# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2024.

Lecture 8

## Summary

**Last Class:**

- Finish up Bloom Filters and optimization of number of hash functions.

- Start on streaming algorithms.

- Introduce the frequent items problem and its applications.

- Start on the Count-Min sketch algorithm for frequent items.

**This Class:**

- Analysis of Count-Min sketch .
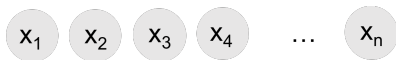
- Start on distinct items counting problem.

$(\epsilon, k)$-**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.

- To solve this problem, it suffices to estimate the frequency $f(x)$ of each item $x$ up to error $\pm \frac{\epsilon n}{k}$.

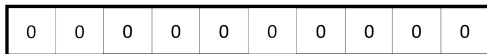- Will discuss later how to maintain the list of top items in small space.

Count-min sketch:

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

random hash function **h**                                                        random hash fur

m length array **A**

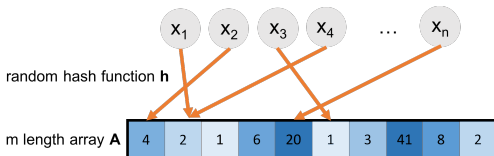| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|

m length array **A**

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

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

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

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

$f(x)$: 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.

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

**Expected Error:**

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

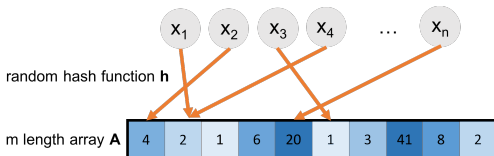$$= \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m}$$

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

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

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,

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

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

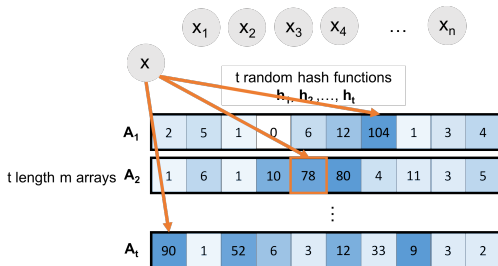> $f(x)$: frequency of $x$ in the stream (i.e., number of items equal to $x$). $h$: random hash function. $m$: size of Count-min sketch array.

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

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

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

- 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 $t = \log_2(1/\delta)$.

**Upshot:** Count-min sketch lets us estimate the frequency of each 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 $(1 - \epsilon)\frac{n}{k}$.
- How should we set $\delta$ 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?

### 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.

Questions on Frequent Items?

## Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream $x_1, \ldots, x_n$, estimate the number of distinct elements in the stream. E.g.,

$$1, 5, 7, 5, 2, 1 \rightarrow 4 \text{ distinct elements}$$

### Applications:

- Distinct IP addresses clicking on an ad or visiting a site.

- Distinct values in a database column (for estimating sizes of joins and group bys).

- Number of distinct search engine queries.

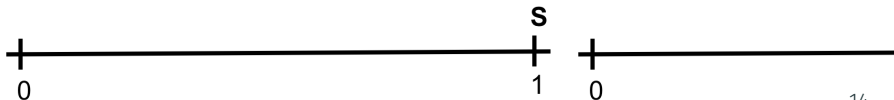- Counting distinct motifs in large DNA sequences.

Google Sawzall, Facebook Presto, Apache Drill, Twitter Algebird

Distinct Elements (Count-Distinct) Problem: Given a stream $x_1, \ldots, x_n$, estimate the number of distinct elements.

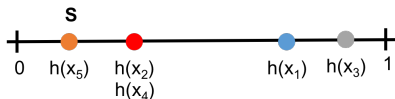Min-Hashing for Distinct Elements (variant of Flajolet-Martin):

- Let $h : U \to [0, 1]$ be a random hash function (with a real valued output)
- $s := 1$
- For $i = 1, \ldots, n$
  - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

**s**

0                                           1   0
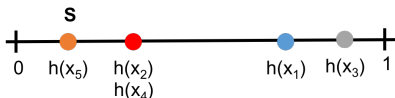
Min-Hashing for Distinct Elements:

- Let $h : U \to [0, 1]$ be a random hash function (with a real valued output)
- $s := 1$
- For $i = 1, \ldots, n$
    - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$



- After all items are processed, $s$ is the minimum of $d$ points chosen uniformly at random on $[0, 1]$. Where $d = \#$ distinct elements.

- Intuition: The larger $d$ is, the smaller we expect $s$ to be.

- Same idea as Flajolet-Martin algorithm and HyperLogLog, except they use discrete hash functions.

# Performance in Expectation

$s$ is the minimum of $d$ points chosen uniformly at random on $[0, 1]$. Where $d = \#$ distinct elements.



$$\mathbb{E}[s] = \frac{1}{d+1} \ \left(\text{using } \mathbb{E}(s) = \int_0^\infty \Pr(s > x)dx\right) + \text{calculus}$$

- So our estimate $\widehat{d} = \frac{1}{s} - 1$ is correct if $s$ exactly equals its expectation. Does this mean $\mathbb{E}[\widehat{d}] = d$? No, but:

- **Approximation is robust:** if $|s - \mathbb{E}[s]| \leq \epsilon \cdot \mathbb{E}[s]$ for any $\epsilon \in (0, 1/2)$ and a small constant $c \leq 4$:

$$(1 - c\epsilon)d \leq \widehat{d} \leq (1 + c\epsilon)d$$