

COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2022.

Lecture 10

- Problem Set 2 is due next Friday at 11:59pm.
- The midterm is in class on Thursday 10/20. Midterm study material will be posted shortly.
- We will have a quiz this week, but not next week.

Summary

Last Class:

- Locality sensitive hashing for near neighbor search.
 - MinHash as a locality sensitive hash function for Jaccard similarity
- $$J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$
- Start on balancing false positives and negatives with LSH signatures and repeated hash tables.

Summary

Last Class:

- Locality sensitive hashing for near neighbor search.
- MinHash as a locality sensitive hash function for Jaccard similarity
- Start on balancing false positives and negatives with LSH signatures and repeated hash tables.

This Class:

- Finish up LSH analysis.
- Frequent Items Estimation
- Count-min sketch algorithm

Upcoming

Next Few Classes:

- Random compression methods for high dimensional vectors.
The Johnson-Lindenstrauss lemma.
- Connections to the weird geometry of high-dimensional space.

Upcoming

Next Few Classes:

- Random compression methods for high dimensional vectors. The Johnson-Lindenstrauss lemma.
- Connections to the weird geometry of high-dimensional space.

After the Midterm: Spectral Methods

- PCA, low-rank approximation, and the singular value decomposition.
- Spectral clustering and spectral graph theory.

Upcoming

Next Few Classes:

- Random compression methods for high dimensional vectors. The Johnson-Lindenstrauss lemma.
- Connections to the weird geometry of high-dimensional space.

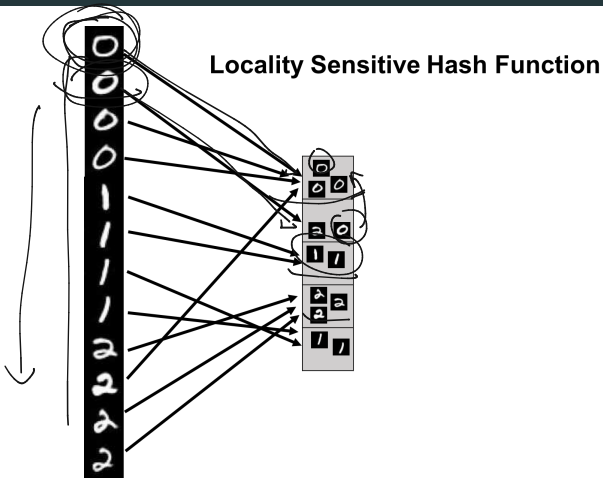
After the Midterm: Spectral Methods

- PCA, low-rank approximation, and the singular value decomposition.
- Spectral clustering and spectral graph theory.

Will use a lot of linear algebra. May be helpful to refresh.

- Vector dot product, addition, Euclidean norm. Matrix vector multiplication.
- Linear independence, column span, orthogonal bases, rank.
- Orthogonal projection, eigendecomposition, linear systems.

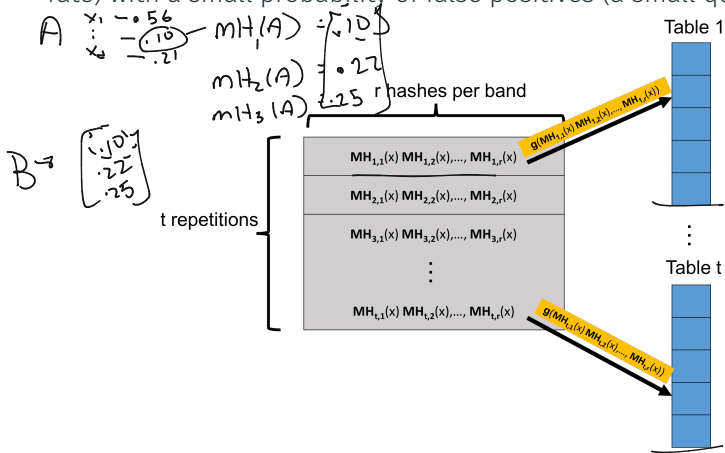
Locality Sensitive Hashing



Speed up near neighbor search by looking only for nearby items that land in the same buckets.

Balancing Hit Rate and Query Time

We want to balance a small probability of false negatives (a high hit rate) with a small probability of false positives (a small query time.)



Create t hash tables. Each is indexed into not with a single MinHash value, but with r values, appended together. A length r **signature**.

Balancing Hit Rate and Query Time

Balancing Hit Rate and Query Time

Consider searching for matches in t hash tables, using MinHash signatures of length r . For x and y with Jaccard similarity $J(x, y) = s$:

$[0, 1]$

- Probability that a single hash matches.

$\Pr [MH_{i,j}(x) = MH_{i,j}(y)] = J(x, y) = s.$

not \cap

- Probability that x and y having matching signatures in repetition i . $\Pr [MH_{i,1}(x), \dots, MH_{i,r}(x) = MH_{i,1}(y), \dots, MH_{i,r}(y)] = s^r.$
- Probability that x and y don't match in repetition i : $1 - s^r.$
- Probability that x and y don't match in *all repetitions*: $(1 - s^r)^t.$
- Probability that x and y match in at least one repetition:

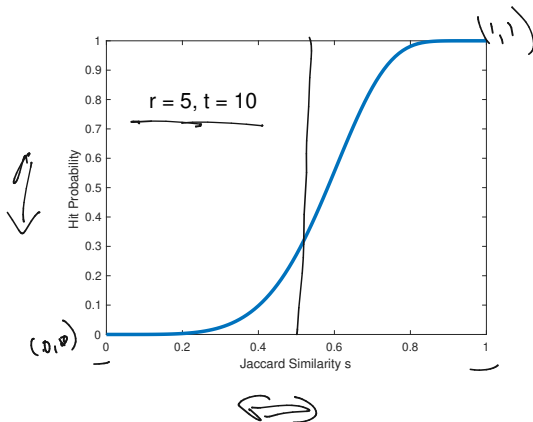
Hit Probability: $1 - (1 - s^r)^t.$

The s-curve

Using t repetitions each with a signature of r MinHash values, the probability that x and y with Jaccard similarity $J(x,y) = s$ match in at least one repetition is: $1 - (1 - s^r)^t$.

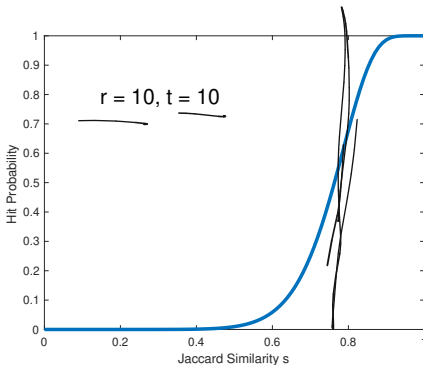
The s-curve

Using t repetitions each with a signature of r MinHash values, the probability that x and y with Jaccard similarity $J(x,y) = s$ match in at least one repetition is: $1 - (1 - s^r)^t$.



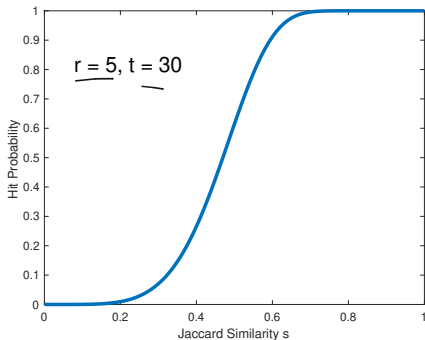
The s-curve

Using t repetitions each with a signature of r MinHash values, the probability that x and y with Jaccard similarity $J(x,y) = s$ match in at least one repetition is: $1 - (1 - s^r)^t$.



The s-curve

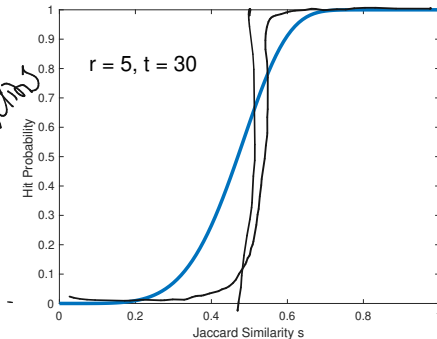
Using t repetitions each with a signature of r MinHash values, the probability that x and y with Jaccard similarity $J(x, y) = s$ match in at least one repetition is: $1 - (1 - s^r)^t$.



The s-curve

Using t repetitions each with a signature of r MinHash values, the probability that x and y with Jaccard similarity $J(x, y) = s$ match in at least one repetition is: $1 - (1 - s^r)^t$.

$t = \# \text{ tables}$
 $r = \# \text{ hash functions}$



r and t are tuned depending on application. 'Threshold' when hit probability is $1/2$ is $\approx (1/t)^{1/r}$. E.g., $\approx (1/30)^{1/5} = .51$ in this case.

s-curve Example

For example: Consider a database with 10,000,000 audio clips. You are given a clip x and want to find any y in the database with

$$\underline{J(x, y) \geq .9.}$$

s-curve Example

For example: Consider a database with 10,000,000 audio clips. You are given a clip x and want to find any y in the database with $J(x, y) \geq .9$.

- There are 10 **true matches** in the database with $J(x, y) \geq .9$.
- There are 10,000 **near matches** with $J(x, y) \in [.7, .9]$.

s-curve Example

For example: Consider a database with 10,000,000 audio clips. You are given a clip x and want to find any y in the database with $J(x, y) \geq .9$.

- There are 10 **true matches** in the database with $J(x, y) \geq .9$.
- There are 10,000 **near matches** with $J(x, y) \in [.7, .9]$.

With signature length $r = 25$ and repetitions $t = 50$, hit probability for $J(x, y) = s$ is $\underline{1 - (1 - s^{25})^{50}}$.

s-curve Example

For example: Consider a database with 10,000,000 audio clips. You are given a clip x and want to find any y in the database with $J(x, y) \geq .9$.

- There are 10 true matches in the database with $J(x, y) \geq .9$.
- There are 10,000 near matches with $J(x, y) \in [.7, .9]$.

With signature length $r = 25$ and repetitions $t = 50$, hit probability for $J(x, y) = s$ is $1 - (1 - s^{25})^{50}$.

- Hit probability for $J(x, y) \geq .9$ is $\geq 1 - (1 - .9^{25})^{50} \approx .98$
- Hit probability for $J(x, y) \in [.7, .9]$ is $\leq 1 - (1 - .9^{25})^{50} \approx .98$
- Hit probability for $J(x, y) \leq .7$ is $\leq 1 - (1 - .7^{25})^{50} \approx .007$

s-curve Example

For example: Consider a database with 10,000,000 audio clips. You are given a clip x and want to find any y in the database with $J(x,y) \geq .9$.

- There are 10 **true matches** in the database with $J(x,y) \geq .9$.
- There are 10,000 **near matches** with $J(x,y) \in [.7, .9]$.

With signature length $r = 25$ and repetitions $t = 50$, hit probability for $J(x,y) = s$ is $1 - (1 - s^{25})^{50}$.

- Similar items*
- Hit probability for $J(x,y) \geq .9$ is $\geq 1 - (1 - .9^{25})^{50} \approx .98$
 - Hit probability for $J(x,y) \in [.7, .9]$ is $\leq 1 - (1 - .9^{25})^{50} \approx .98$
 - Hit probability for $J(x,y) \leq .7$ is $\leq 1 - (1 - .7^{25})^{50} \approx .007$
- 0.999*

Expected Number of Items Scanned: (proportional to query time)

$$\leq 10 + .98 * 10,000 + .007 * 9,989,990 \approx 80,000$$

s-curve Example

For example: Consider a database with 10,000,000 audio clips. You are given a clip x and want to find any y in the database with $J(x, y) \geq .9$.

- There are 10 **true matches** in the database with $J(x, y) \geq .9$.
- There are 10,000 **near matches** with $J(x, y) \in [.7, .9]$.

With signature length $r = 25$ and repetitions $t = 50$, hit probability for $J(x, y) = s$ is $1 - (1 - s^{25})^{50}$.

- Hit probability for $J(x, y) \geq .9$ is $\geq 1 - (1 - .9^{25})^{50} \approx .98$
- Hit probability for $J(x, y) \in [.7, .9]$ is $\leq 1 - (1 - .9^{25})^{50} \approx .98$
- Hit probability for $J(x, y) \leq .7$ is $\leq 1 - (1 - .7^{25})^{50} \approx .007$

Expected Number of Items Scanned: (proportional to query time)

$$\leq 10 + .98 * 10,000 + .007 * 9,989,990 \approx 80,000 \ll 10,000,000.$$

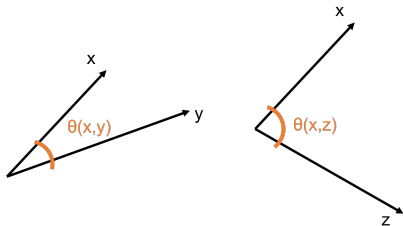
Hashing for Duplicate Detection

	Hash Table	Bloom Filters	^{LSH} MinHash Similarity Search	Distinct Elements
Goal	Check if x is a duplicate of any y in database <u>and return y.</u>	Check if x is a duplicate of y in database.	Check if x is a duplicate of any y in database <u>and return y.</u>	Count # of items, excluding duplicates.
Space	<u>$O(n)$ items</u>	$O(n)$ bits	$O(n \cdot t)$ items (when t tables used)	$O\left(\frac{\log \log n}{\epsilon^2}\right)$
Query Time	$O(1)$	$O(1)$	Potentially $o(n)$	NA
Approximate Duplicates?	✘	✘	✔	✘

All different variants of detecting duplicates/finding matches in large datasets. An important problem in many contexts!

SimHash for Cosine Similarity

Repetition and s-curve tuning can be used for fast similarity search with any similarity metric, given a locality sensitive hash function for that metric.



Cosine Similarity: $\cos(\theta(x, y)) = \frac{\langle x, y \rangle}{\|x\|_2 \cdot \|y\|_2}$.

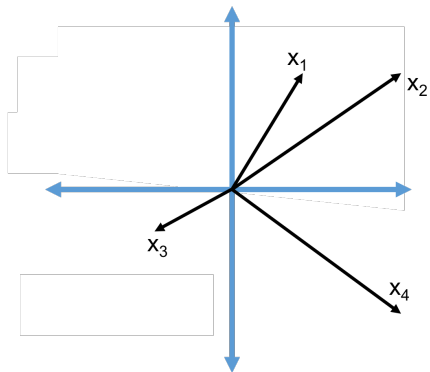
- $\cos(\theta(x, y)) = 1$ when $\theta(x, y) = 0^\circ$ and $\cos(\theta(x, y)) = 0$ when $\theta(x, y) = 90^\circ$, and $\cos(\theta(x, y)) = -1$ when $\theta(x, y) = 180^\circ$.

SimHash for Cosine Similarity

SimHash: LSH for cosine similarity.

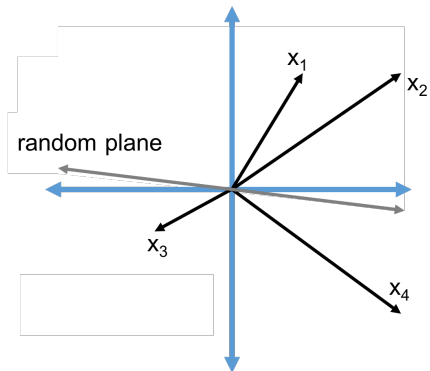
SimHash for Cosine Similarity

SimHash: LSH for cosine similarity.



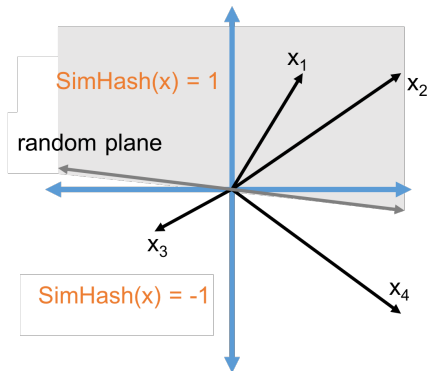
SimHash for Cosine Similarity

SimHash: LSH for cosine similarity.



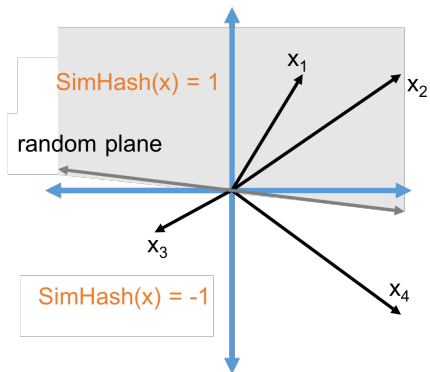
SimHash for Cosine Similarity

SimHash: LSH for cosine similarity.



SimHash for Cosine Similarity

SimHash: LSH for cosine similarity.



$\text{SimHash}(x) = \text{sign}(\langle x, t \rangle)$ for a random vector t .

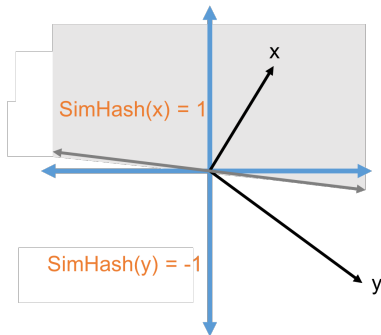
SimHash for Cosine Similarity

What is $\Pr[\text{SimHash}(x) = \text{SimHash}(y)]$?

SimHash for Cosine Similarity

What is $\Pr[\text{SimHash}(x) = \text{SimHash}(y)]$?

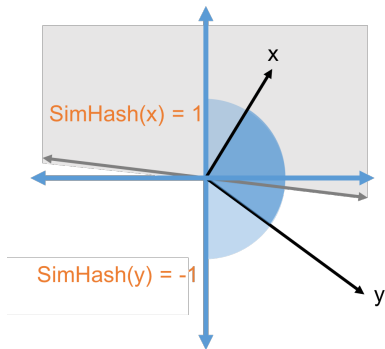
$\text{SimHash}(x) \neq \text{SimHash}(y)$ when the plane separates x from y .



SimHash for Cosine Similarity

What is $\Pr[\text{SimHash}(x) = \text{SimHash}(y)]$?

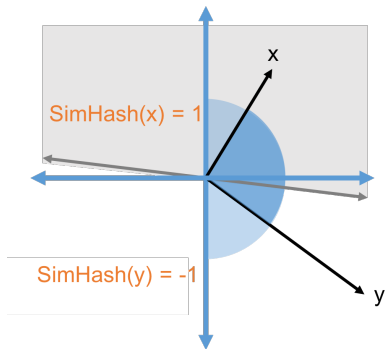
$\text{SimHash}(x) \neq \text{SimHash}(y)$ when the plane separates x from y .



SimHash for Cosine Similarity

What is $\Pr[\text{SimHash}(x) = \text{SimHash}(y)]$?

$\text{SimHash}(x) \neq \text{SimHash}(y)$ when the plane separates x from y .

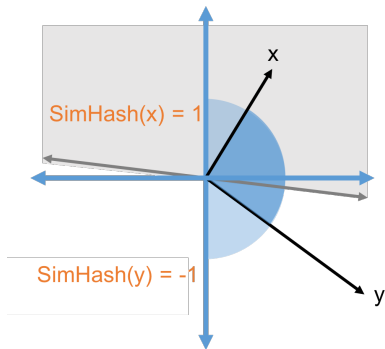


$$\cdot \Pr[\text{SimHash}(x) \neq \text{SimHash}(y)] = \frac{\theta(x,y)}{\pi}$$

SimHash for Cosine Similarity

What is $\Pr[\text{SimHash}(x) = \text{SimHash}(y)]$?

$\text{SimHash}(x) \neq \text{SimHash}(y)$ when the plane separates x from y .



- $\Pr[\text{SimHash}(x) \neq \text{SimHash}(y)] = \frac{\theta(x,y)}{\pi}$
- $\Pr[\text{SimHash}(x) = \text{SimHash}(y)] = 1 - \frac{\theta(x,y)}{\pi} \approx \frac{\cos(\theta(x,y))+1}{2}$

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.

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.

3 - frequent

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.

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) $\log n$

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.

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

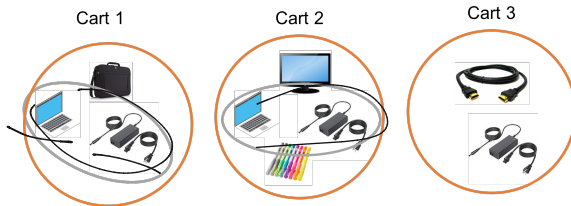
Association rule learning: A very common task in data mining is to identify common associations between different events.



- 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 $\underline{o(n)}$ 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 $\underline{n/k - 1}$ (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**.

$n = 1000$
 $k = 10$
 $\epsilon = .1$

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**.

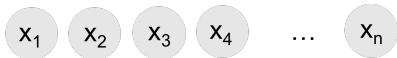
- An example of relaxing to a 'promise problem': for items with frequencies in $[(1 - \epsilon) \cdot \frac{n}{k}, \frac{n}{k}]$ no output guarantee.

Frequent Elements with Count-Min Sketch

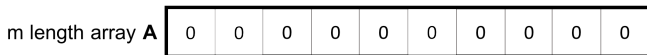
Today: Count-min sketch – a random hashing based method closely related to bloom filters.

Frequent Elements with Count-Min Sketch

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

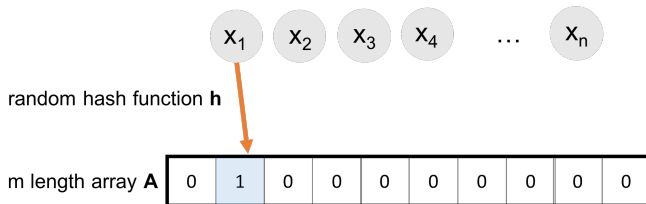


random hash function h



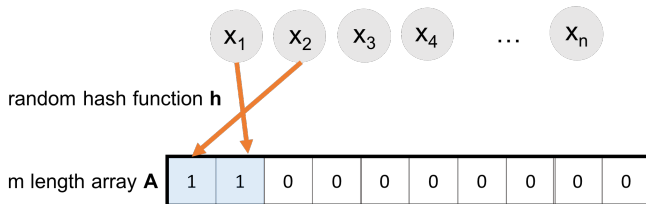
Frequent Elements with Count-Min Sketch

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



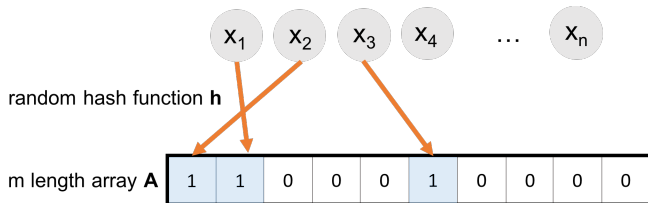
Frequent Elements with Count-Min Sketch

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



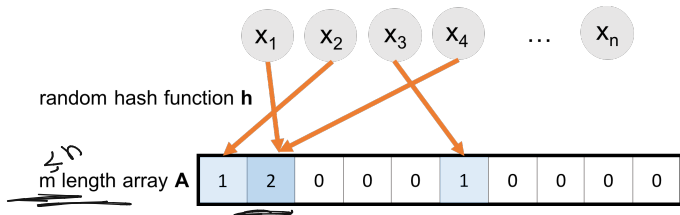
Frequent Elements with Count-Min Sketch

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



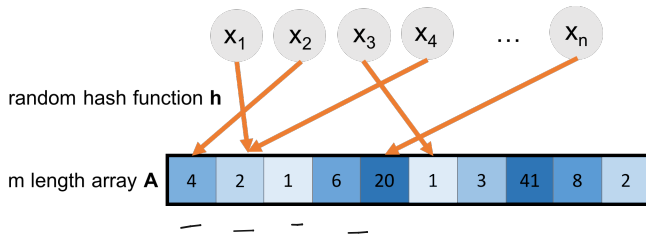
Frequent Elements with Count-Min Sketch

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



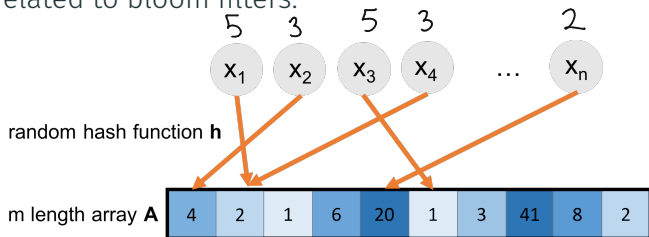
Frequent Elements with Count-Min Sketch

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



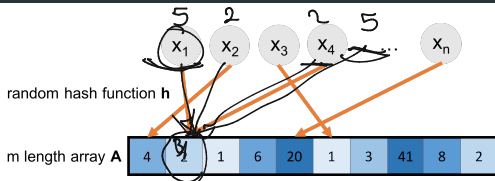
Frequent Elements with Count-Min Sketch

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



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

Count-Min Sketch Accuracy

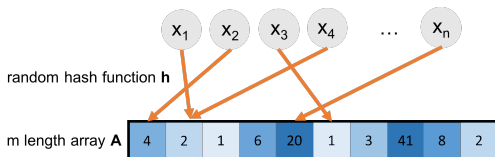


Use $A[h(x)]$ to estimate $f(x)$.

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



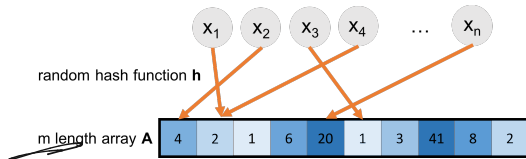
Use $A[h(x)]$ to estimate $f(x)$.

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

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

$f(x)$: 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

m counters

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

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:

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

$f(x)$: 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:

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

$f(x)$: 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

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

Expected Error:

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

$f(x)$: 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

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

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 $\geq \frac{2n}{m}$?

$f(x)$: 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: $\Pr \left[\underbrace{\sum_{y \neq x: h(y)=h(x)} f(y)}_{\text{error}} \geq \frac{2n}{m} \right] \leq \frac{1}{2}$.

$f(x)$: 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:

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

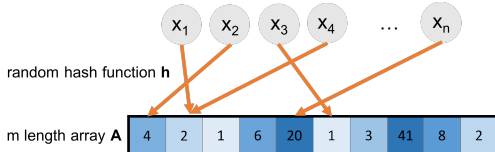
What is a bound on probability that the error is $\geq \frac{2n}{m}$?

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

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

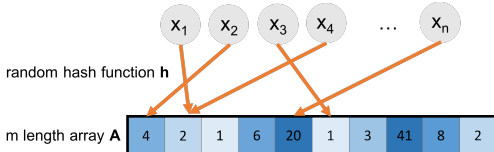


Claim: For any x , with probability at least $1/2$,

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

$f(x)$: 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 $1/2$,

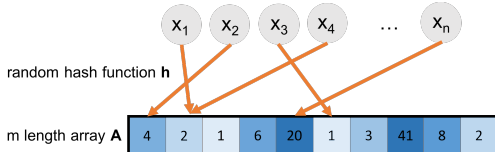
$$f(x) \leq A[h(x)] \leq f(x) + \frac{2n}{m}.$$

$$\frac{2n}{m} = \epsilon \cdot \frac{1}{\epsilon}$$

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



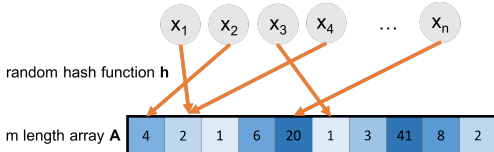
Claim: For any x , with probability at least $1/2$,

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

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

$f(x)$: 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



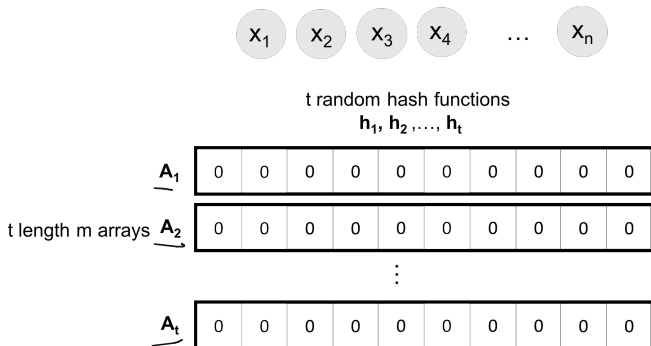
Claim: For any x , with probability at least $1/2$,

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

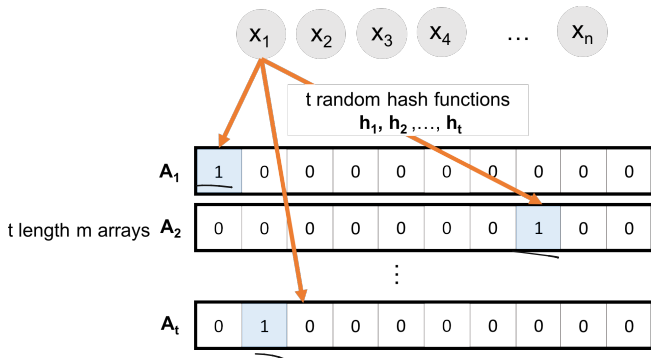
To solve the (ϵ, k) -Frequent elements problem, set $m = \frac{2k}{\epsilon}$. How can we improve the success probability? Repetition.

$f(x)$: 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.

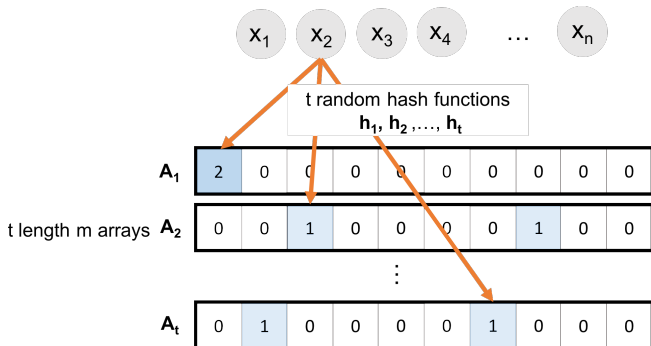
Count-Min Sketch Repetition



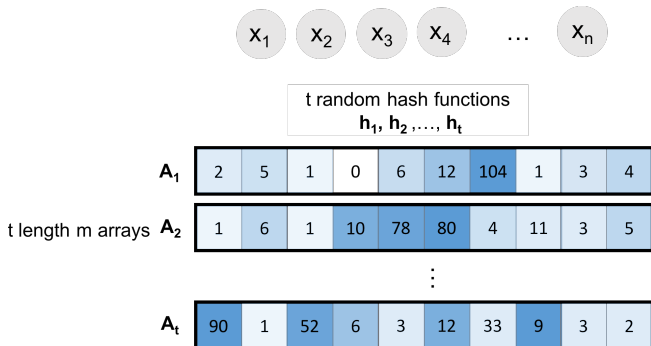
Count-Min Sketch Repetition



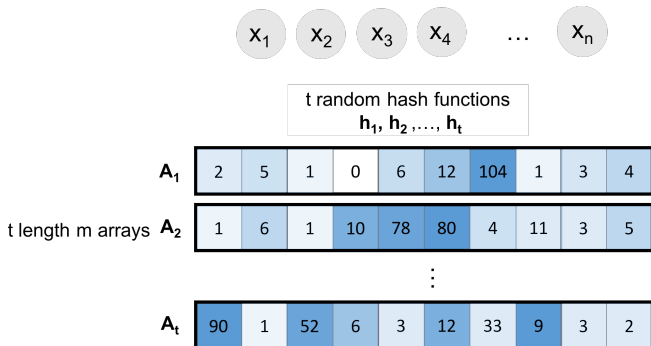
Count-Min Sketch Repetition



Count-Min Sketch Repetition

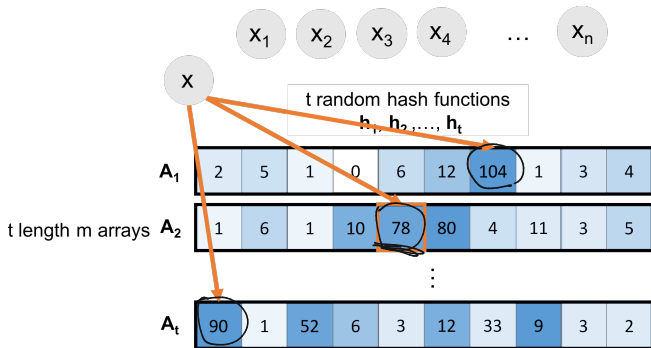


Count-Min Sketch Repetition



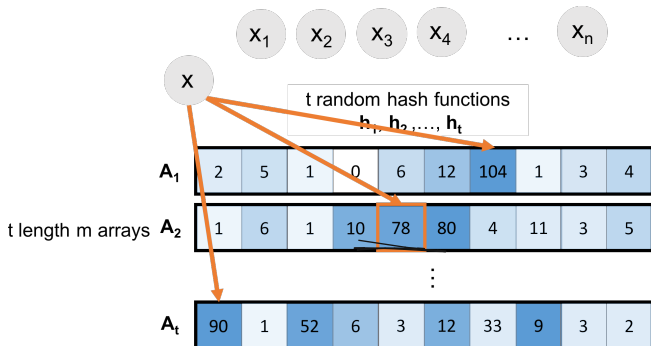
Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$. (count-min sketch)

Count-Min Sketch Repetition



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

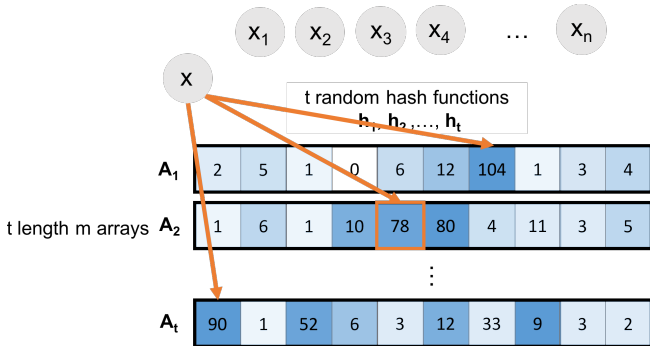
Count-Min Sketch Repetition



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?

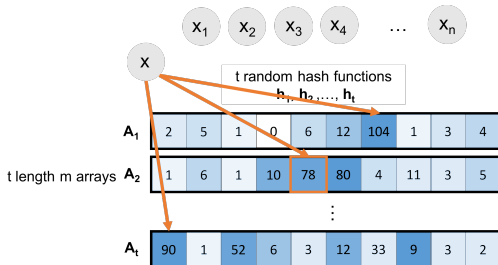
Count-Min Sketch Repetition



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

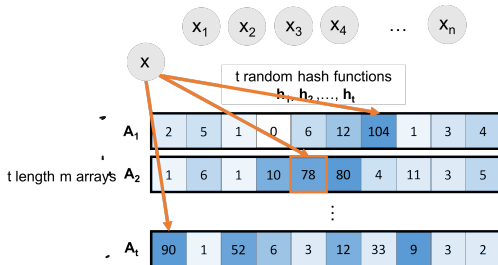
Why min instead of mean or median? The minimum estimate is always the most accurate since they are all overestimates of the true frequency!

Count-Min Sketch Analysis



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

Count-Min Sketch Analysis

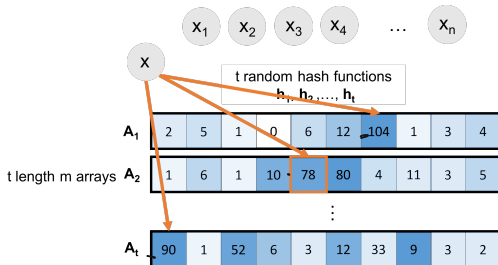


Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

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

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

Count-Min Sketch Analysis



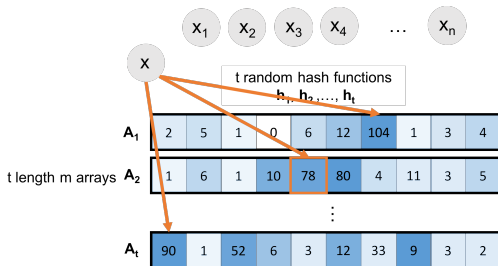
Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- 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



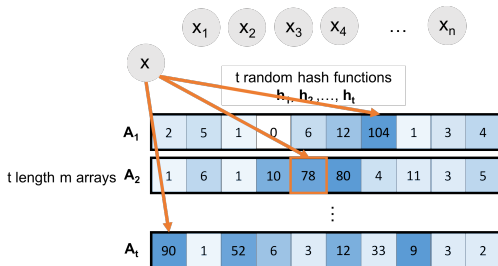
Estimate $f(x)$ by $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$

- 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 $\geq 1/2$:

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

prob all estimates "below"

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

Count-Min Sketch

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

Space $\approx m$

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

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

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

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?

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?

One approach:

- 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 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.