## COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2020.
Lecture 9

- Problem Set 2 is due this upcoming Monday. Get an early start on it.
- Problem Set 1 grades have been released. Mean: 34/41, Median 36/41.
- If you are unhappy with your grade, ping me and let's chat about strategies going forward. If you believe there is a grading error, send a private message to the instructors on Piazza or ask during office hours.
- The midterm will be any 2 hour slot on 10/8-10/9. We won't have class on 10/8.
- Study guide/practice questions will be released this week.

Last Class:

- MinHash as a locality sensitive hash function for Jaccard similarity
- Near neighbor search with LSH signatures and repeated hash tables..
- SimHash for cosine similarity.

**Last Class:**

$g(mH(A))$ → $\boxed{\quad}$

$Pr\left(mH(A) = mH(B)\right) = J(A,B)$

$Pr\left(g(mH(A)) = g(mH(B))\right) = J(A,B)$

- MinHash as a locality sensitive hash function for Jaccard similarity
- Near neighbor search with LSH signatures and repeated hash tables..

  $h$ is a random hash function

  $Pr\left(h(x) = h(y)\right) = \frac{1}{m}$

- SimHash for cosine similarity.

  $h(x)$
  $h(x)$

A locality sensitive hash function can be: (check all that apply)

**Select one or more:**

- ☐ a. Randomized
- ☒ b. Pairwise-Independent
- ☐ c. Sensitive to Jaccard Similarity
- ☐ d. Have the distribution of h(x) independent of x

pairwise ind. $Pr\left(h(x) = h(y) = i\right) = \frac{1}{n} \cdot \frac{1}{n}$ for all $x, y$

(minHash (x), simHash(x))

2

$$A = \{ a_1, a_2 \ldots a_n \}$$

**Last Class:**

$$MH(A) = \min_i h(a_i)$$

$a_i$ will be 1 shingle

- MinHash as a locality sensitive hash function for Jaccard similarity
- Near neighbor search with LSH signatures and repeated hash tables..
- SimHash for cosine similarity.

**This Class:** Frequent Items Estimation

- Count-min sketch (random hashing for frequent element estimation).

**Next Few Classes:**

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

### Next Few Classes:

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

### After That: Spectral Methods

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

### Next Few Classes:

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

### After That: 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.

3

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

*k*-Frequent Items (Heavy-Hitters) Problem: Consider a stream of $n$ items $x_1, \ldots, x_n$ (with possible duplicates). Return any item at 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 |

$k$-**Frequent Items (Heavy-Hitters) Problem**: Consider a stream of $n$ items $x_1, \ldots, x_n$ (with possible duplicates). Return any item at 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 |

$k$-**Frequent Items (Heavy-Hitters) Problem**: Consider a stream of $n$ items $x_1, \ldots, x_n$ (with possible duplicates). Return any item at 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~~ *can* be returned?

   a)   $n$    b)   $k$    c)   $n/k$    d)   $\log n$

\#HHs $\leq k$

each heavy hitter appears $\frac{n}{k}$ times

so total frequency is $\frac{n}{k} \cdot \#$ HHs $\leq n$

4

$k$-**Frequent Items (Heavy-Hitters) Problem**: Consider a stream of $n$ items $x_1, \ldots, x_n$ (with possible duplicates). Return any item at 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)   $\log n$

$k$-**Frequent Items (Heavy-Hitters) Problem**: Consider a stream of $n$ items $x_1, \ldots, x_n$ (with possible duplicates). Return any item at 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. I.e., want to maintain a running list of frequent items that appear in a stream.

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

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

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



Cart 1      Cart 2      Cart 3

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



Cart 1  Cart 2  Cart 3

- Identified via frequent itemset counting. Find all sets of $k$ items that appear many times in the same basket.

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



Cart 1　　Cart 2　　Cart 3

- Identified via frequent itemset counting. Find all sets of $k$ items that appear many times in the same basket.
- Frequency of an itemset is known as its support.

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



Cart 1    Cart 2    Cart 3

- Identified via frequent itemset counting. Find all sets of $k$ 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.

**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).

$o(n)$

**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

**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

$(\epsilon, k)$-**Frequent Items Problem**: Consider a stream of $n$ items $x_1, \ldots, x_n$. 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.

$$\epsilon = .1 \qquad k = 100$$

**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

$(\epsilon, k)$-**Frequent Items Problem**: Consider a stream of $n$ items $x_1, \ldots, x_n$. 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.

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

$$x_1 \quad x_2 \quad x_3 \quad x_4 \quad \ldots \quad x_n$$

random hash function **h**

m length array **A** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

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

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

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

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

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

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



random hash function **h**

m length array **A**

$A[2] = 2$

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

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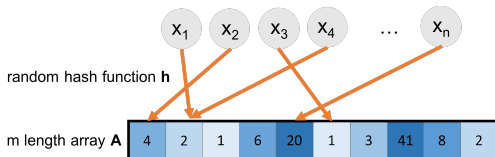
Use $A[h(x)]$ to estimate $f(x)$.

**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.

Use $A[\mathsf{h}(x)]$ to estimate $f(x)$.

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

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

*(handwritten annotations: "random error" above the sum; "Fixed" below $f(x)$)*

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

9

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

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

## COUNT-MIN SKETCH ACCURACY

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

**Expected Error:**

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

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

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

**Expected Error:**

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

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

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

m buckets counters

**Expected Error:**

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

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

*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.

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

**Expected Error:**

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

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

$$\frac{1}{m} \sum_{y \neq x} f(y)$$

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

10

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

3   3   5  6  3  10
X = 3

**Expected Error:**

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

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

$\sum_{y \neq x} f(y) = n - f(x) \leq n$

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.

10

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

**Expected Error:**

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

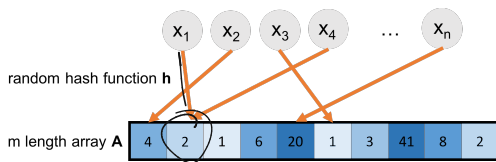$$= \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m}$$

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

**Markov's inequality:** $\Pr\left[\underbrace{\sum_{y \neq x: \mathbf{h}(y) = \mathbf{h}(x)} f(y)}_{\alpha \cdot \mathbb{E}(\text{error})} \geq \frac{2n}{m}\right] \leq \frac{1}{2}.$

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

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

**Expected Error:**

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

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

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

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

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.

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

**Expected Error:**

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

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

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

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

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.

10

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

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

$f(x)$: 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.
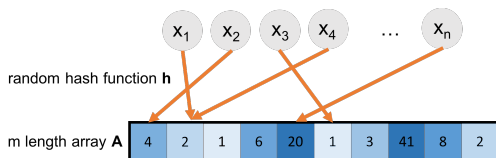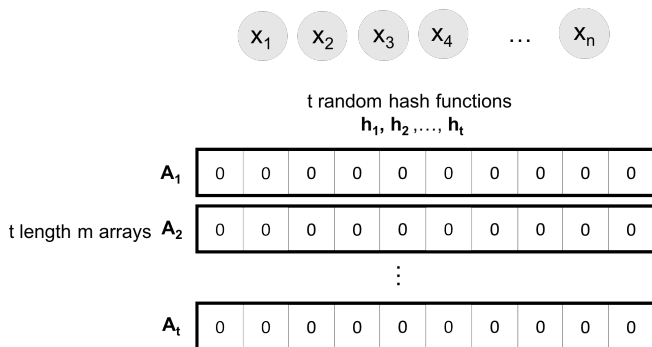
11

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

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

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

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

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

$f(x)$: 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.

11

**Claim:** For any *x*, with probability at least 1/2,
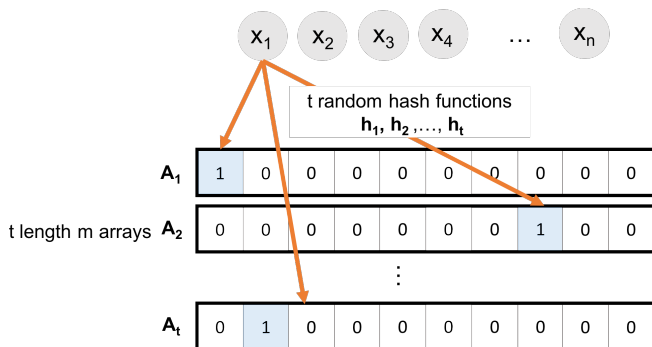
$$f(x) \leq A[\mathsf{h}(x)] \leq f(x) + \frac{2n}{m}.$$

To solve the $(\epsilon, k)$-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*). $\mathsf{h}$: random hash function. *m*: size of Count-min sketch array.
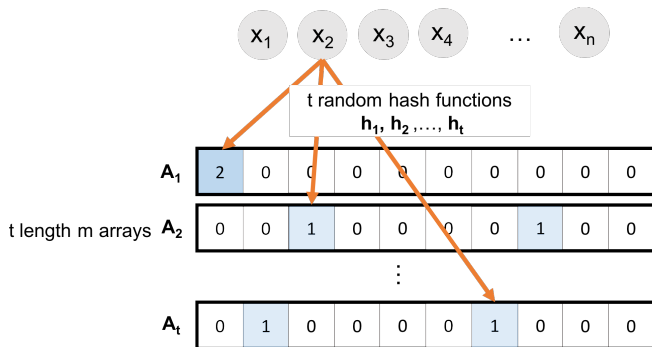
**Claim:** For any *x*, with probability at least 1/2,

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

To solve the $(\epsilon, k)$-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.

$x_1$ $x_2$ $x_3$ $x_4$ … $x_n$

t random hash functions
$h_1, h_2, \ldots, h_t$

$A_1$

| 2 | 5 | 1 | 0 | 6 | 12 | 104 | 1 | 3 | 4 |

t length m arrays $A_2$

| 1 | 6 | 1 | 10 | 78 | 80 | 4 | 11 | 3 | 5 |

$\vdots$

$A_t$

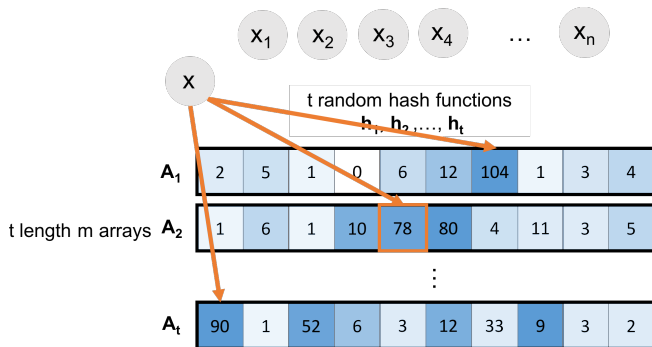| 90 | 1 | 52 | 6 | 3 | 12 | 33 | 9 | 3 | 2 |

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

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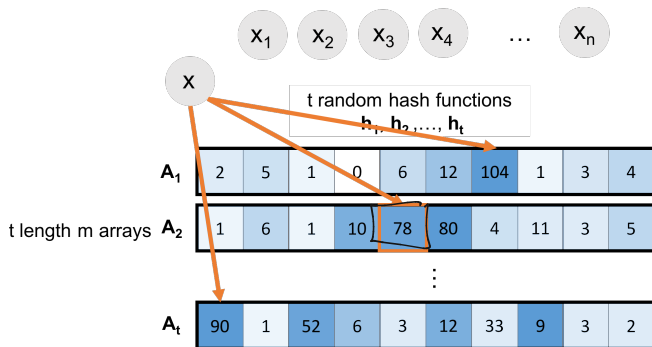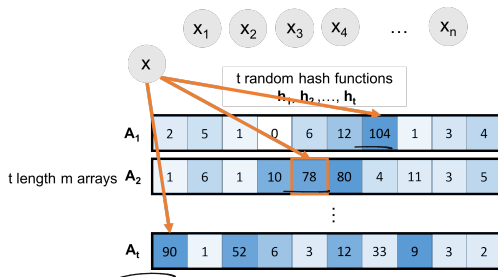
Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$. (count-min sketch)

Why min instead of mean or median?

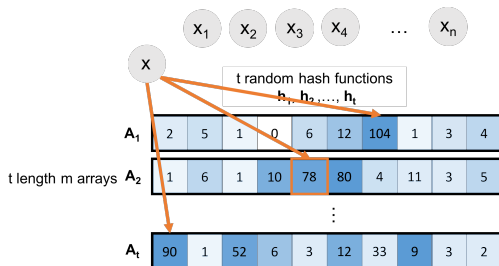Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[\mathbf{h}_i(x)]$. (count-min sketch)

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

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

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

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

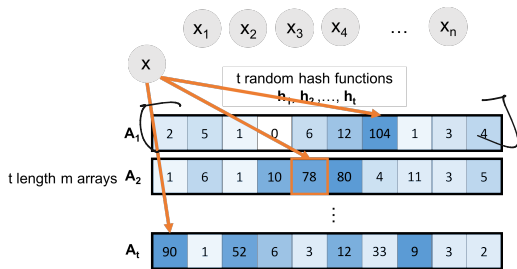$$f(x) \leq A_i[\mathbf{h}_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

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

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

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

$$f(x) \leq \underbrace{A_i[\mathbf{h}_i(x)]} \leq f(x) + \frac{\epsilon n}{k}.$$

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

$$Pr\left( \tilde{f}(x) \notin \left[ f(x), f(x) + \frac{\epsilon n}{k} \right] \right) = \frac{1}{2^t}$$
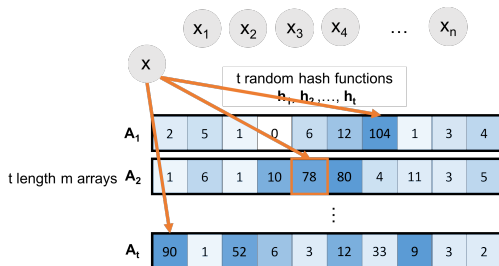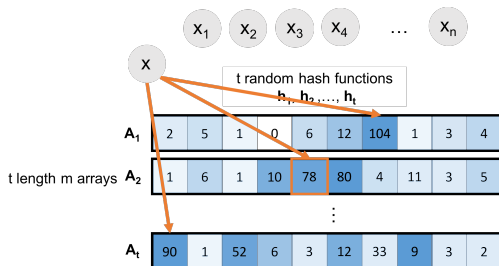
13

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 $t = \log(1/\delta)$.

13

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 $O\left(\log(1/\delta) \cdot k/\epsilon\right)$ 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 $O\left(\log(1/\delta) \cdot k/\epsilon\right)$ space.

· Accurate enough to solve the $(\epsilon, k)$-Frequent elements problem – distinquish between items with frequency $\frac{n}{k}$ and those with frequency $(1 - \epsilon)\frac{n}{k}$.      $\frac{\epsilon n}{k}$

set $\delta = .01$

$t = \log(1/\delta) = \delta(\log n)$

apply union bound.

$f(x_i) < \hat{f}(x_i) \leq f(x_i) + \frac{\epsilon n}{k}$ w.p $1-\delta$

For all $i$ at once, $f(x_i) \leq \hat{f}(x_i) \leq f(x_i) + \frac{\epsilon n}{k}$ w.p $1-n\delta$

**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 $O\left(\log(1/\delta) \cdot k/\epsilon\right)$ space.

$O\left(\log(n) \cdot \frac{k}{\epsilon}\right)$

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

· How should we set $\delta$ if we want a good estimate for all items at once, with 99% probability?

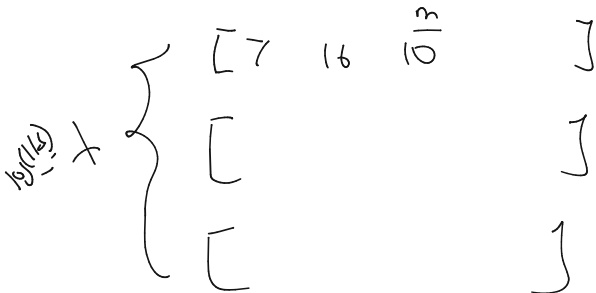$n$ possible items, we want to estimate all frequencies.

let $E_i$ be the event that I **fail** to estimate $f(x_i)$ well

$Pr(E_i) \leq \delta$      $Pr(E_1 \cup E_2 \cup \ldots, E_n) \leq \sum_{i=1}^{n} Pr(E_i) \leq n\delta$

14

$$O\left(\log n \cdot \frac{k}{\varepsilon}\right) \text{ space.} \qquad \frac{n}{k} \quad (1-\varepsilon)\frac{n}{k}$$

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?

$$m = \frac{k}{\varepsilon} \qquad \frac{\varepsilon n}{k}$$

$$\log(1/\delta) \times \left\{ \begin{array}{c} [\, 7 \quad 16 \quad \frac{m}{10} \qquad ] \\[2em] [ \qquad\qquad\qquad ] \\[2em] [ \qquad\qquad\qquad ] \end{array} \right.$$

total space =
$$t \cdot m = \log(1/\delta) \cdot m$$
$$= \log(1/\delta) \cdot \frac{k}{\varepsilon}$$

$\delta$ = failure prob.        $(\varepsilon, k)$ are parameters of freq items prob.

15

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: $X_1, X_2 . X_i, \quad X_n$

- 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.

Time complexity: $O(k)$ time to hash $x_i$
an increment buckets

$O(k)$ time to check if any of stored
items have freq $\leq \frac{i}{k}$

$O\left(t \cdot \frac{k}{\epsilon}\right)$

$+ t \cdot$ size hash function.

- Defined frequent items

- Relaxed it to the $(\epsilon, k)$-Freq. item

- used count-min sketch (variant on Bloom filters)
to solve this problem in $O\left(\log n \cdot \frac{k}{\epsilon}\right)$ space

$\ll O(n)$

16