

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Andrew McGregor

Lecture 9

THE FREQUENT ITEMS DATA STREAM PROBLEM

k -Frequent Items (Heavy-Hitters) Problem: Consider a stream of n items x_1, \dots, x_n (with possible duplicates). Return any item that appears at least $\frac{n}{k}$ times.

THE FREQUENT ITEMS DATA STREAM PROBLEM

k -Frequent Items (Heavy-Hitters) Problem: Consider a stream of n items x_1, \dots, x_n (with possible duplicates). Return any item that appears at least $\frac{n}{k}$ times.

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9
5	12	3	3	4	5	5	10	3

THE FREQUENT ITEMS DATA STREAM PROBLEM

k -Frequent Items (Heavy-Hitters) Problem: Consider a stream of n items x_1, \dots, x_n (with possible duplicates). Return any item that appears at least $\frac{n}{k}$ times.

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9
5	12	3	3	4	5	5	10	3

THE FREQUENT ITEMS DATA STREAM PROBLEM

k -Frequent Items (Heavy-Hitters) Problem: Consider a stream of n items x_1, \dots, x_n (with possible duplicates). Return any item that appears at least $\frac{n}{k}$ times.

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9
5	12	3	3	4	5	5	10	3

- What is the maximum number of items that must be returned?
a) n b) k c) n/k d) $\log n$

THE FREQUENT ITEMS DATA STREAM PROBLEM

k -Frequent Items (Heavy-Hitters) Problem: Consider a stream of n items x_1, \dots, x_n (with possible duplicates). Return any item that appears at least $\frac{n}{k}$ times.

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9
5	12	3	3	4	5	5	10	3

- What is the maximum number of items that must be returned?
a) n b) k c) n/k d) $\log n$

THE FREQUENT ITEMS DATA STREAM PROBLEM

k -Frequent Items (Heavy-Hitters) Problem: Consider a stream of n items x_1, \dots, x_n (with possible duplicates). Return any item that appears at least $\frac{n}{k}$ times.

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9
5	12	3	3	4	5	5	10	3

- What is the maximum number of items that must be returned?
a) n b) k c) n/k d) $\log n$
- Trivial with $O(n)$ space: Store the count for each item and return the one that appears $\geq n/k$ times.
- Can we do it with less space? I.e., without storing all n items?

Applications of Frequent Items:

Applications of Frequent Items:

- Finding top/viral items (i.e., products on Amazon, videos watched on Youtube, Google searches, etc.)

Applications of Frequent Items:

- Finding top/viral items (i.e., products on Amazon, videos watched on Youtube, Google searches, etc.)
- Finding very frequent IP addresses sending requests (to detect DoS attacks/network anomalies).

Applications of Frequent Items:

- Finding top/viral items (i.e., products on Amazon, videos watched on Youtube, Google searches, etc.)
- Finding very frequent IP addresses sending requests (to detect DoS attacks/network anomalies).
- 'Iceberg queries' for all items in a database with frequency above some threshold.

Applications of Frequent Items:

- Finding top/viral items (i.e., products on Amazon, videos watched on Youtube, Google searches, etc.)
- Finding very frequent IP addresses sending requests (to detect DoS attacks/network anomalies).
- 'Iceberg queries' for all items in a database with frequency above some threshold.

Generally want very fast detection, without having to scan through database/logs. That is we want to maintain a running list of frequent items that appear in a stream.


APPROXIMATE FREQUENT ELEMENTS

Issue: No algorithm using $o(n)$ space can output just the items with frequency $\geq n/k$. Hard to tell between an item with frequency n/k (should be output) and $n/k - 1$ (should not be output).

APPROXIMATE FREQUENT ELEMENTS

Issue: No algorithm using $o(n)$ space can output just the items with frequency $\geq n/k$. Hard to tell between an item with frequency n/k (should be output) and $n/k - 1$ (should not be output).

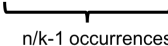
x_1	x_2	x_3	x_4	x_5	x_6	...	$x_{n-n/k+1}$...	x_n
3	12	9	27	4	101		3		3


n/k-1 occurrences

APPROXIMATE FREQUENT ELEMENTS

Issue: No algorithm using $o(n)$ space can output just the items with frequency $\geq n/k$. Hard to tell between an item with frequency n/k (should be output) and $n/k - 1$ (should not be output).

x_1	x_2	x_3	x_4	x_5	x_6	...	$x_{n-n/k+1}$...	x_n
3	12	9	27	4	101		3		3


n/k-1 occurrences

(ϵ, k) -Frequent Items Problem: Consider a stream of n items x_1, \dots, x_n . Return a set F of items, including all items that appear at least $\frac{n}{k}$ times and only items that appear at least $(1 - \epsilon) \cdot \frac{n}{k}$ times.

APPROXIMATE FREQUENT ELEMENTS

Issue: No algorithm using $o(n)$ space can output just the items with frequency $\geq n/k$. Hard to tell between an item with frequency n/k (should be output) and $n/k - 1$ (should not be output).

x_1	x_2	x_3	x_4	x_5	x_6	...	$x_{n-n/k+1}$...	x_n
3	12	9	27	4	101		3		3

└──────────────────┘
n/k-1 occurrences

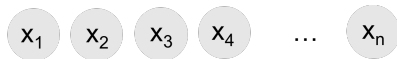
(ϵ, k) -Frequent Items Problem: Consider a stream of n items x_1, \dots, x_n . Return a set F of items, including all items that appear at least $\frac{n}{k}$ times and only items that appear at least $(1 - \epsilon) \cdot \frac{n}{k}$ times.

- An example of relaxing to a 'promise problem': for items with frequencies in $[(1 - \epsilon) \cdot \frac{n}{k}, \frac{n}{k}]$ no output guarantee.

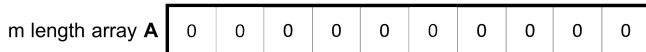
Today: Count-min sketch – a random hashing based method closely related to bloom filters.

FREQUENT ELEMENTS WITH COUNT-MIN SKETCH

Today: Count-min sketch – a random hashing based method closely related to bloom filters.

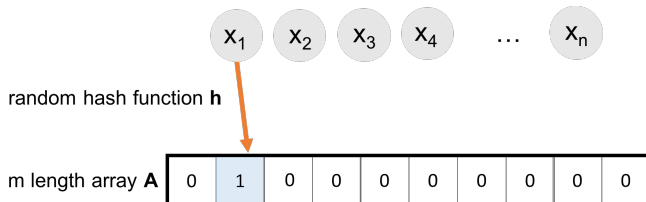


random hash function h



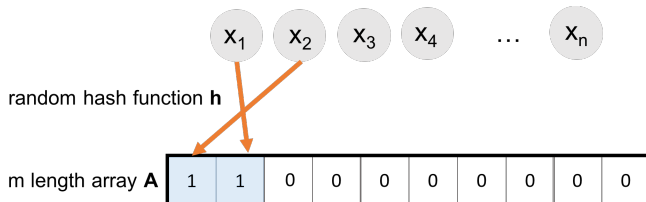
FREQUENT ELEMENTS WITH COUNT-MIN SKETCH

Today: Count-min sketch – a random hashing based method closely related to bloom filters.



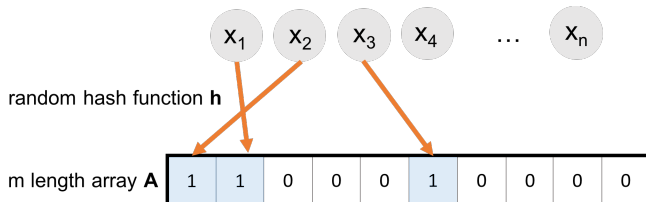
FREQUENT ELEMENTS WITH COUNT-MIN SKETCH

Today: Count-min sketch – a random hashing based method closely related to bloom filters.



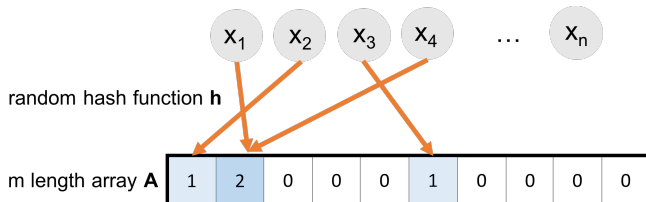
FREQUENT ELEMENTS WITH COUNT-MIN SKETCH

Today: Count-min sketch – a random hashing based method closely related to bloom filters.



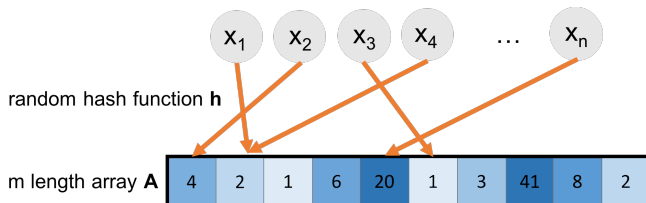
FREQUENT ELEMENTS WITH COUNT-MIN SKETCH

Today: Count-min sketch – a random hashing based method closely related to bloom filters.



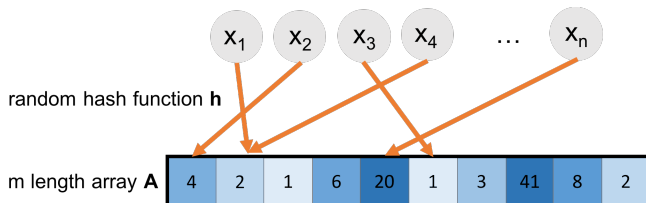
FREQUENT ELEMENTS WITH COUNT-MIN SKETCH

Today: Count-min sketch – a random hashing based method closely related to bloom filters.



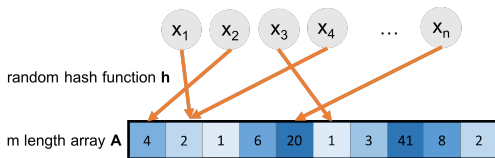
FREQUENT ELEMENTS WITH COUNT-MIN SKETCH

Today: Count-min sketch – a random hashing based method closely related to bloom filters.



Will use $A[h(x)]$ to estimate $f(x) = |\{i : x_i = x\}|$, the frequency of x in the stream.

COUNT-MIN SKETCH ACCURACY

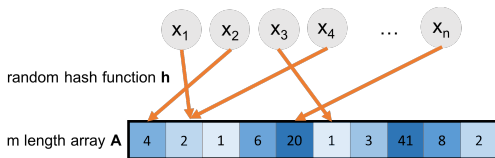


Use $A[h(x)]$ to estimate $f(x)$.

Claim 1: We always have $A[h(x)] \geq f(x)$. Why?

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY



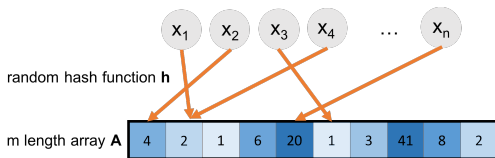
Use $A[h(x)]$ to estimate $f(x)$.

Claim 1: We always have $A[h(x)] \geq f(x)$. Why?

- $A[h(x)]$ counts the number of occurrences of any y with $h(y) = h(x)$, including x itself.

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY



Use $A[h(x)]$ to estimate $f(x)$.

Claim 1: We always have $A[h(x)] \geq f(x)$. Why?

- $A[h(x)]$ counts the number of occurrences of any y with $h(y) = h(x)$, including x itself.
- $A[h(x)] = f(x) + \sum_{y \neq x: h(y) = h(x)} f(y)$.

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] =$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] = \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y)$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &\leq \sum_{y \neq x} \frac{1}{m} \cdot f(y) \end{aligned}$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &\leq \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m} \end{aligned}$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &\leq \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m} \end{aligned}$$

What is a bound on probability that the error is $\geq \frac{2n}{m}$?

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &\leq \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m} \end{aligned}$$

What is a bound on probability that the error is $\geq \frac{2n}{m}$?

Markov's inequality: $\Pr \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \geq \frac{2n}{m} \right] \leq \frac{1}{2}$.

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &\leq \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m} \end{aligned}$$

What is a bound on probability that the error is $\geq \frac{2n}{m}$?

Markov's inequality: $\Pr \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \geq \frac{2n}{m} \right] \leq \frac{1}{2}$.

What property of \mathbf{h} is required to show this bound? a) fully random b) pairwise independent c) 2-universal d) locality sensitive

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY

$$A[\mathbf{h}(x)] = f(x) + \underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\text{error in frequency estimate}} .$$

Expected Error:

$$\begin{aligned} \mathbb{E} \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \right] &= \sum_{y \neq x} \Pr(\mathbf{h}(y) = \mathbf{h}(x)) \cdot f(y) \\ &\leq \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m} \end{aligned}$$

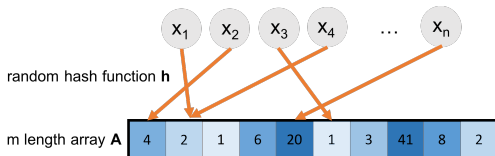
What is a bound on probability that the error is $\geq \frac{2n}{m}$?

Markov's inequality: $\Pr \left[\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y) \geq \frac{2n}{m} \right] \leq \frac{1}{2}$.

What property of \mathbf{h} is required to show this bound? a) fully random b) pairwise independent c) 2-universal d) locality sensitive

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). \mathbf{h} : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY

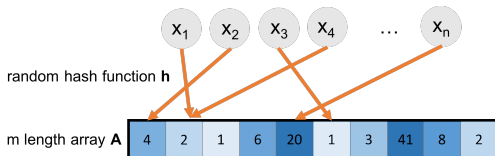


Claim: For any x , with probability at least $1/2$,

$$f(x) \leq A[h(x)] \leq f(x) + \frac{2n}{m}.$$

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY



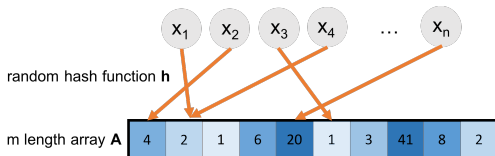
Claim: For any x , with probability at least $1/2$,

$$f(x) \leq A[h(x)] \leq f(x) + \frac{2n}{m}.$$

How can we improve the success probability?

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY



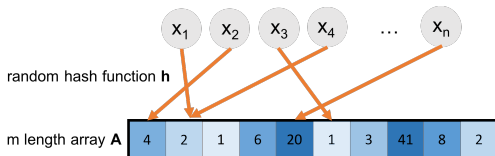
Claim: For any x , with probability at least $1/2$,

$$f(x) \leq A[h(x)] \leq f(x) + \frac{2n}{m}.$$

How can we improve the success probability?

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of Count-min sketch array.

COUNT-MIN SKETCH ACCURACY



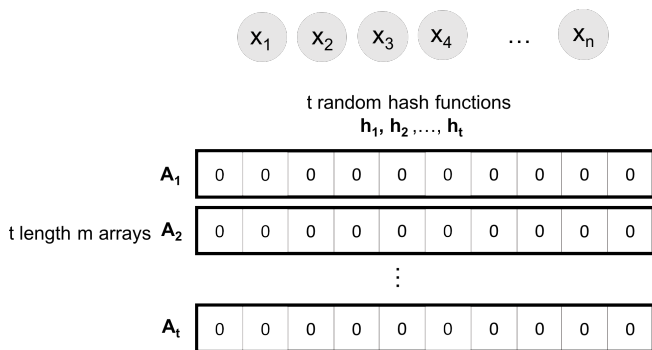
Claim: For any x , with probability at least $1/2$,

$$f(x) \leq A[h(x)] \leq f(x) + \frac{2n}{m}.$$

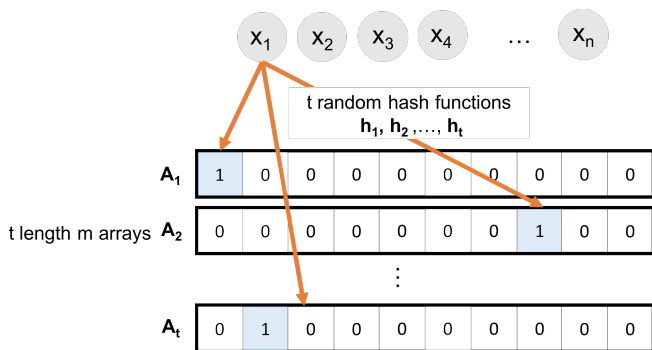
How can we improve the success probability? **Repetition.**

$f(x)$: frequency of x in the stream (i.e., number of items equal to x). h : random hash function. m : size of Count-min sketch array.

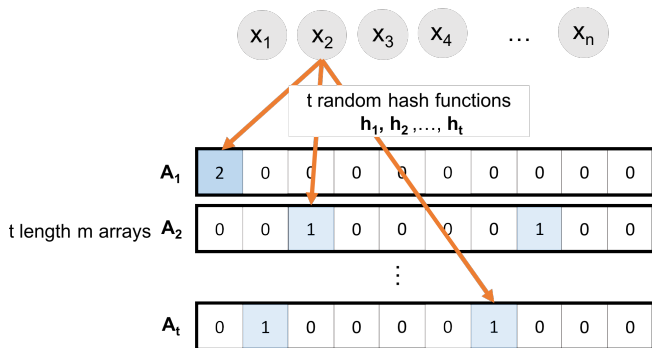
COUNT-MIN SKETCH ACCURACY



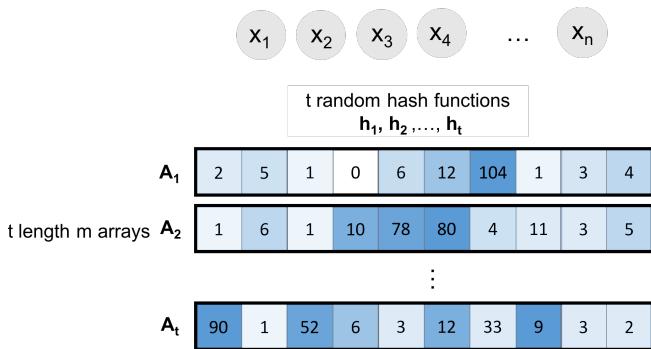
COUNT-MIN SKETCH ACCURACY



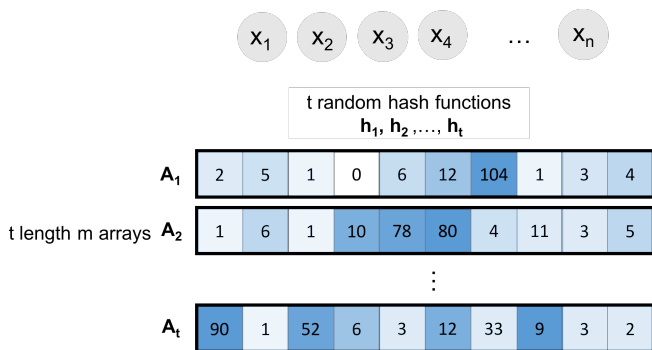
COUNT-MIN SKETCH ACCURACY



COUNT-MIN SKETCH ACCURACY

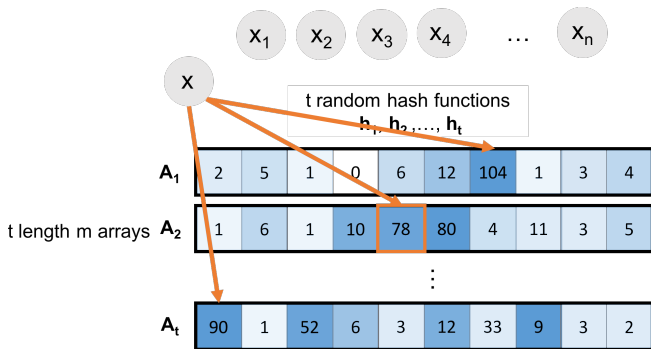


COUNT-MIN SKETCH ACCURACY



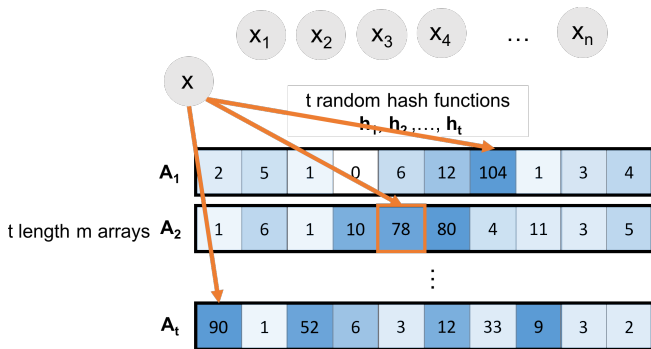
Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$. (count-min sketch)

COUNT-MIN SKETCH ACCURACY



Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$. (count-min sketch)

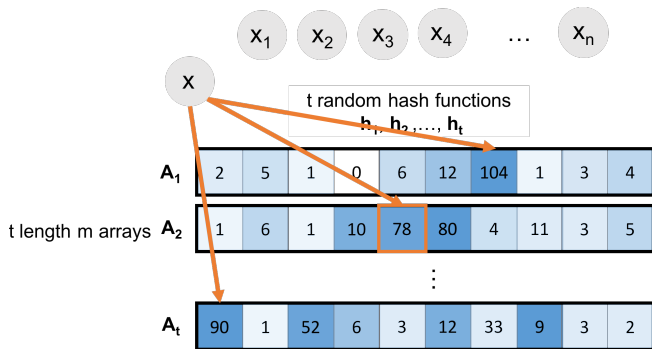
COUNT-MIN SKETCH ACCURACY



Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$. (count-min sketch)

Why min instead of mean or median?

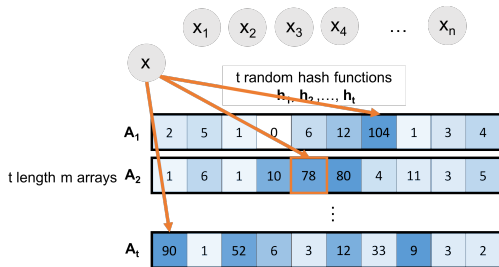
COUNT-MIN SKETCH ACCURACY



Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$. (count-min sketch)

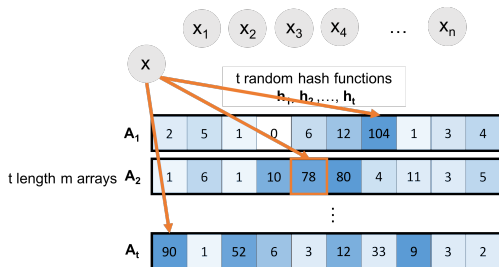
Why min instead of mean or median? The minimum estimate is always the most accurate since they are all overestimates of the true frequency!

COUNT-MIN SKETCH ANALYSIS



Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

COUNT-MIN SKETCH ANALYSIS

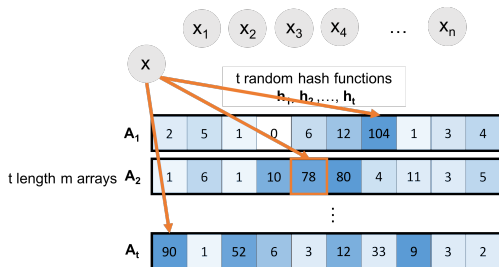


Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- For every x and $i \in [t]$, we know that for $m = \frac{2k}{\epsilon}$, with probability $\geq 1/2$:

$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

COUNT-MIN SKETCH ANALYSIS



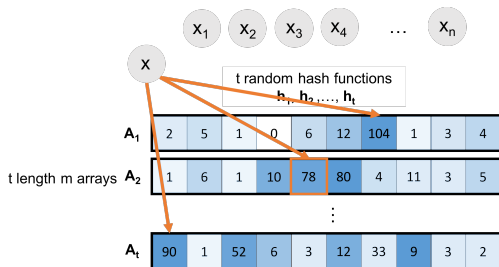
Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- For every x and $i \in [t]$, we know that for $m = \frac{2k}{\epsilon}$, with probability $\geq 1/2$:

$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

- What is $\Pr[f(x) \leq \tilde{f}(x) \leq f(x) + \frac{\epsilon n}{k}]$?

COUNT-MIN SKETCH ANALYSIS



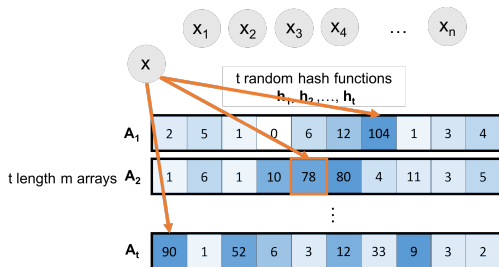
Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- For every x and $i \in [t]$, we know that for $m = \frac{2k}{\epsilon}$, with probability $\geq 1/2$:

$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

- What is $\Pr[f(x) \leq \tilde{f}(x) \leq f(x) + \frac{\epsilon n}{k}]$? $1 - 1/2^t$.

COUNT-MIN SKETCH ANALYSIS



Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- For every x and $i \in [t]$, we know that for $m = \frac{2k}{\epsilon}$, with probability $\geq 1/2$:

$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

- What is $\Pr[f(x) \leq \tilde{f}(x) \leq f(x) + \frac{\epsilon n}{k}]$? $1 - 1/2^t$.
- To get a good estimate with probability $\geq 1 - \delta$, set $t = \log(1/\delta)$.

Upshot: Count-min sketch lets us estimate the frequency of every item in a stream up to error $\frac{\epsilon n}{k}$ with probability $\geq 1 - \delta$ in $O(\log(1/\delta) \cdot k/\epsilon)$ space.

Upshot: Count-min sketch lets us estimate the frequency of every item in a stream up to error $\frac{\epsilon n}{k}$ with probability $\geq 1 - \delta$ in $O(\log(1/\delta) \cdot k/\epsilon)$ space.

- Accurate enough to solve the (ϵ, k) -Frequent elements problem: Can distinguish between items with frequency $\frac{n}{k}$ and those with frequency $< (1 - \epsilon)\frac{n}{k}$.

Upshot: Count-min sketch lets us estimate the frequency of every item in a stream up to error $\frac{\epsilon n}{k}$ with probability $\geq 1 - \delta$ in $O(\log(1/\delta) \cdot k/\epsilon)$ space.

- Accurate enough to solve the (ϵ, k) -Frequent elements problem: Can distinguish between items with frequency $\frac{n}{k}$ and those with frequency $< (1 - \epsilon)\frac{n}{k}$.
- How should we set δ if we want a good estimate for all items at once, with 99% probability?

Upshot: Count-min sketch lets us estimate the frequency of every item in a stream up to error $\frac{\epsilon n}{k}$ with probability $\geq 1 - \delta$ in $O(\log(1/\delta) \cdot k/\epsilon)$ space.

- Accurate enough to solve the (ϵ, k) -Frequent elements problem: Can distinguish between items with frequency $\frac{n}{k}$ and those with frequency $< (1 - \epsilon)\frac{n}{k}$.
- How should we set δ if we want a good estimate for all items at once, with 99% probability? $\delta = 0.01/|U|$ ensures

$$\begin{aligned} & \Pr[\text{there exists } x \in U \text{ with a bad estimate}] \\ & \leq \sum_{x \in U} \Pr[\text{estimate for } x \text{ is bad}] \leq \sum_{x \in U} 0.01/|U| = 0.01 \end{aligned}$$

Count-min sketch gives an accurate frequency estimate for every item in the stream. But how do we identify the frequent items without having to look up the estimated frequency for $x \in U$?

Count-min sketch gives an accurate frequency estimate for every item in the stream. But how do we identify the frequent items without having to look up the estimated frequency for $x \in U$?

One approach:

- Maintain a set F while processing the stream:
- At step i :
 - Add i th stream element to F if its estimated frequency is $\geq i/k$ and it isn't already in F .
 - Remove any element from F whose estimated frequency is $< i/k$.
- Store at most k items at once and have all items with frequency $\geq n/k$ stored at the end of the stream.

Questions on Frequent Elements?