CS 473: Algorithms

Ruta Mehta

University of Illinois, Urbana-Champaign

Spring 2018

CS 473: Algorithms, Spring 2018

Streaming Algorithms

Lecture 12 March 1, 2018

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.

Examples

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
- which connections are using more than 0.1% of the bandwidth.

Computing aggregative information about data streams.

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.

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].



Computing Sum

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

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



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.



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.



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?



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



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?



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!

Streaming Algorithms: Framework

```
 \begin{array}{l} & \langle \mbox{Initialize summary information} \rangle \\ & \mbox{While stream $\pmb{S}$ is not done} \\ & \pmb{x} \leftarrow \mbox{next element in $\pmb{S}$} \\ & \langle \mbox{Do something with $\pmb{x}$ and update summary information} \rangle \\ & \langle \mbox{Output something if needed} \rangle \\ \end{array}
```

Return \langle summary \rangle

```
⟨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⟩
```

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!

Puzzle Finding a Majority Element

Treasure hunt

Once upon a time...

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.

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?

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= \emptyset 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=Ø and counter=0

When element 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=Ø 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=Ø 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=Ø 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.

Finding the Majority Element: Example

R. Boyer and J. S. Moore Algorithm

Initialize: mem=Ø 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.
```

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

a _t	Е	D	В	D	D	D	В	В	В	В	В	•••
mem	E	E	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

ϵ -Heavy Hitters

Definition

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

$$\operatorname{count}_t(e) = |\{i \leq t \mid a_i = e\}|$$

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

ϵ -Heavy Hitters

Definition

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

$$\operatorname{count}_t(e) = |\{i \leq t \mid a_i = e\}|$$

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

Goal:

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

ϵ -Heavy Hitters

Definition

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

$$\operatorname{count}_t(e) = |\{i \leq t \mid a_i = e\}|$$

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

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.

Ruta (UIUC)

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

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 $T[1, \ldots, k]$ 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.

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 $j \le k$), then C[j] + +. Else if (C[j] == 0 for some $j \le k$), 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 $j \le k$), then C[j] + +. Else if (C[j] == 0 for some $j \le k$), 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[j] = 0 and $T[j] = \emptyset$ for all i.

When element a_t arrives, If $(a_t == T[j]$ for some $j \le k$), then C[j] + +. Else if (C[j] == 0 for some $j \le k$), 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 $k = \lceil 1/\epsilon \rceil - 1$. (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 $j \le k$), then C[j] + +. Else if (C[j] == 0 for some $j \le k$), 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 $\epsilon = 1/2$.

Ruta (UIUC)

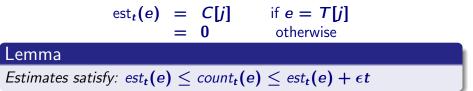
Lemma

At any time *t*, our estimates are:

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

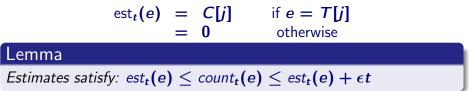
Estimates satisfy: $est_t(e) \leq count_t(e) \leq est_t(e) + \epsilon t$

At any time *t*, our estimates are:



For each element, count is maintained up to ϵt error!

At any time *t*, our estimates are:



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.

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 $\operatorname{count}_t(e) > \epsilon t$.

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,

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

Ruta (UIUC)

22

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,

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

Ruta (UIUC)

22

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 $e, \therefore \operatorname{est}_t(e) \leq \operatorname{count}_t(e)$.

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_t $(e) \le \text{count}_t(e)$. We want $\text{count}_t(e) - \text{est}_t(e) \le \epsilon t$. It increases by one,

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

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_t $(e) \le \text{count}_t(e)$. We want $\text{count}_t(e) - \text{est}_t(e) \le \epsilon t$. It increases by one,

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

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_t $(e) \le \text{count}_t(e)$. We want $\text{count}_t(e) - \text{est}_t(e) \le \epsilon t$. It increases by one,

- ullet when we decrease all $m{k}$ counters, and see an element outside $m{ au}$
- 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 increases.

ϵ -Heavy Hitters: Algorithm Space usage

Set
$$k = \lfloor 1/\epsilon \rfloor - 1$$
. (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

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

ϵ -Heavy Hitters: Algorithm Space usage

Set
$$k = \lfloor 1/\epsilon \rfloor - 1$$
. (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

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

ϵ -Heavy Hitters: Algorithm Space usage

Set
$$k = \lfloor 1/\epsilon \rfloor - 1$$
. (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

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

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

Recall: maintains counts for all elements up to ϵt 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 ϵt is OK, then we can use ϵ -heavy hitter algorithm.

Problem Statement:

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

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

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

Problem Statement:

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

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

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

Can we do better?

Yes – Bloom filter like idea

Storage for inserts and lookups

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

$$\begin{array}{c} \text{Insert}(e) \\ \text{For } i = 1...d \\ \text{Set } T_i[h_i(e)] \leftarrow 1 \end{array} \end{array} \begin{array}{c} \text{Lookup}(e) \\ \text{For } i = 1...d \\ \text{If } (T_i[h_i(e)] == 0) \text{ then return "No"} \\ \text{Return "Yes"} \end{array}$$

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

Storage for inserts and lookups

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

Insert(e)Lookup(e)For
$$i = 1...d$$
For $i = 1...d$ Set $T_i[h_i(e)] \leftarrow 1$ 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∼H}[e not inserted and T_i[h_i(e)] = 1] ≤ α, then combined error probability would be at most α^d.

Ruta (UIUC)

Count Min-Sketch By G. Cormode and S. M. Muthukrishnan'05

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

Keep d arrays C_1, \ldots, 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} .

Keep d arrays C_1, \ldots, 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 $C_i[h_i(e)] + +$

Keep d arrays C_1, \ldots, 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...dDo $C_i[h_i(e)] + +$ CMEstimate(e) $est \leftarrow \infty$ For i = 1...d $est \leftarrow \min\{est, C_i[h_i(e)]\}$ Return est

Keep d arrays C_1, \ldots, 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...dDo $C_i[h_i(e)] + +$ CMEstimate(e) $est \leftarrow \infty$ For i = 1...d $est \leftarrow \min\{est, C_i[h_i(e)]\}$ Return est

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)

$$CMInsert(e)$$

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

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

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

$$CMInsert(e)$$

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 $\operatorname{est}_t(e) = \operatorname{CMEstimate}(e) = \min_{i=1}^d C_i[h_i(e)].$ Observation: $\operatorname{est}_t(e) \ge \operatorname{count}_t(e).$

$$CMInsert(e)$$

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 $est_t(e) = CMEstimate(e) = \min_{i=1}^d C_i[h_i(e)]$. Observation: $est_t(e) \ge count_t(e)$.

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

$$CMInsert(e)$$

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 $est_t(e) = CMEstimate(e) = \min_{i=1}^d C_i[h_i(e)]$. Observation: $est_t(e) \ge count_t(e)$.

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

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

Let $f'_e = \operatorname{est}_t(e)$ and $f_e = \operatorname{count}_t(e)$. We want to bound $(f'_e - f_e)$.

Observations:

Let
$$f'_e = \operatorname{est}_t(e)$$
 and $f_e = \operatorname{count}_t(e)$. We want to bound $(f'_e - f_e)$.

Observations:

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

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

Let $f'_e = \operatorname{est}_t(e)$ and $f_e = \operatorname{count}_t(e)$. We want to bound $(f'_e - f_e)$.

Observations:

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

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

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

Let
$$f'_e = \operatorname{est}_t(e)$$
 and $f_e = \operatorname{count}_t(e)$. We want to bound $(f'_e - f_e)$.

Observations:

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

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

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,

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

Let $f'_e = \operatorname{est}_t(e)$ and $f_e = \operatorname{count}_t(e)$. We want to bound $(f'_e - f_e)$.

Observations:

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

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

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,

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

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

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

Then, for $\epsilon > 0$

 $\Pr[C_i[h_i(e)] - f_e \ge \epsilon t] = \Pr[X_{i,e} \ge \epsilon t]$ [definition]

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

Then, for $\epsilon > 0$

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

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

Then, for $\epsilon > 0$

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

We have $C_i[h_i(e)] = X_{i,e} + f_e$ and $E[X_{i,e}] \le \frac{t}{m}$. Then, for $\epsilon > 0$

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

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

We have $C_i[h_i(e)] = X_{i,e} + f_e$ and $E[X_{i,e}] \le \frac{t}{m}$. Then, for $\epsilon > 0$

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

Recall: $f'_e = \operatorname{est}_t(e) = \min_{i=1}^d C_i[h_i(e)].$ $\Pr[f'_e - f_e \ge \epsilon t] =$

We have $C_i[h_i(e)] = X_{i,e} + f_e$ and $E[X_{i,e}] \le \frac{t}{m}$. Then, for $\epsilon > 0$

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

Recall: $f'_e = \operatorname{est}_t(e) = \min_{i=1}^d C_i[h_i(e)].$ $\Pr[f'_e - f_e \ge \epsilon t] = \Pr[C_i[h_i(e)] - f_e \ge \epsilon t \text{ for all } i]$

We have $C_i[h_i(e)] = X_{i,e} + f_e$ and $E[X_{i,e}] \le \frac{t}{m}$. Then, for $\epsilon > 0$

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

Recall: $f'_e = \operatorname{est}_t(e) = \min_{i=1}^d C_i[h_i(e)].$ $\Pr[f'_e - f_e \ge \epsilon t] = \Pr[C_i[h_i(e)] - f_e \ge \epsilon t \text{ for all } i]$ $= \Pr[X_{i,e} \ge \epsilon t \text{ for all } i]$

We have $C_i[h_i(e)] = X_{i,e} + f_e$ and $E[X_{i,e}] \le \frac{t}{m}$. Then, for $\epsilon > 0$

$$\begin{aligned} \mathsf{Pr}[\mathcal{C}_{i}[h_{i}(e)] - f_{e} \geq \epsilon t] &= \mathsf{Pr}[X_{i,e} \geq \epsilon t] & \text{[definition]} \\ &\leq \frac{\underline{\mathsf{f}}[x_{i,e}]}{\frac{t/\overline{m}}{\epsilon t}} & \text{[Markov's inequality]} \\ &\leq \frac{t/\overline{m}}{\epsilon t} = \frac{1}{m\epsilon} & \text{[derived above]} \end{aligned}$$

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

$$\begin{aligned} \mathsf{Pr}\big[f'_e - f_e \geq \epsilon t\big] &= \mathsf{Pr}[C_i[h_i(e)] - f_e \geq \epsilon t \text{ for all } i] \\ &= \mathsf{Pr}[X_{i,e} \geq \epsilon t \text{ for all } i] \\ &= \mathsf{\Pi}^d_{i=1} \mathsf{Pr}[X_{i,e} \geq \epsilon t] \quad [\text{independence of } h_i's] \end{aligned}$$

We have $C_i[h_i(e)] = X_{i,e} + f_e$ and $E[X_{i,e}] \le \frac{t}{m}$. Then, for $\epsilon > 0$

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

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

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

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

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

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

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

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

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

Space:

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

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.

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 $O(\frac{1}{\epsilon} \lg \frac{1}{\delta})$ many counters.