


COMPSCI 514: Algorithms for Data Science

Cameron Musco

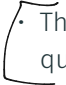
University of Massachusetts Amherst. Fall 2023.

Lecture 11

- Problem Set 2 is due on Monday at 11:59pm.

A hand-drawn black bracket on the left side of the slide, grouping the second and third bullet points. It consists of a vertical line with a horizontal top bar and a horizontal bottom bar, both ending in small inward-pointing hooks.

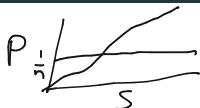
- Midterm is in class Tuesday, 10/24. Thursday 10/19 will be devoted to midterm review.

A hand-drawn black bracket on the left side of the slide, grouping the third bullet point. It consists of a vertical line with a horizontal top bar and a horizontal bottom bar, both ending in small inward-pointing hooks.

- The grading on this week's quiz regarding the extra credit question had a bug. We will fix manually in the next few days.

Summary

Last Class:



- Locality sensitive hashing to solve the similarity search problem efficiently. S-curve
- MinHash as a locality sensitive hash function for Jaccard similarity.
- Brief look at SimHash as a locality sensitive hash function for cosine (dot product) similarity.

This Class:

Can a locality sensitive hash function be pairwise independent?

- Introduce the frequent elements problem and its applications.
- Solution via the Count-Min sketch randomized data structure.

$$P(h(x)=i \wedge h(y)=j) = \frac{1}{n^2}$$
$$\left(P(h(x)=h(y)) = \frac{1}{n} \right)$$

The Frequent Items Problems

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.

$$k = 10$$

The Frequent Items Problems

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 Problems

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. $k = 3$

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

all items that appeared $\frac{n}{k} = \frac{9}{3} = 3$ the

The Frequent Items Problems

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. $k=10$

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 can be returned?
a) n b) k c) n/k d) $\log n$

$$\Rightarrow \frac{n}{k} \geq n$$

The Frequent Items Problems

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 can 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?

The Frequent Items Problem

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. — $\frac{D}{K}$

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

Frequent Itemset Mining

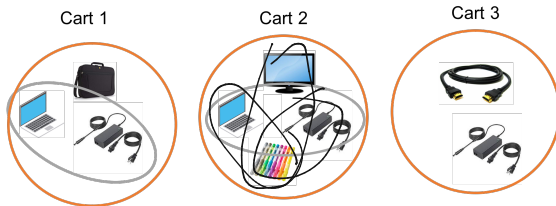
Association rule learning: A very common task in data mining is to identify common associations between different events.



- Identified via **frequent itemset** counting. Find all sets of t items that appear many times in the same basket.
- Frequency of an itemset is known as its support.
- A single basket includes many different itemsets, and with many different baskets an efficient approach is critical. E.g., baskets are Twitter users and itemsets are subsets of who they follow.

Frequent Itemset Mining

Association rule learning: A very common task in data mining is to identify common associations between different events.



- Identified via **frequent itemset** counting. Find all sets of t items that appear many times in the same basket.
- Frequency of an itemset is known as its support.
- A single basket includes many different itemsets, and with many different baskets an efficient approach is critical. E.g., baskets are Twitter users and itemsets are subsets of who they follow.

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.

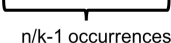
$k = 10 \rightarrow$ return all items that show up at least $0.1 \cdot n$ times

$\epsilon = 0.1 \rightarrow$ should not return any item showing up $\leq (1 - \epsilon) \cdot 0.1 \cdot n = 0.09 \cdot n$ 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**.

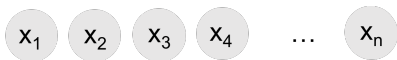
- An example of relaxing to a ‘promise problem’: for items with frequencies in $[(1 - \epsilon) \cdot \frac{n}{k}, \frac{n}{k}]$ no output guarantee.

Frequent Elements with Count-Min Sketch

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

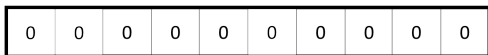
Frequent Elements with Count-Min Sketch

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



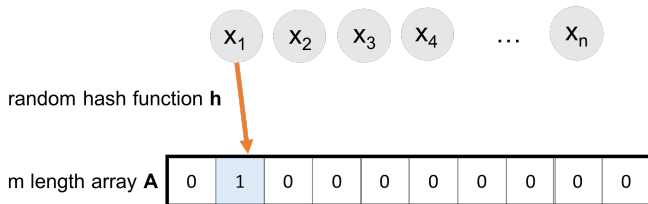
random hash function h

m length array \mathbf{A}



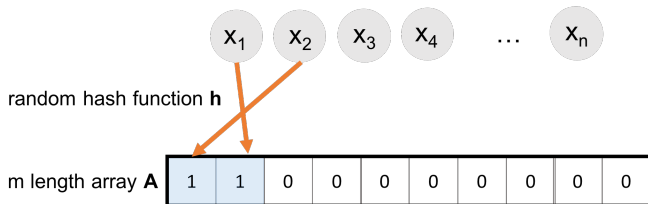
Frequent Elements with Count-Min Sketch

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



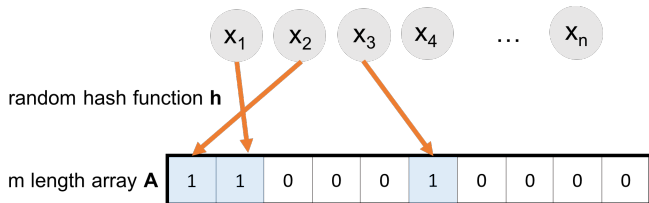
Frequent Elements with Count-Min Sketch

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



Frequent Elements with Count-Min Sketch

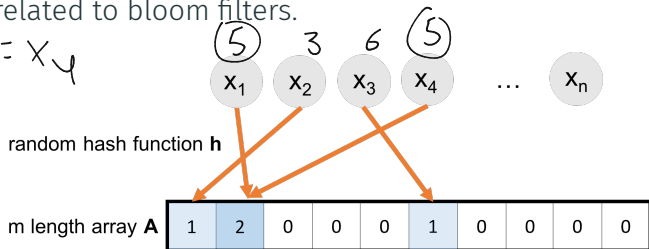
Today: Count-min sketch – a random hashing based method closely related to bloom filters.



Frequent Elements with Count-Min Sketch

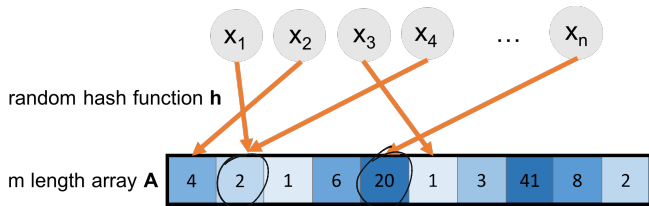
Today: Count-min sketch – a random hashing based method closely related to bloom filters.

$$x_1 = x_4$$



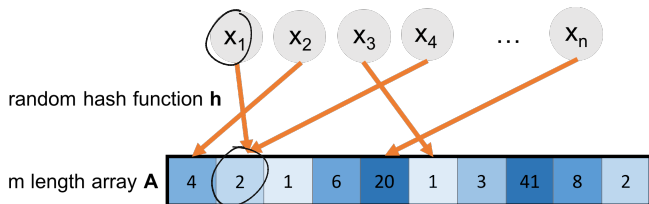
Frequent Elements with Count-Min Sketch

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



Frequent Elements with Count-Min Sketch

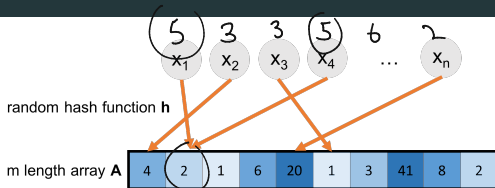
Today: Count-min sketch – a random hashing based method closely related to bloom filters.



Will use $A[\overset{2}{h}(x)]$ to estimate $f(x)$, the frequency of x in the stream. I.e., $|\{x_i : x_i = x\}|$.

Count-Min Sketch Accuracy

$$h(5) = 2$$
$$A[2] \approx 2$$



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

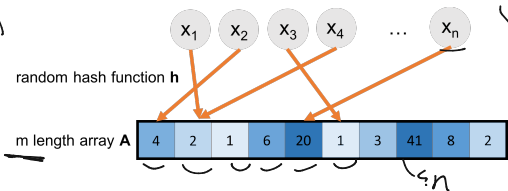
Claim 1: We always have $\underline{A[h(x)]} \geq \underline{f(x)}$. Why?

$f(x)$: 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

$$m \ll n$$

$$m \cdot \log(n)$$



Store m
integers of
size $\leq n$

this takes
 $m \cdot \log_2 n$
bits

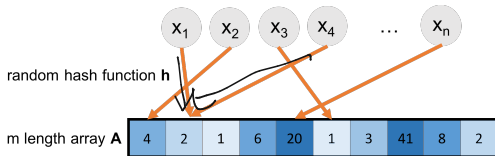
Use $A[h(x)]$ to estimate $f(x)$.

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

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

Expected Error:

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

$f(x)$: 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

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

Expected Error:

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

$\underbrace{\hspace{10em}}_{\frac{1}{m}}$

$f(x)$: 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

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

Expected Error:

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

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

Count-Min Sketch Accuracy

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

Expected Error:

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

Handwritten notes:
A bracket under the term $\frac{1}{m} \sum_{y \neq x} f(y)$ is labeled with $\frac{1}{m} \sum_{y \neq x} f(y)$.

$f(x)$: 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

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

Expected Error:

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

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

↳ Markov

$f(x)$: 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

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

Expected Error:

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

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

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

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

Count-Min Sketch Accuracy

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

Expected Error:

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

Handwritten notes: a bracket under the sum is labeled '1/m', and a '3' is written above the sum.

$$\leq \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m}$$

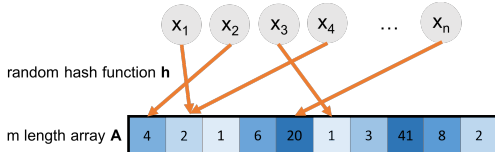
What is a bound on probability that the error is $\geq \frac{2n}{m}$?

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

What property of 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). h : random hash function. m : size of Count-min sketch array.

Count-Min Sketch Accuracy

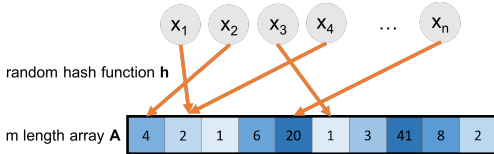


Claim: For any x , with probability at least $1/2$,

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

$f(x)$: 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



$$\geq \frac{n}{k}$$

$$\leq \frac{(1-\epsilon)n}{k}$$

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

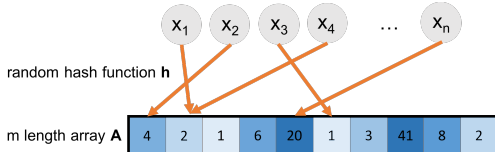
$$f(x) \leq A[h(x)] \leq f(x) + \frac{2n}{m}$$

$$\frac{2n}{m} = \frac{2n}{\frac{2k}{\epsilon}} = \epsilon \cdot \frac{n}{k}$$

To solve the (ϵ, k) -Frequent elements problem, set $m = \frac{2k}{\epsilon}$.

$f(x)$: 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



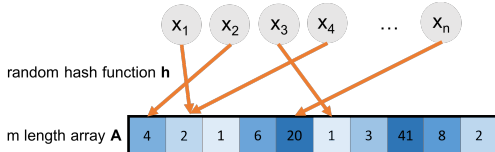
Claim: For any x , with probability at least $1/2$,

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

To solve the (ϵ, k) -Frequent elements problem, set $m = \frac{2k}{\epsilon}$. 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



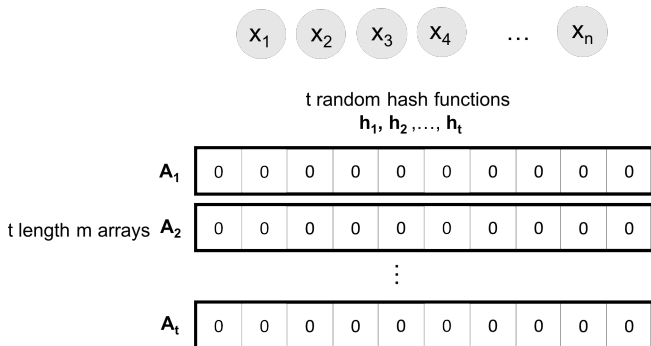
Claim: For any x , with probability at least $1/2$,

$$f(x) \leq \underline{A[h(x)]} \leq f(x) + \frac{2n}{m} \cdot \frac{\epsilon n}{k}$$

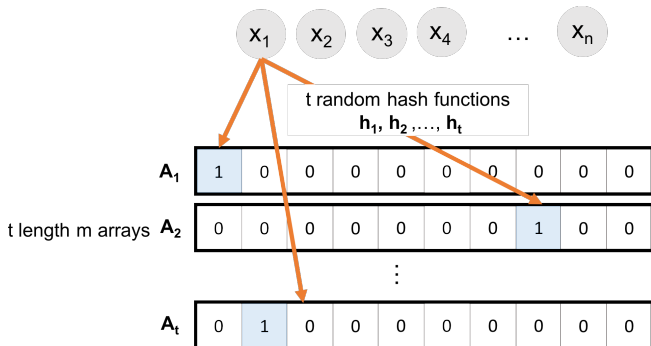
To solve the (ϵ, k) -Frequent elements problem, set $m = \frac{2k}{\epsilon}$. How can we improve the success probability? **Repetition.**

$f(x)$: 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.

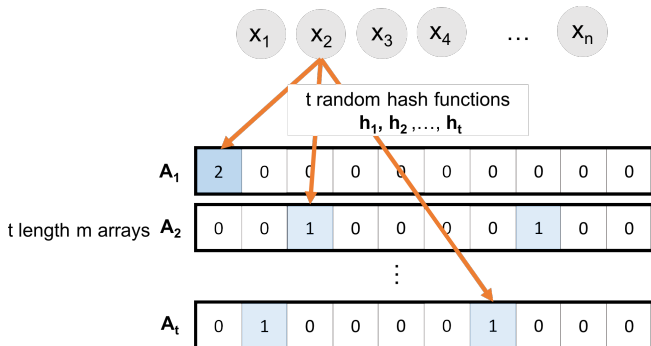
Count-Min Sketch Repetition



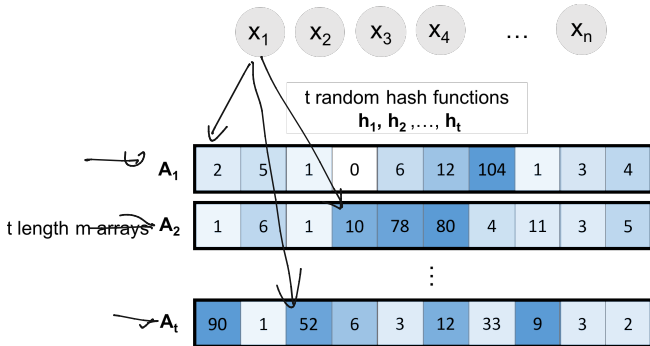
Count-Min Sketch Repetition



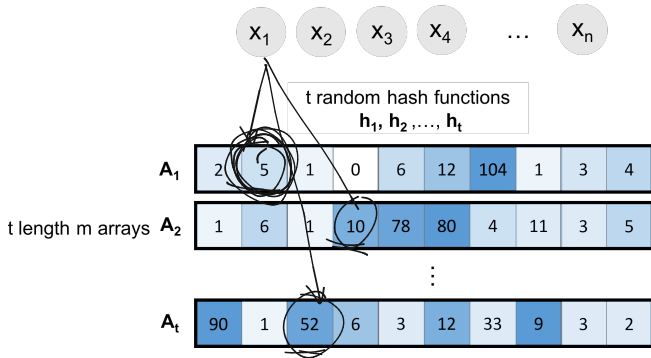
Count-Min Sketch Repetition



Count-Min Sketch Repetition



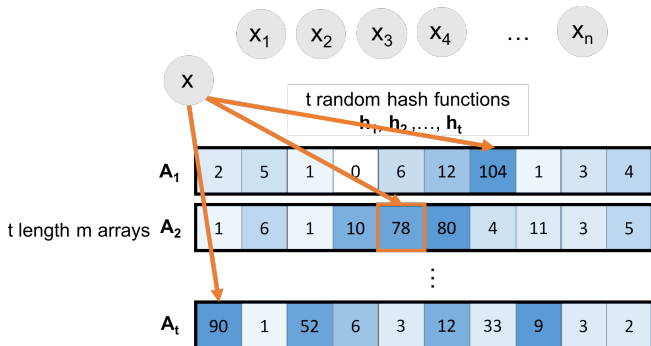
Count-Min Sketch Repetition



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

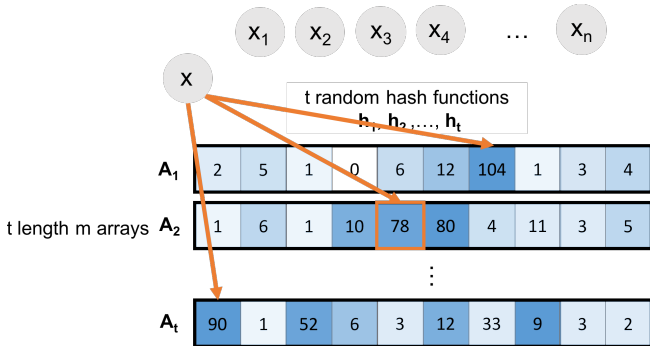
≈ 5

Count-Min Sketch Repetition



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

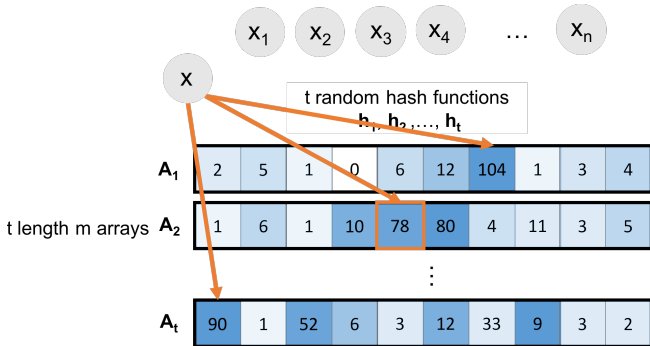
Count-Min Sketch Repetition



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?

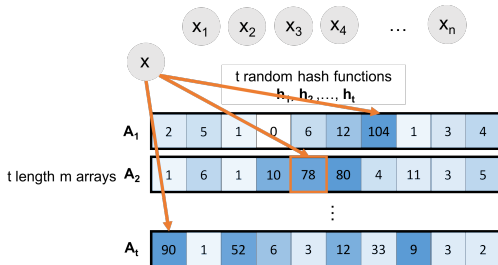
Count-Min Sketch Repetition



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

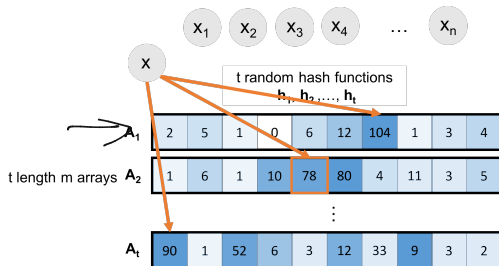
Why min instead of mean or median? The minimum estimate is always the most accurate since they are all overestimates of the true frequency!

Count-Min Sketch Analysis



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

Count-Min Sketch Analysis

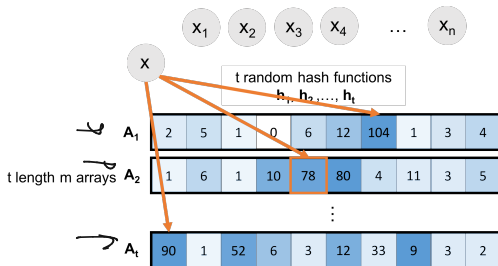


Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

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

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

Count-Min Sketch Analysis



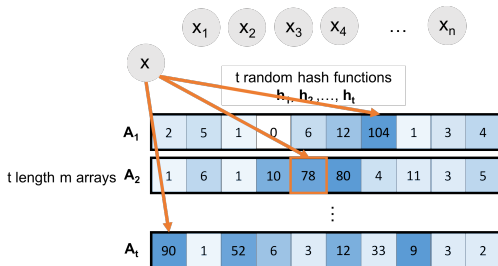
Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

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

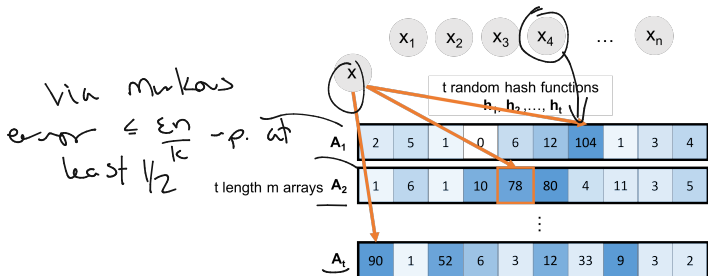
- For every x and $i \in [t]$, we know that for $m = \frac{2k}{\epsilon}$, with probability $\geq 1/2$:

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

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

$$\tilde{F}(x) \geq f(x) + \frac{\epsilon n}{k} \text{ only if } A_i[h_i(x)] \geq f(x) + \frac{\epsilon n}{k} \text{ for all } i \text{ arrays}$$

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 $\geq 1/2$:

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

$\geq k\delta, \delta$

- What is $\Pr[f(x) \leq \tilde{f}(x) \leq f(x) + \frac{\epsilon n}{k}]$? $1 - \frac{1}{2^t}$.
- To get a good estimate with probability $\geq 1 - \delta$, set $t = \log_2(1/\delta)$.

Count-Min Sketch

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.

$\ll n$

~~$O(\log(1/\delta) \cdot k/\epsilon)$~~ $m = \text{length of each array}$

↓
array 1

Count-Min Sketch

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 – distinguish between items with frequency $\frac{n}{k}$ and those with frequency $(1 - \epsilon)\frac{n}{k}$.

Count-Min Sketch

E_i event that estimate for item i is bad

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 – 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? $\log(1/\delta) = \log(100) = \frac{0.01}{n}$

union bound

$\delta = \frac{0.01}{n}$

$$\Pr(E_1 \cup E_2 \dots \cup E_n) \leq n \cdot \Pr(E_i) \leq n \cdot \delta \leq 0.01$$

Identifying Frequent Elements

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 store/look up the estimated frequency for all elements in the stream?

Identifying Frequent Elements

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 store/look up the estimated frequency for all elements in the stream?

don't know n ahead of
w. don't know $\frac{n}{k}$

One approach:

$\frac{2}{k}$

- When a new item comes in at step i , check if its estimated frequency is $\geq i/k$ and store it if so.
- At step i remove any stored items whose estimated frequency drops below $\underline{i/k}$.
- Store at most $O(k)$ items at once and have all items with frequency $\geq n/k$ stored at the end of the stream.

$$(1-\epsilon)\frac{i}{k}$$

$$(1-\epsilon)\frac{n}{k}, \text{ violate up to } \delta$$