof.

Counter for e increases only when we see e, \therefore est_t $(e) \le count_t(e)$. count_t $(e) - est_t(e)$ increases by one,

- 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: $0 \le count_t(e) - est_t(e) \le \frac{t}{k+1} \le \epsilon t$

Proof.

Counter for **e** increases only when we see **e**, \therefore est_t(**e**) \leq count_t(**e**). count_t(**e**) — est_t(**e**) increases by one,

- 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) such increases.

Space usage

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

Algorithm

Keep an array $T[1, \ldots, k]$ to hold elements Keep an array $C[1, \ldots, k]$ to hold their counters

÷

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

Space usage

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

Algorithm

Keep an array $T[1, \ldots, k]$ to hold elements Keep an array $C[1, \ldots, k]$ to hold their counters

:

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

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

Space usage

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

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

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

Recall: maintains counts for all elements up to ϵt error.

Part II

Use of Hash Functions

Problem Statement:

At any time **t**, given an element **e**, estimate the number of times an **e** appeared so far.

Problem Statement:

At any time **t**, given an element **e**, estimate the number of times an **e** appeared so far.

If error up to ϵt is OK, then we can use ϵ -heavy hitter algorithm.

Problem Statement:

At any time **t**, given an element **e**, estimate the number of times an **e** appeared so far.

If error up to ϵt is OK, then we can use ϵ -heavy hitter algorithm.

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

Problem Statement:

At any time **t**, given an element **e**, estimate the number of times an **e** appeared so far.

If error up to ϵt is OK, then we can use ϵ -heavy hitter algorithm.

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

Can we do better?

Yes – Bloom filter like idea

Recall: Bloom Filter

Storage for inserts and lookups

Sample hash functions $\mathbf{h_1}, \dots, \mathbf{h_d}$ independently and uniformly at random.

```
\begin{aligned} & \text{Insert(e)} \\ & \text{For } \mathbf{i} = \mathbf{1}...\mathbf{d} \\ & \text{Set } \mathbf{T_i[h_i(e)]} \leftarrow \mathbf{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 $\mathbf{h_1}, \dots, \mathbf{h_d}$ independently and uniformly at random.

```
\begin{aligned} & \text{Insert(e)} \\ & \text{For i} = 1...d \\ & \text{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".

- 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 $Pr[e \text{ not inserted and } T_i[h_i(e)] = 1] \leq \alpha$, then combined error probability would be at most α^d .

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 $h: U \to \{0, \dots, m-1\}$. Sample h_1, \dots, h_d independently and uniformly at random from \mathcal{H} .

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 $h:U\to\{0,\ldots,m-1\}$. Sample h_1,\ldots,h_d independently and uniformly at random from \mathcal{H} .

```
\begin{aligned} &\mathsf{CMInsert}(e) \\ &\mathsf{For}\ \mathbf{i} = \mathbf{1}...\mathbf{d} \\ &\mathsf{Do}\ \mathbf{C_i[h_i(e)]} + + \end{aligned}
```

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

As element $\mathbf{a_t}$ arrives at time \mathbf{t} , call CMInsert($\mathbf{a_t}$).

To get count of \mathbf{e} at any time \mathbf{t} , call CMEstimate(\mathbf{e}).

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

```
\begin{aligned} &\mathsf{CMInsert}(e) \\ &\mathsf{For}\ \mathbf{i} = \mathbf{1}...\mathbf{d} \\ &\mathsf{Do}\ \mathbf{C_i}[\mathbf{h_i}(e)] + + \end{aligned}
```

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

```
At time \mathbf{t}, let \mathsf{est}_{\mathbf{t}}(\mathbf{e}) = \mathsf{CMEstimate}(\mathbf{e}). Observation: \mathsf{est}_{\mathbf{t}}(\mathbf{e}) \geq \mathsf{count}_{\mathbf{t}}(\mathbf{e}).
```

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

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

```
\begin{aligned} &\mathsf{CMInsert}(e) \\ &\mathsf{For}\ \mathbf{i} = \mathbf{1}...\mathbf{d} \\ &\mathsf{Do}\ \mathbf{C_i[h_i(e)]} + + \end{aligned}
```

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

```
At time \mathbf{t}, let \operatorname{est}_{\mathbf{t}}(\mathbf{e}) = \operatorname{CMEstimate}(\mathbf{e}).
Observation: \operatorname{est}_{\mathbf{t}}(\mathbf{e}) \geq \operatorname{count}_{\mathbf{t}}(\mathbf{e}).
```

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

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

Count Min-Sketch: Analysis

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

For simplicity let $f_e' = \mathsf{est}_t(e)$ and $f_e = \mathsf{count}_t(e)$. Bound $f_e' - f_e$.

Observations:

Count Min-Sketch: Analysis

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

For simplicity let $f_e' = \mathsf{est}_t(e)$ and $f_e = \mathsf{count}_t(e)$. Bound $f_e' - f_e$.

Observations:

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

$$\mathsf{E}[\mathsf{X}_{\mathsf{i},\mathsf{e},\mathsf{y}}] = \mathsf{Pr}[\mathsf{h}_\mathsf{i}(\mathsf{y}) = \mathsf{h}_\mathsf{i}(\mathsf{e})] = 1/\mathsf{m}$$

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

For simplicity let $f_e' = \mathsf{est}_t(e)$ and $f_e = \mathsf{count}_t(e)$. Bound $f_e' - f_e$.

Observations:

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

$$\mathsf{E}[\mathsf{X}_{\mathsf{i},\mathsf{e},\mathsf{y}}] = \mathsf{Pr}[\mathsf{h}_\mathsf{i}(\mathsf{y}) = \mathsf{h}_\mathsf{i}(\mathsf{e})] = 1/\mathsf{m}$$

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

$$C_i[h_i(e)] = X_{i,e} + f_e$$

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

For simplicity let $f_e' = \mathsf{est}_t(e)$ and $f_e = \mathsf{count}_t(e)$. Bound $f_e' - f_e$.

Observations:

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

$$\mathsf{E}[\mathsf{X}_{\mathsf{i},\mathsf{e},\mathsf{y}}] = \mathsf{Pr}[\mathsf{h}_\mathsf{i}(\mathsf{y}) = \mathsf{h}_\mathsf{i}(\mathsf{e})] = 1/\mathsf{m}$$

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

$$C_i[h_i(e)] = X_{i,e} + f_e$$

and since at most t elements have arrived so far,

$$E[X_{i,e}] = \sum_{y \neq e} E[X_{i,e,y}] \, f_y = \frac{1}{m} \sum_{y \neq e} f_y \leq \frac{t}{m}$$

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

$$egin{aligned} \mathsf{C_i}[\mathsf{h_i}(\mathsf{e})] &= \mathsf{X_{i,e}} + \mathsf{f_e} \text{ and } \mathsf{E}[\mathsf{X_{i,e}}] \leq \frac{\mathsf{t}}{\mathsf{m}}. \ \end{aligned}$$
 For $\epsilon > 0$

$$\text{Pr}[\textbf{C}_{i}[\textbf{h}_{i}(\textbf{e})] - \textbf{f}_{\textbf{e}} \geq \epsilon \textbf{t}] \ = \ \text{Pr}[\textbf{X}_{i,\textbf{e}} \geq \epsilon \textbf{t}]$$

[definition]

$$egin{aligned} \mathsf{C_i}[\mathsf{h_i}(\mathsf{e})] &= \mathsf{X_{i,e}} + \mathsf{f_e} \text{ and } \mathsf{E}[\mathsf{X_{i,e}}] \leq \frac{\mathsf{t}}{\mathsf{m}}. \ \end{aligned}$$
 For $\epsilon > 0$

$$\begin{array}{lll} \text{Pr}[C_i[h_i(e)] - f_e \geq \epsilon t] & = & \text{Pr}[X_{i,e} \geq \epsilon t] & \text{[definition]} \\ & \leq & \frac{\underline{t}[X_{i,e}]}{\epsilon t} & \text{[Markov's inequality]} \end{array}$$

$$egin{aligned} \mathsf{C_i}[\mathsf{h_i}(\mathsf{e})] &= \mathsf{X_{i,e}} + \mathsf{f_e} \text{ and } \mathsf{E}[\mathsf{X_{i,e}}] \leq \frac{\mathsf{t}}{\mathsf{m}}. \ \end{aligned}$$
 For $\epsilon > 0$

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

$$egin{aligned} \mathsf{C_i}[\mathsf{h_i}(\mathsf{e})] &= \mathsf{X_{i,e}} + \mathsf{f_e} \text{ and } \mathsf{E}[\mathsf{X_{i,e}}] \leq \frac{\mathsf{t}}{\mathsf{m}}. \ \end{aligned}$$
 For $\epsilon > 0$

$$\begin{array}{lll} \text{Pr}[\mathsf{C}_{\mathsf{i}}[\mathsf{h}_{\mathsf{i}}(\mathsf{e})] - \mathsf{f}_{\mathsf{e}} \geq \epsilon \mathsf{t}] & & [\text{definition}] \\ & \leq & \frac{\mathsf{g}[\mathsf{X}_{\mathsf{i},\mathsf{e}}]}{\epsilon \mathsf{t}} & & [\text{Markov's inequality}] \\ & \leq & \frac{\mathsf{t}/\mathsf{m}}{\epsilon \mathsf{t}} = \frac{1}{\mathsf{m}\epsilon} & & [\text{derived above}] \end{array}$$

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

$$egin{aligned} \mathbf{C_i}[\mathbf{h_i}(\mathbf{e})] &= \mathbf{X_{i,e}} + \mathbf{f_e} \text{ and } \mathbf{E}[\mathbf{X_{i,e}}] \leq \frac{t}{m}. \ \end{aligned}$$
 For $\epsilon > 0$

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

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

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

$$egin{aligned} \mathbf{C_i}[\mathbf{h_i}(\mathbf{e})] &= \mathbf{X_{i,e}} + \mathbf{f_e} \text{ and } \mathbf{E}[\mathbf{X_{i,e}}] \leq \frac{t}{m}. \ \end{aligned}$$
 For $\epsilon > 0$

$$\begin{array}{lll} \text{Pr}[\mathsf{C_i}[\mathsf{h_i}(\mathsf{e})] - \mathsf{f_e} \geq \epsilon \mathsf{t}] & = & \text{Pr}[\mathsf{X_{i,e}} \geq \epsilon \mathsf{t}] & \text{[definition]} \\ & \leq & \frac{\mathsf{g}[\mathsf{X_{i,e}}]}{\epsilon \mathsf{t}} & \text{[Markov's inequality]} \\ & \leq & \frac{\mathsf{t/m}}{\epsilon \mathsf{t}} = \frac{1}{\mathsf{m}\epsilon} & \text{[derived above]} \end{array}$$

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

$$\begin{array}{lll} \Pr \big[f_e' - f_e \geq \epsilon t \big] & = & \Pr [\mathsf{C_i}[\mathsf{h_i}(e)] - f_e \geq \epsilon t \text{ for all } i] \\ & = & \Pr [\mathsf{X_{i,e}} \geq \epsilon t \text{ for all } i] \end{array}$$

$$egin{aligned} \mathsf{C_i}[\mathsf{h_i}(\mathsf{e})] &= \mathsf{X_{i,e}} + \mathsf{f_e} \text{ and } \mathsf{E}[\mathsf{X_{i,e}}] \leq \frac{\mathsf{t}}{\mathsf{m}}. \ \end{aligned}$$
 For $\epsilon > 0$

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

$$\text{Recall: } f_e' = \mathsf{est}_t(e) = \mathsf{min}_{i=1}^d \, C_i[h_i(e)].$$

$$\begin{array}{ll} \text{Pr}\big[f_e' - f_e \geq \epsilon t\big] &=& \text{Pr}\big[C_i[h_i(e)] - f_e \geq \epsilon t \text{ for all } i\big] \\ &=& \text{Pr}\big[X_{i,e} \geq \epsilon t \text{ for all } i\big] \\ &=& \Pi_{i=1}^d \, \text{Pr}\big[X_{i,e} \geq \epsilon t\big] \quad [\text{independence of } h_i's] \end{array}$$

$$egin{aligned} \mathsf{C_i}[\mathsf{h_i}(\mathsf{e})] &= \mathsf{X_{i,e}} + \mathsf{f_e} \text{ and } \mathsf{E}[\mathsf{X_{i,e}}] \leq \frac{\mathsf{t}}{\mathsf{m}}. \ \end{aligned}$$
 For $\epsilon > 0$

$$\begin{array}{lll} \text{Pr}[\mathsf{C_i}[\mathsf{h_i}(\mathsf{e})] - \mathsf{f_e} \geq \epsilon \mathsf{t}] & = & \text{Pr}[\mathsf{X_{i,e}} \geq \epsilon \mathsf{t}] & \text{[definition]} \\ & \leq & \frac{\mathsf{g}[\mathsf{X_{i,e}}]}{\epsilon \mathsf{t}} & \text{[Markov's inequality]} \\ & \leq & \frac{\mathsf{t/m}}{\epsilon \mathsf{t}} = \frac{1}{\mathsf{m}\epsilon} & \text{[derived above]} \end{array}$$

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

$$\begin{array}{ll} \mathsf{Pr}\big[f_e' - f_e \geq \epsilon t\big] &=& \mathsf{Pr}\big[\mathsf{C_i}[\mathsf{h_i}(e)] - f_e \geq \epsilon t \text{ for all } i\big] \\ &=& \mathsf{Pr}\big[\mathsf{X_{i,e}} \geq \epsilon t \text{ for all } i\big] \\ &=& \Pi_{i=1}^d \, \mathsf{Pr}\big[\mathsf{X_{i,e}} \geq \epsilon t\big] \quad [\text{independence of } \mathsf{h_i}\text{'s}] \\ &\leq& \left(\frac{1}{\epsilon m}\right)^d \qquad \qquad [\text{derived above}] \end{array}$$

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

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

Set
$$\mathbf{m} = \lceil \mathbf{2}/\epsilon \rceil$$
 and $\mathbf{d} = \lceil \lg \mathbf{1}/\delta \rceil$ gives us

$$\text{Pr}[\text{est}_{t}(\mathbf{e}) - \text{count}_{t}(\mathbf{e}) \geq \epsilon t] \leq \delta$$

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

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

Set $\mathbf{m} = \lceil \mathbf{2}/\epsilon \rceil$ and $\mathbf{d} = \lceil \lg 1/\delta \rceil$ gives us

$$\text{Pr}[\mathsf{est}_\mathsf{t}(\mathbf{e}) - \mathsf{count}_\mathsf{t}(\mathbf{e}) \geq \epsilon \mathbf{t}] \leq \delta$$

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

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

Lemma

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