# 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 .
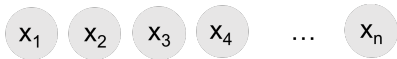- Start on distinct items counting problem.

($\epsilon, k$)-**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.
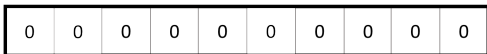
Count-min sketch:

Count-min sketch:

$x_1$ $x_2$ $x_3$ $x_4$ ... $x_n$

random hash function **h**

m length array **A**

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

Count-min sketch:

Count-min sketch:



random hash function **h**

m length array **A**

Count-min sketch:



$x_1$ $x_2$ $x_3$ $x_4$ ... $x_n$

random hash function **h**

m length array **A**

| 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |

Count-min sketch:

Count-min sketch:

Count-min sketch:



random hash function **h**

m length array **A**

| 4 | 2 | 1 | 6 | 20 | 1 | 3 | 41 | 8 | 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 $A[\mathbf{h}(x)]$ 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.

> $f(x)$: 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



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.

## Count-Min Sketch Accuracy

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

$f(x)$: 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

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

**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:**

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

$f(x)$: 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

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

**Expected Error:**

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

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

---

$f(x)$: 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

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

**Expected Error:**

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

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

---

$f(x)$: 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.

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

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

$f(x)$: 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.

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

**Expected Error:**

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

$$= \sum_{y \neq x} \frac{1}{m} \cdot f(y) = \frac{1}{m} \cdot (n - f(x)) \leq \frac{n}{m}$$
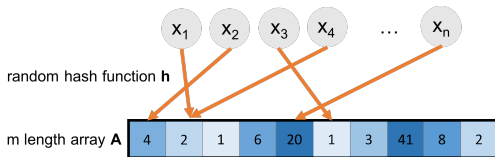
What is a bound on probability that the error is $\geq \frac{2n}{m}$?

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

$f(x)$: 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

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

**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 $\mathsf{h}$ 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 $x$). $\mathsf{h}$: random
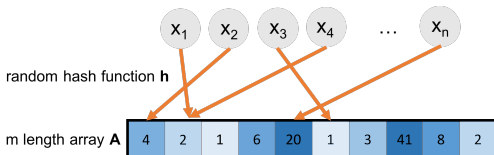> hash function. $m$: size of Count-min sketch array.

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

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

$f(x)$: 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.

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

---

$f(x)$: 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.

**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 $m = \frac{2k}{\epsilon}$. How can we improve the success probability?

$f(x)$: 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.
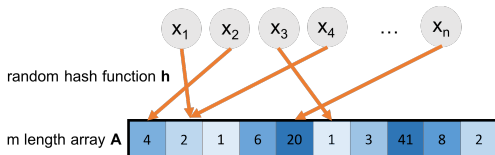
# Count-Min Sketch Accuracy



random hash function **h**

m length array **A**

**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.
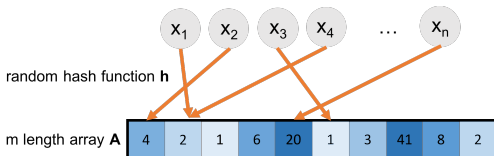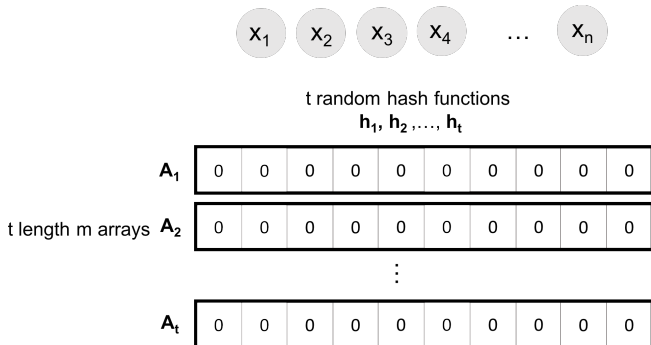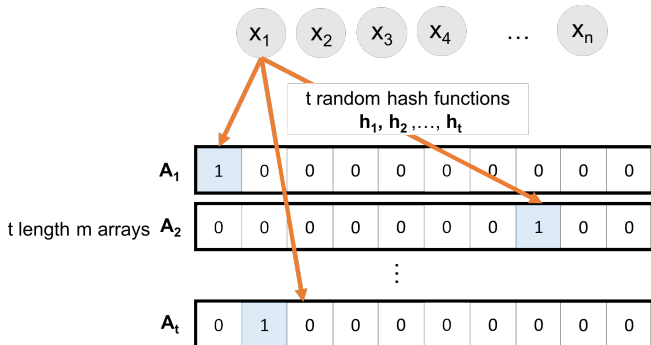
# Count-Min Sketch Repetition

x$_1$  x$_2$  x$_3$  x$_4$  …  x$_n$

t random hash functions
**h$_1$, h$_2$,…, h$_t$**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**A$_1$**

t length m arrays **A$_2$**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |

**A$_t$**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

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 taking the average?

# Count-Min Sketch Repetition



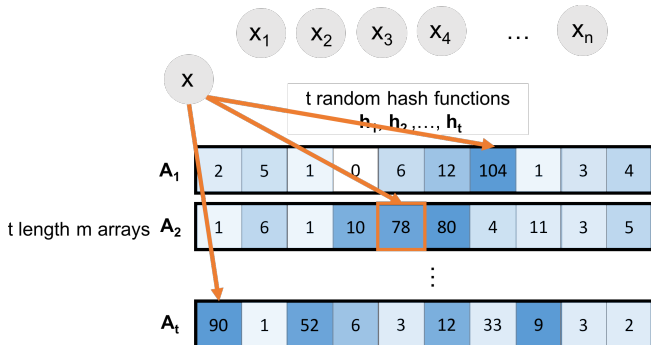Estimate $f(x)$ with $\tilde{f}(x) = \min_{i \in [t]} A_i[h_i(x)]$. (Count-min sketch)

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

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

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 $\geq 1/2$:

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

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

- For every $x$ and $i \in [t]$, we know that for $m = \frac{2k}{\epsilon}$, with probability $\geq 1/2$:
$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

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

# Count-Min Sketch Analysis



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

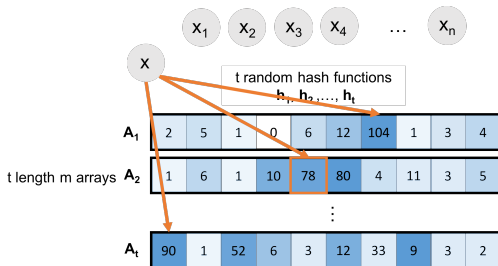- For every $x$ and $i \in [t]$, we know that for $m = \frac{2k}{\epsilon}$, with probability $\geq 1/2$:
$$f(x) \leq A_i[h_i(x)] \leq f(x) + \frac{\epsilon n}{k}.$$

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

- To get a good estimate with probability $\geq 1 - \delta$, set $t = \log_2(1/\delta)$.

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.

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

**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?

## 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.,

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

## 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.

## 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

## Hashing for Distinct Elements

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

## Hashing for Distinct Elements

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 $\tilde{d} = \frac{1}{s} - 1$

## Hashing for Distinct Elements

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 $\tilde{d} = \frac{1}{s} - 1$

## Hashing for Distinct Elements

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 $\tilde{d} = \frac{1}{s} - 1$

## Hashing for Distinct Elements

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

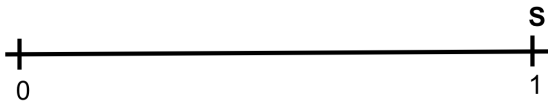Min-Hashing for Distinct Elements (variant of Flajolet-Martin):

- 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 $\tilde{d} = \frac{1}{s} - 1$

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 $\tilde{d} = \frac{1}{s} - 1$

# Hashing for Distinct Elements

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 $i = 1, \ldots, n$
  - $s := \min(s, h(x_i))$
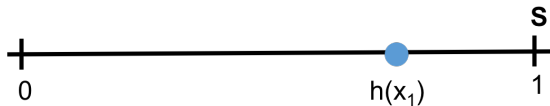- Return $\tilde{d} = \frac{1}{s} - 1$

## Hashing for Distinct Elements

**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 $i = 1, \ldots, n$
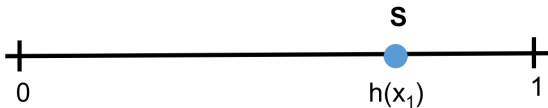    - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

# Hashing for Distinct Elements

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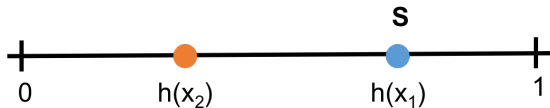- Return $\tilde{d} = \frac{1}{s} - 1$

# Hashing for Distinct Elements

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 $i = 1, \ldots, n$
    - $s := \min(s, h(x_i))$
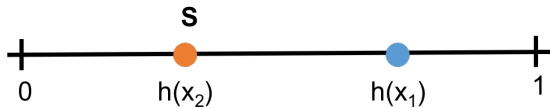- Return $\tilde{d} = \frac{1}{s} - 1$

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 $\tilde{d} = \frac{1}{s} - 1$

Min-Hashing for Distinct Elements:

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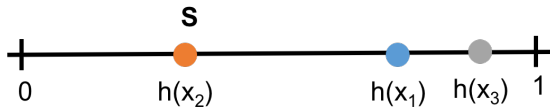- Return $\tilde{d} = \frac{1}{s} - 1$

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.

Min-Hashing for Distinct Elements:

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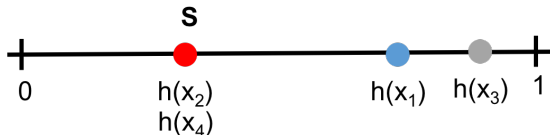- Return $\tilde{d} = \frac{1}{s} - 1$



- 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.

Min-Hashing for Distinct Elements:

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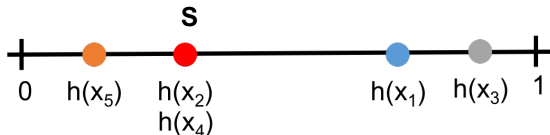    - $s := \min(s, h(x_i))$
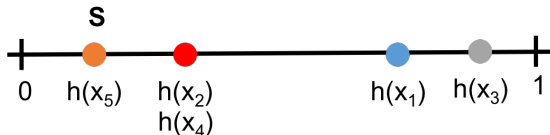- Return $\tilde{d} = \frac{1}{s} - 1$



- 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.

15

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

## Performance in Expectation

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



$$\mathbb{E}[s] =$$

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



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

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



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

· So our estimate $\widehat{\mathsf{d}} = \frac{1}{\mathsf{s}} - 1$ is correct if **s** exactly equals its 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 $\mathbb{E}[\widehat{d}] = d$?

$s$ 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:

$s$ 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 $c \leq 4$:
$$(1 - c\epsilon)d \leq \widehat{d} \leq (1 + c\epsilon)d$$

So question is how well **s** concentrates around its mean.

$$\mathbb{E}[\mathsf{s}] = \frac{1}{d+1}$$

.

.

> **s**: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$, computed by hashing algorithm. $\widehat{\mathsf{d}} = \frac{1}{\mathsf{s}} - 1$: estimate of # distinct elements $d$.

So question is how well **s** concentrates around its mean.

$$\mathbb{E}[\mathsf{s}] = \frac{1}{d+1} \text{ and } \mathsf{Var}[\mathsf{s}] \leq \frac{1}{(d+1)^2} \text{ (}also \text{ }via \text{ }calculus\text{)}.$$

.

> **s**: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$, computed by hashing algorithm. $\widehat{\mathsf{d}} = \frac{1}{\mathsf{s}} - 1$: estimate of # distinct elements $d$.

So question is how well **s** concentrates around its mean.

$$\mathbb{E}[\mathsf{s}] = \frac{1}{d+1} \text{ and } \mathsf{Var}[\mathsf{s}] \leq \frac{1}{(d+1)^2} \text{ (}\textit{also via calculus}\text{)}.$$

Chebyshev's Inequality:

$$\Pr\left[|\mathsf{s} - \mathbb{E}[\mathsf{s}]| \geq \epsilon \mathbb{E}[\mathsf{s}]\right] \leq \frac{\mathsf{Var}[\mathsf{s}]}{(\epsilon \mathbb{E}[\mathsf{s}])^2} \qquad .$$

---

**s**: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$, computed by hashing algorithm. $\widehat{\mathsf{d}} = \frac{1}{\mathsf{s}} - 1$: estimate of # distinct elements $d$.

So question is how well $\mathbf{s}$ concentrates around its mean.

$$\mathbb{E}[\mathbf{s}] = \frac{1}{d+1} \text{ and } \text{Var}[\mathbf{s}] \leq \frac{1}{(d+1)^2} \text{ (also via calculus).}$$

Chebyshev's Inequality:

$$\Pr\left[|\mathbf{s} - \mathbb{E}[\mathbf{s}]| \geq \epsilon\mathbb{E}[\mathbf{s}]\right] \leq \frac{\text{Var}[\mathbf{s}]}{(\epsilon\mathbb{E}[\mathbf{s}])^2} = \frac{1}{\epsilon^2}.$$

$\mathbf{s}$: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$, computed by hashing algorithm. $\widehat{\mathbf{d}} = \frac{1}{\mathbf{s}} - 1$: estimate of # distinct elements $d$.

So question is how well **s** concentrates around its mean.

$$\mathbb{E}[\mathsf{s}] = \frac{1}{d+1} \text{ and } \mathsf{Var}[\mathsf{s}] \le \frac{1}{(d+1)^2} \text{ (also via calculus).}$$

**Chebyshev's Inequality:**

$$\Pr\left[|\mathsf{s} - \mathbb{E}[\mathsf{s}]| \ge \epsilon\mathbb{E}[\mathsf{s}]\right] \le \frac{\mathsf{Var}[\mathsf{s}]}{(\epsilon\mathbb{E}[\mathsf{s}])^2} = \frac{1}{\epsilon^2}.$$

Bound is vacuous for any $\epsilon < 1$.

> **s**: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$, computed by hashing algorithm. $\widehat{\mathsf{d}} = \frac{1}{\mathsf{s}} - 1$: estimate of # distinct elements $d$.

So question is how well **s** concentrates around its mean.

$$\mathbb{E}[s] = \frac{1}{d+1} \text{ and } \mathsf{Var}[s] \leq \frac{1}{(d+1)^2} \text{ (also via calculus).}$$

Chebyshev's Inequality:

$$\Pr\left[|s - \mathbb{E}[s]| \geq \epsilon\mathbb{E}[s]\right] \leq \frac{\mathsf{Var}[s]}{(\epsilon\mathbb{E}[s])^2} = \frac{1}{\epsilon^2}.$$

Bound is vacuous for any $\epsilon < 1$. How can we improve accuracy?

> **s**: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$, computed by hashing algorithm. $\widehat{d} = \frac{1}{s} - 1$: estimate of # distinct elements $d$.

## Improving Performance

Leverage the law of large numbers: improve accuracy via repeated independent trials.

## Improving Performance

Leverage the law of large numbers: improve accuracy via repeated independent trials.

### Hashing for Distinct Elements (Improved):

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

## Improving Performance

Leverage the law of large numbers: improve accuracy via repeated independent trials.

Hashing for Distinct Elements (Improved):

- Let $h_1, h_2, \ldots, h_k : U \to [0, 1]$ be random hash functions
- $s := 1$
- For $i = 1, \ldots, n$
    - $s := \min(s, h(x_i))$
- Return $\widehat{d} = \frac{1}{s} - 1$

## Improving Performance

Leverage the law of large numbers: improve accuracy via repeated independent trials.

### Hashing for Distinct Elements (Improved):

- Let $h_1, h_2, \ldots, h_k : U \rightarrow [0, 1]$ be random hash functions
- $s_1, s_2, \ldots, s_k := 1$
- For $i = 1, \ldots, n$
    - $s := \min(s, h(x_i))$
- Return $\widehat{d} = \frac{1}{s} - 1$

## Improving Performance

Leverage the law of large numbers: improve accuracy via repeated independent trials.

### Hashing for Distinct Elements (Improved):

- Let $h_1, h_2, \ldots, h_k : U \to [0, 1]$ be random hash functions
- $s_1, s_2, \ldots, s_k := 1$
- For $i = 1, \ldots, n$
  - For j=1,…,k, $s_j := \min(s_j, h_j(x_i))$
- Return $\widehat{d} = \frac{1}{s} - 1$

## Improving Performance

Leverage the law of large numbers: improve accuracy via repeated independent trials.

### Hashing for Distinct Elements (Improved):

- Let $h_1, h_2, \ldots, h_k : U \to [0, 1]$ be random hash functions
- $s_1, s_2, \ldots, s_k := 1$
- For $i = 1, \ldots, n$
    - For j=1,…,k, $s_j := \min(s_j, h_j(x_i))$
- $s := \frac{1}{k} \sum_{j=1}^{k} s_j$
- Return $\widehat{d} = \frac{1}{s} - 1$

Leverage the law of large numbers: improve accuracy via repeated independent trials.

## Hashing for Distinct Elements (Improved):

- Let $h_1, h_2, \ldots, h_k : U \to [0, 1]$ be random hash functions
- $s_1, s_2, \ldots, s_k := 1$
- For $i = 1, \ldots, n$
    - For j=1,…,k, $s_j := \min(s_j, h_j(x_i))$
- $s := \frac{1}{k} \sum_{j=1}^{k} s_j$
- Return $\widehat{d} = \frac{1}{s} - 1$

## Analysis

$\mathsf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathsf{s}_j$. Have already shown that for $j = 1, \ldots, k$:

$$\mathbb{E}[\mathsf{s}_j] = \frac{1}{d+1}$$

$$\mathsf{Var}[\mathsf{s}_j] \leq \frac{1}{(d+1)^2}$$

$\mathsf{s}_j$: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$. $\mathsf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathsf{s}_j$.
$\widehat{\mathsf{d}} = \frac{1}{\mathsf{s}} - 1$: estimate of # distinct elements $d$.

## Analysis

$\mathbf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathbf{s}_j$. Have already shown that for $j = 1, \ldots, k$:

$$\mathbb{E}[\mathbf{s}_j] = \frac{1}{d+1} \implies \mathbb{E}[\mathbf{s}]$$

$$\mathrm{Var}[\mathbf{s}_j] \leq \frac{1}{(d+1)^2}$$

$\mathbf{s}_j$: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$. $\mathbf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathbf{s}_j$.
$\widehat{d} = \frac{1}{\mathbf{s}} - 1$: estimate of # distinct elements $d$.

## Analysis

$\mathsf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathsf{s}_j$. Have already shown that for $j = 1, \dots, k$:

$$\mathbb{E}[\mathsf{s}_j] = \frac{1}{d+1} \implies \mathbb{E}[\mathsf{s}] = \frac{1}{d+1} \text{ (linearity of expectation)}$$

$$\mathsf{Var}[\mathsf{s}_j] \leq \frac{1}{(d+1)^2}$$

---

$\mathsf{s}_j$: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$. $\mathsf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathsf{s}_j$.
$\widehat{\mathsf{d}} = \frac{1}{\mathsf{s}} - 1$: estimate of # distinct elements $d$.

$s = \frac{1}{k} \sum_{j=1}^{k} s_j$. Have already shown that for $j = 1, \ldots, k$:

$$\mathbb{E}[s_j] = \frac{1}{d+1} \implies \mathbb{E}[s] = \frac{1}{d+1} \text{ (linearity of expectation)}$$

$$\text{Var}[s_j] \leq \frac{1}{(d+1)^2} \implies \text{Var}[s]$$

> $s_j$: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$. $s = \frac{1}{k} \sum_{j=1}^{k} s_j$.
> $\widehat{d} = \frac{1}{s} - 1$: estimate of $\#$ distinct elements $d$.

## Analysis

$\mathsf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathsf{s}_j$. Have already shown that for $j = 1, \ldots, k$:

$$\mathbb{E}[\mathsf{s}_j] = \frac{1}{d+1} \implies \mathbb{E}[\mathsf{s}] = \frac{1}{d+1} \text{ (linearity of expectation)}$$

$$\mathsf{Var}[\mathsf{s}_j] \leq \frac{1}{(d+1)^2} \implies \mathsf{Var}[\mathsf{s}] \leq \frac{1}{k \cdot (d+1)^2} \text{ (linearity of variance)}$$

$\mathsf{s}_j$: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$. $\mathsf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathsf{s}_j$.
$\widehat{\mathsf{d}} = \frac{1}{\mathsf{s}} - 1$: estimate of # distinct elements $d$.

$\mathsf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathsf{s}_j$. Have already shown that for $j = 1, \ldots, k$:

$$\mathbb{E}[\mathsf{s}_j] = \frac{1}{d+1} \implies \mathbb{E}[\mathsf{s}] = \frac{1}{d+1} \text{ (linearity of expectation)}$$

$$\mathsf{Var}[\mathsf{s}_j] \leq \frac{1}{(d+1)^2} \implies \mathsf{Var}[\mathsf{s}] \leq \frac{1}{k \cdot (d+1)^2} \text{ (linearity of variance)}$$

Chebyshev Inequality:

$$\Pr\left[|\mathsf{s} - \mathbb{E}[\mathsf{s}]| \geq \epsilon \mathbb{E}[\mathsf{s}]\right] \leq \frac{\mathsf{Var}[\mathsf{s}]}{(\epsilon \mathbb{E}[\mathsf{s}])^2}$$

$\mathsf{s}_j$: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$. $\mathsf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathsf{s}_j$.
$\widehat{\mathsf{d}} = \frac{1}{\mathsf{s}} - 1$: estimate of # distinct elements $d$.

$s = \frac{1}{k} \sum_{j=1}^{k} s_j$. Have already shown that for $j = 1, \ldots, k$:

$$\mathbb{E}[s_j] = \frac{1}{d+1} \implies \mathbb{E}[s] = \frac{1}{d+1} \text{ (linearity of expectation)}$$

$$\mathrm{Var}[s_j] \leq \frac{1}{(d+1)^2} \implies \mathrm{Var}[s] \leq \frac{1}{k \cdot (d+1)^2} \text{ (linearity of variance)}$$

Chebyshev Inequality:

$$\Pr\left[|s - \mathbb{E}[s]| \geq \epsilon \mathbb{E}[s]\right] \leq \frac{\mathrm{Var}[s]}{(\epsilon \mathbb{E}[s])^2} = \frac{\mathbb{E}[s]^2 / k}{\epsilon^2 \mathbb{E}[s]^2} = \frac{1}{k \cdot \epsilon^2}$$

---

$s_j$: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$. $s = \frac{1}{k} \sum_{j=1}^{k} s_j$.
$\widehat{d} = \frac{1}{s} - 1$: estimate of # distinct elements $d$.

$\mathbf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathbf{s}_j$. Have already shown that for $j = 1, \ldots, k$:

$$\mathbb{E}[\mathbf{s}_j] = \frac{1}{d+1} \implies \mathbb{E}[\mathbf{s}] = \frac{1}{d+1} \text{ (linearity of expectation)}$$

$$\mathsf{Var}[\mathbf{s}_j] \leq \frac{1}{(d+1)^2} \implies \mathsf{Var}[\mathbf{s}] \leq \frac{1}{k \cdot (d+1)^2} \text{ (linearity of variance)}$$

Chebyshev Inequality:

$$\Pr\left[|\mathbf{s} - \mathbb{E}[\mathbf{s}]| \geq \epsilon \mathbb{E}[\mathbf{s}]\right] \leq \frac{\mathsf{Var}[\mathbf{s}]}{(\epsilon \mathbb{E}[\mathbf{s}])^2} = \frac{\mathbb{E}[\mathbf{s}]^2 / k}{\epsilon^2 \mathbb{E}[\mathbf{s}]^2} = \frac{1}{k \cdot \epsilon^2}$$

How should we set $k$ if we want an error with probability at most $\delta$?

---

$\mathbf{s}_j$: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$. $\mathbf{s} = \frac{1}{k} \sum_{j=1}^{k} \mathbf{s}_j$.
$\widehat{\mathbf{d}} = \frac{1}{\mathbf{s}} - 1$: estimate of # distinct elements $d$.

$\mathsf{s} = \frac{1}{k}\sum_{j=1}^{k}\mathsf{s}_j$. Have already shown that for $j = 1, \ldots, k$:

$$\mathbb{E}[\mathsf{s}_j] = \frac{1}{d+1} \implies \mathbb{E}[\mathsf{s}] = \frac{1}{d+1} \text{ (linearity of expectation)}$$

$$\mathsf{Var}[\mathsf{s}_j] \leq \frac{1}{(d+1)^2} \implies \mathsf{Var}[\mathsf{s}] \leq \frac{1}{k \cdot (d+1)^2} \text{ (linearity of variance)}$$

Chebyshev Inequality:

$$\Pr\left[|\mathsf{s} - \mathbb{E}[\mathsf{s}]| \geq \epsilon\mathbb{E}[\mathsf{s}]\right] \leq \frac{\mathsf{Var}[\mathsf{s}]}{(\epsilon\mathbb{E}[\mathsf{s}])^2} = \frac{\mathbb{E}[\mathsf{s}]^2/k}{\epsilon^2\mathbb{E}[\mathsf{s}]^2} = \frac{1}{k \cdot \epsilon^2}$$

How should we set $k$ if we want an error with probability at most $\delta$?
$k = \frac{1}{\epsilon^2 \cdot \delta}$.

---

$\mathsf{s}_j$: minimum of $d$ distinct hashes chosen randomly over $[0, 1]$. $\mathsf{s} = \frac{1}{k}\sum_{j=1}^{k}\mathsf{s}_j$.
$\widehat{\mathsf{d}} = \frac{1}{\mathsf{s}} - 1$: estimate of # distinct elements $d$.

$\mathbf{s} = \frac{1}{k}\sum_{j=1}^{k}\mathbf{s}_j$. Have already shown that for $j = 1,\ldots,k$:

$$\mathbb{E}[\mathbf{s}_j] = \frac{1}{d+1} \implies \mathbb{E}[\mathbf{s}] = \frac{1}{d+1} \text{ (linearity of expectation)}$$

$$\mathrm{Var}[\mathbf{s}_j] \leq \frac{1}{(d+1)^2} \implies \mathrm{Var}[\mathbf{s}] \leq \frac{1}{k\cdot(d+1)^2} \text{ (linearity of variance)}$$

Chebyshev Inequality:

$$\Pr\left[|\mathbf{s} - \mathbb{E}[\mathbf{s}]| \geq \epsilon\mathbb{E}[\mathbf{s}]\right] \leq \frac{\mathrm{Var}[\mathbf{s}]}{(\epsilon\mathbb{E}[\mathbf{s}])^2} = \frac{\mathbb{E}[\mathbf{s}]^2/k}{\epsilon^2\mathbb{E}[\mathbf{s}]^2} = \frac{1}{k\cdot\epsilon^2} = \frac{\epsilon^2\cdot\delta}{\epsilon^2} = \delta.$$

How should we set $k$ if we want an error with probability at most $\delta$?
$k = \frac{1}{\epsilon^2\cdot\delta}$.

---

$\mathbf{s}_j$: minimum of $d$ distinct hashes chosen randomly over $[0,1]$. $\mathbf{s} = \frac{1}{k}\sum_{j=1}^{k}\mathbf{s}_j$.
$\widehat{\mathbf{d}} = \frac{1}{\mathbf{s}} - 1$: estimate of # distinct elements $d$.

# Space Complexity

Hashing for Distinct Elements:

- Let $h_1, h_2, \ldots, h_k : U \to [0, 1]$ be random hash functions
- $s_1, s_2, \ldots, s_k := 1$
- For $i = 1, \ldots, n$
    - For j=1,…, k, $s_j := \min(s_j, h_j(x_i))$
- $s := \frac{1}{k} \sum_{j=1}^{k} s_j$
- Return $\widehat{d} = \frac{1}{s} - 1$



- Setting $k = \frac{1}{\epsilon^2 \cdot \delta}$, algorithm returns $\widehat{d}$ with $|d - \widehat{d}| \le 4\epsilon \cdot d$ with probability at least $1 - \delta$.

Hashing for Distinct Elements:

- Let $h_1, h_2, \ldots, h_k : U \to [0, 1]$ be random hash functions
- $s_1, s_2, \ldots, s_k := 1$
- For $i = 1, \ldots, n$
    - For j=1,…, k, $s_j := \min(s_j, h_j(x_i))$
- $s := \frac{1}{k} \sum_{j=1}^{k} s_j$
- Return $\widehat{d} = \frac{1}{s} - 1$



- Setting $k = \frac{1}{\epsilon^2 \cdot \delta}$, algorithm returns $\widehat{d}$ with $|d - \widehat{d}| \le 4\epsilon \cdot d$ with probability at least $1 - \delta$.
- Space complexity is $k = \frac{1}{\epsilon^2 \cdot \delta}$ real numbers $s_1, \ldots, s_k$.

Hashing for Distinct Elements:

- Let $h_1, h_2, \ldots, h_k : U \to [0, 1]$ be random hash functions
- $s_1, s_2, \ldots, s_k := 1$
- For $i = 1, \ldots, n$
    - For j=1,..., k, $s_j := \min(s_j, h_j(x_i))$
- $s := \frac{1}{k} \sum_{j=1}^{k} s_j$
- Return $\widehat{d} = \frac{1}{s} - 1$



- Setting $k = \frac{1}{\epsilon^2 \cdot \delta}$, algorithm returns $\widehat{d}$ with $|d - \widehat{d}| \leq 4\epsilon \cdot d$ with probability at least $1 - \delta$.
- Space complexity is $k = \frac{1}{\epsilon^2 \cdot \delta}$ real numbers $s_1, \ldots, s_k$.
- $\delta = 5\%$ failure rate gives a factor 20 overhead in space complexity.

How can we improve our dependence on the failure rate $\delta$?

## Improved Failure Rate

How can we improve our dependence on the failure rate $\delta$?

**The median trick:** Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/5$ – each using $k = \frac{1}{\delta'\epsilon^2} = \frac{5}{\epsilon^2}$ hash functions.

## Improved Failure Rate

How can we improve our dependence on the failure rate $\delta$?

**The median trick:** Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/5$ – each using $k = \frac{1}{\delta'\epsilon^2} = \frac{5}{\epsilon^2}$ hash functions.

- Letting $\widehat{d}_1, \ldots, \widehat{d}_t$ be the outcomes of the $t$ trials, return $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

How can we improve our dependence on the failure rate $\delta$?

**The median trick:** Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/5$ – each using $k = \frac{1}{\delta'\epsilon^2} = \frac{5}{\epsilon^2}$ hash functions.

- Letting $\widehat{d}_1, \ldots, \widehat{d}_t$ be the outcomes of the $t$ trials, return $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$

# Improved Failure Rate

How can we improve our dependence on the failure rate $\delta$?

**The median trick:** Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/5$ – each using $k = \frac{1}{\delta'\epsilon^2} = \frac{5}{\epsilon^2}$ hash functions.
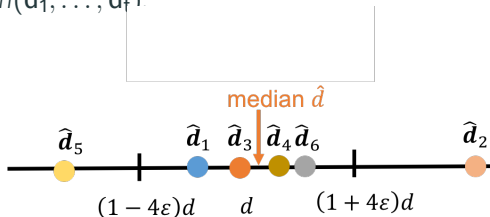
- Letting $\widehat{d}_1, \ldots, \widehat{d}_t$ be the outcomes of the $t$ trials, return $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$



- If $> 1/2$ of trials fall in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$, then the median will.

How can we improve our dependence on the failure rate $\delta$?

**The median trick:** Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/5$ – each using $k = \frac{1}{\delta' \epsilon^2} = \frac{5}{\epsilon^2}$ hash functions.

- Letting $\widehat{d}_1, \ldots, \widehat{d}_t$ be the outcomes of the $t$ trials, return $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$



- If $> 1/2$ of trials fall in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$, then the median will.
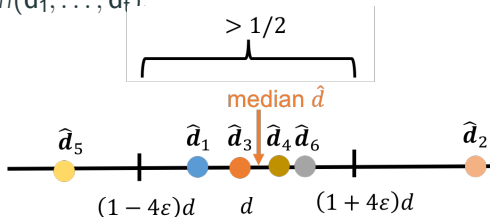
- Have $< 1/2$ of trials on both the left and right.

How can we improve our dependence on the failure rate $\delta$?

**The median trick:** Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/5$ – each using $k = \frac{1}{\delta' \epsilon^2} = \frac{5}{\epsilon^2}$ hash functions.

- Letting $\widehat{d}_1, \ldots, \widehat{d}_t$ be the outcomes of the $t$ trials, return $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$
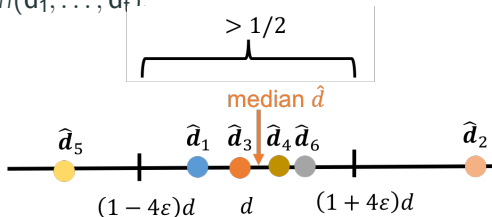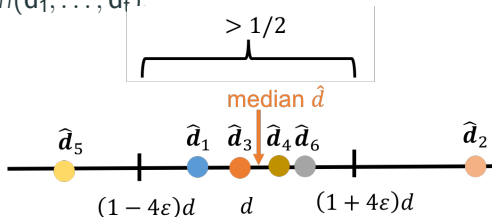


- If $> 2/3$ of trials fall in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$, then the median will.
- Have $< 1/3$ of trials on both the left and right.

## The Median Trick

- $\widehat{d}_1, \ldots, \widehat{d}_t$ are the outcomes of the $t$ trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $4/5$.
- $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

What is the probability that the median $\widehat{d}$ falls in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$?

- $\widehat{d}_1, \ldots, \widehat{d}_t$ are the outcomes of the $t$ trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $4/5$.
- $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

What is the probability that the median $\widehat{d}$ falls in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$?

- Let $X$ be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$.

- $\widehat{d}_1, \ldots, \widehat{d}_t$ are the outcomes of the $t$ trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $4/5$.
- $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

What is the probability that the median $\widehat{d}$ falls in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$?

- Let $X$ be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$.

$$\Pr\left(\widehat{d} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]\right) \leq \Pr\left(X < \frac{2}{3} \cdot t\right)$$

- $\widehat{d}_1, \ldots, \widehat{d}_t$ are the outcomes of the $t$ trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $4/5$.
- $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

What is the probability that the median $\widehat{d}$ falls in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$?

- Let $X$ be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$.
  $\mathbb{E}[X] =$

$$\Pr\left(\widehat{d} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]\right) \leq \Pr\left(X < \frac{2}{3} \cdot t\right)$$

- $\widehat{d}_1, \ldots, \widehat{d}_t$ are the outcomes of the $t$ trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $4/5$.
- $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

What is the probability that the median $\widehat{d}$ falls in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$?

- Let $X$ be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$. $\mathbb{E}[X] = \frac{4}{5} \cdot t$.

$$\Pr\left(\widehat{d} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]\right) \leq \Pr\left(X < \frac{2}{3} \cdot t\right)$$

- $\widehat{d}_1, \ldots, \widehat{d}_t$ are the outcomes of the $t$ trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $4/5$.
- $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

What is the probability that the median $\widehat{d}$ falls in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$?

- Let $X$ be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$. $\mathbb{E}[X] = \frac{4}{5} \cdot t$.

$$\Pr\left(\widehat{d} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]\right) \leq \Pr\left(X < \frac{5}{6} \cdot \mathbb{E}[X]\right)$$

- $\widehat{d}_1, \ldots, \widehat{d}_t$ are the outcomes of the $t$ trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $4/5$.
- $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

What is the probability that the median $\widehat{d}$ falls in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$?

- Let $X$ be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$. $\mathbb{E}[X] = \frac{4}{5} \cdot t$.

$$\Pr\left(\widehat{d} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]\right) \leq \Pr\left(X < \frac{5}{6} \cdot \mathbb{E}[X]\right) \leq \Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{6}\mathbb{E}[X]\right)$$

- $\widehat{d}_1, \ldots, \widehat{d}_t$ are the outcomes of the $t$ trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $4/5$.
- $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

What is the probability that the median $\widehat{d}$ falls in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$?

- Let $X$ be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$. $\mathbb{E}[X] = \frac{4}{5} \cdot t$.

$$\Pr\left(\widehat{d} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]\right) \leq \Pr\left(X < \frac{5}{6} \cdot \mathbb{E}[X]\right) \leq \Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{6}\mathbb{E}[X]\right)$$

Apply Chernoff bound:

# The Median Trick

- $\widehat{d}_1, \ldots, \widehat{d}_t$ are the outcomes of the $t$ trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $4/5$.
- $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

What is the probability that the median $\widehat{d}$ falls in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$?

- Let $X$ be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$. $\mathbb{E}[X] = \frac{4}{5} \cdot t$.

$$\Pr\left(\widehat{d} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]\right) \leq \Pr\left(X < \frac{5}{6} \cdot \mathbb{E}[X]\right) \leq \Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{6}\mathbb{E}[X]\right)$$

Apply Chernoff bound:

$$\Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{6}\mathbb{E}[X]\right) \leq 2\exp\left(-\frac{\frac{1}{6}^2 \cdot \frac{4}{5}t}{2 + 1/6}\right) = O\left(e^{-ct}\right).$$

- $\widehat{d}_1, \ldots, \widehat{d}_t$ are the outcomes of the $t$ trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $4/5$.
- $\widehat{d} = median(\widehat{d}_1, \ldots, \widehat{d}_t)$.

What is the probability that the median $\widehat{d}$ falls in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$?

- Let $X$ be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$. $\mathbb{E}[X] = \frac{4}{5} \cdot t$.

$$\Pr\left(\widehat{d} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]\right) \leq \Pr\left(X < \frac{5}{6} \cdot \mathbb{E}[X]\right) \leq \Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{6}\mathbb{E}[X]\right)$$

Apply Chernoff bound:

$$\Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{6}\mathbb{E}[X]\right) \leq 2\exp\left(-\frac{\frac{1}{6}^2 \cdot \frac{4}{5}t}{2 + 1/6}\right) = O\left(e^{-ct}\right).$$

- Setting $t = O(\log(1/\delta))$ gives failure probability $e^{-\log(1/\delta)} = \delta$.

## Median Trick

**Upshot:** The median of $t = O(\log(1/\delta))$ independent runs of the hashing algorithm for distinct elements returns $\widehat{d} \in [(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $1 - \delta$.

## Median Trick

**Upshot:** The median of $t = O(\log(1/\delta))$ independent runs of the hashing algorithm for distinct elements returns $\widehat{d} \in [(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $1 - \delta$.

**Total Space Complexity:** $t$ trials, each using $k = \frac{1}{\epsilon^2 \delta'}$ hash functions, for $\delta' = 1/5$. Space is $\frac{5t}{\epsilon^2} = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ real numbers (the minimum value of each hash function).

## Median Trick

**Upshot:** The median of $t = O(\log(1/\delta))$ independent runs of the hashing algorithm for distinct elements returns $\widehat{d} \in [(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $1 - \delta$.

**Total Space Complexity:** $t$ trials, each using $k = \frac{1}{\epsilon^2 \delta'}$ hash functions, for $\delta' = 1/5$. Space is $\frac{5t}{\epsilon^2} = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ real numbers (the minimum value of each hash function).

No dependence on the number of distinct elements $d$ or the number of items in the stream $n$! Both of these numbers are typically very large.

## Median Trick

**Upshot:** The median of $t = O(\log(1/\delta))$ independent runs of the hashing algorithm for distinct elements returns $\widehat{\mathsf{d}} \in [(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $1 - \delta$.

**Total Space Complexity:** $t$ trials, each using $k = \frac{1}{\epsilon^2 \delta'}$ hash functions, for $\delta' = 1/5$. Space is $\frac{5t}{\epsilon^2} = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ real numbers (the minimum value of each hash function).

No dependence on the number of distinct elements $d$ or the number of items in the stream $n$! Both of these numbers are typically very large.

**A note on the median:** The median is often used as a robust alternative to the mean, when there are outliers (e.g., heavy tailed distributions, corrupted data).