

COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Spring 2026.

Lecture 8

Logistics

- Problem Set 2 is due Sunday 3/8 at 11:59pm.
- ~~Problem Set 1 grades will be released later today.~~

Logistics

- Problem Set 2 is due Sunday 3/8 at 11:59pm.
- Problem Set 1 grades will be released later today.
- Midterm is **in class** next Thursday, 3/12.
- Study guide is posted on the course webpage (under the midterm row in the schedule tab).
- Past midterms are posted in Canvas. Note that some cover a bit more material than we have seen – e.g., we will not cover the **Johnson-Lindenstrauss lemma** or **low-distortion embeddings** before Midterm 1.
- At least some of next Tuesday's lecture will be for midterm review. I will also hold additional review office hours next **Wednesday, 11am-12pm**.
- If you need extended time on the midterm, please reach out by **this Friday** at the latest to make a plan.

Summary

Last Class:

$X_1 \quad X_2 \quad X_3 \dots$

- Overview of streaming algorithms.
- The (ϵ, k) -frequent items problem and its applications.

- Count-Min sketch algorithm for frequent items.

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

n items

identify all items that appear $\geq \frac{n}{k}$ times

not identify " $\leq (1-\epsilon)\frac{n}{k}$ times

Summary

Last Class:

- Overview of streaming algorithms.
- The (ϵ, k) -frequent items problem and its applications.
- Count-Min sketch algorithm for frequent items.

This Class:

- Distinct items counting via min-hashing.
- Success probability boosting via the **median trick**.

Distinct Elements

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

E.g.,

$(1, 5, 7, 5, 2)1 \rightarrow 4$ distinct elements

Distinct Elements

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

E.g.,

1, 5, 7, 5, 2, 1 \rightarrow 4 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, \dots, x_n , estimate the number of distinct elements in the stream.

E.g.,

1, 5, 7, 5, 2, 1 \rightarrow 4 distinct elements

*O(n) space
↓
exactly solve
by storing a set*

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, \dots, x_n , estimate the number of distinct elements.

Hashing for Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream x_1, \dots, 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)

$s := 1$
minimum hash value

- For $i = 1, \dots, n$

$s := \min(s, h(x_i))$

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

$$h(x_1) = 0.07381\dots$$

$$h(x_2) = 0.372\dots$$

$$h(x_3) = 0.73\dots$$

\vdots

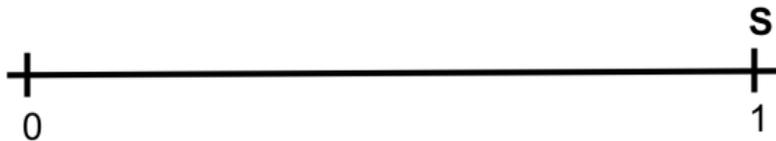
$$h(x_n) = 0.521\dots$$

Hashing for Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream x_1, \dots, 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 $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

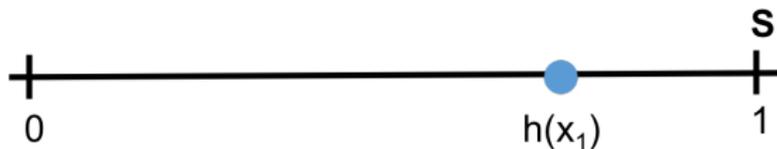


Hashing for Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream x_1, \dots, 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 $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

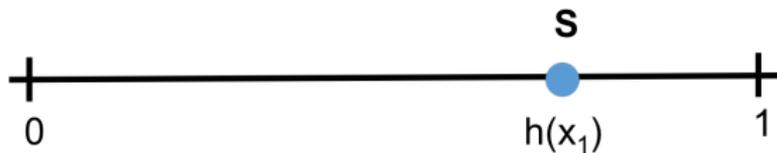


Hashing for Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream x_1, \dots, 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 $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

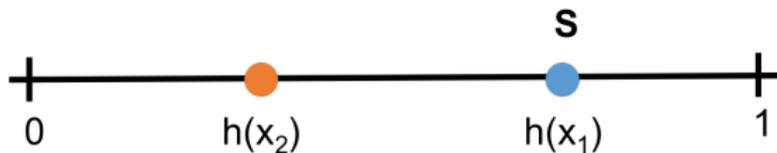


Hashing for Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream x_1, \dots, 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 $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

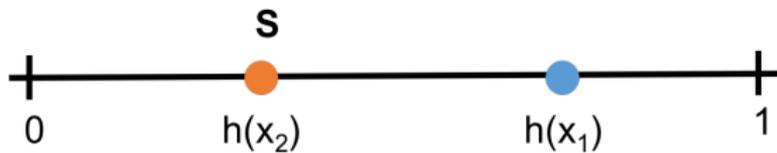


Hashing for Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream x_1, \dots, 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 $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

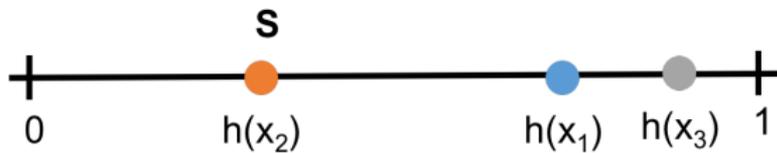


Hashing for Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream x_1, \dots, 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 $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

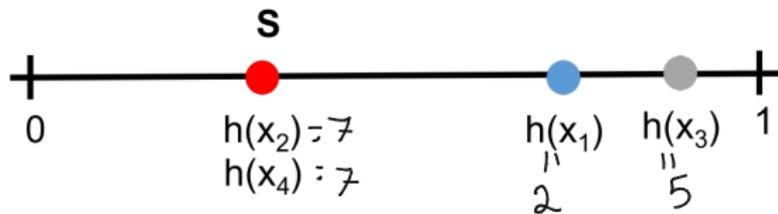


Hashing for Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream x_1, \dots, 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 $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

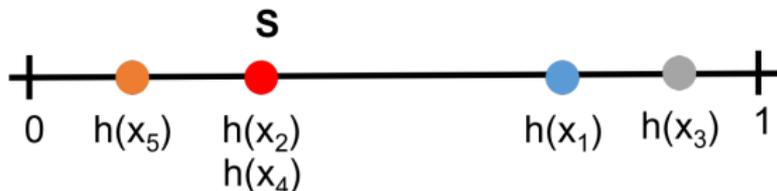


Hashing for Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream x_1, \dots, 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 $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

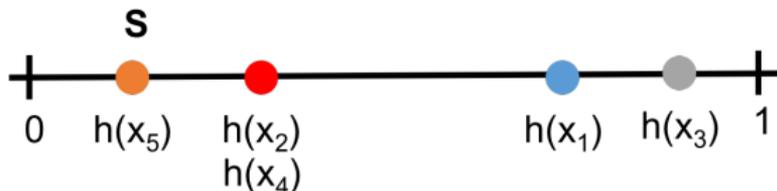


Hashing for Distinct Elements

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

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

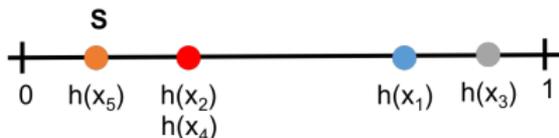
- Let $h : U \rightarrow [0, 1]$ be a random hash function (with a real valued output)
- $s := 1$
- For $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$



Hashing for Distinct Elements

Min-Hashing for Distinct Elements:

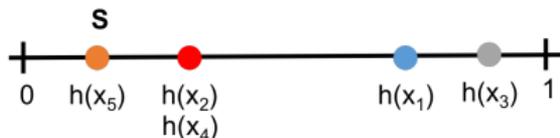
- Let $h : U \rightarrow [0, 1]$ be a random hash function (with a real valued output)
- $s := 1$
- For $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$



Hashing for Distinct Elements

Min-Hashing for Distinct Elements:

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



distinct items

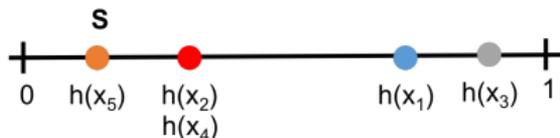
- After all items are processed, s is the minimum of d points chosen uniformly at random on $[0, 1]$. Where $d = \#$ distinct elements.

$d \uparrow$ $s \downarrow$
 $d \downarrow$ $s \uparrow$

Hashing for Distinct Elements

Min-Hashing for Distinct Elements:

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

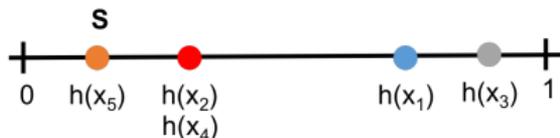


- After all items are processed, s is the minimum of d points chosen uniformly at random on $[0, 1]$. Where $d = \#$ distinct elements.
- Intuition: The larger d is, the smaller we expect s to be.

Hashing for Distinct Elements

Min-Hashing for Distinct Elements:

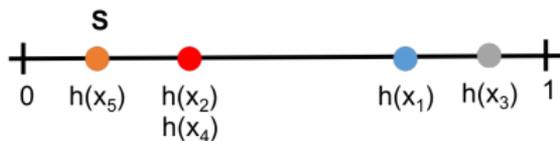
- Let $h : U \rightarrow [0, 1]$ be a random hash function (with a real valued output)
- $s := 1$
- For $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$



- 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.

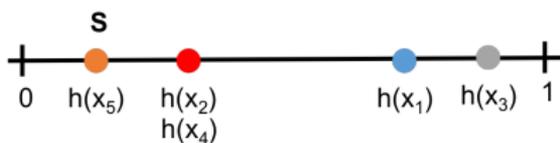
Performance in Expectation

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



Performance in Expectation

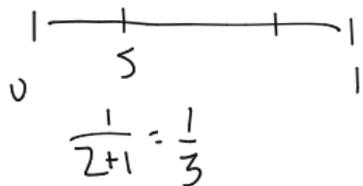
s is the minimum of d points chosen uniformly at random on $[0, 1]$.
Where $d = \#$ distinct elements.



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



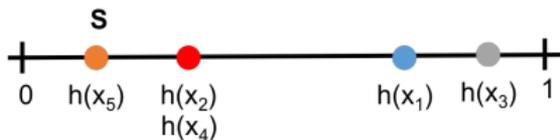
$$\frac{1}{1+1} = \frac{1}{2}$$



$$\frac{1}{2+1} = \frac{1}{3}$$

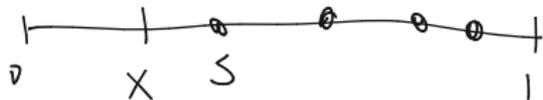
Performance in Expectation

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



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

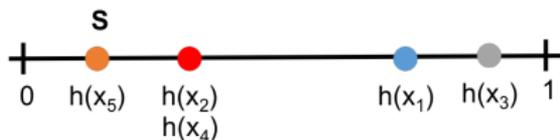
$$\Pr(s > x) = (1-x)^d$$



$$\int_0^\infty (1-x)^d = \frac{1}{d+1}$$

Performance in Expectation

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



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

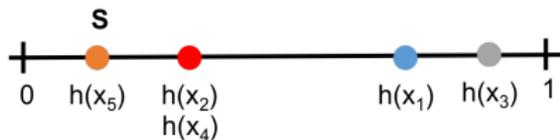
- So our estimate $\tilde{d} = \frac{1}{s} - 1$ is correct if s exactly equals its expectation.

$$s = \frac{1}{d+1}$$

$$\begin{aligned} (d+1)s &= 1 \\ d+1 &= 1/s \\ d &= 1/s - 1 \end{aligned}$$

Performance in Expectation

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



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

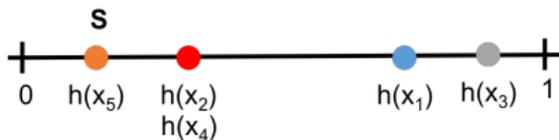
- So our estimate $\hat{d} = \frac{1}{s} - 1$ is correct if s exactly equals its expectation. Does this mean $\mathbb{E}[\hat{d}] = d$? No!

$$\mathbb{E}\left[\frac{1}{s} - 1\right] = d?$$

~~$\frac{1}{\mathbb{E}s} - 1 = d$~~

Performance in Expectation

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

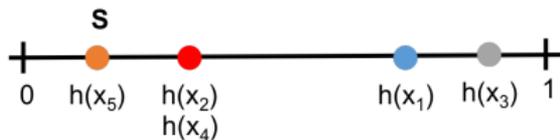


$$\mathbb{E}[s] = \frac{1}{d+1} \text{ (using } \mathbb{E}(s) = \int_0^\infty \Pr(s > x) dx \text{ + calculus)}$$

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

Performance in Expectation

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



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

- So our estimate $\hat{d} = \frac{1}{s} - 1$ is correct if s exactly equals its expectation. Does this mean $\mathbb{E}[\hat{d}] = d$? 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 $c \leq 4$.

$$(1 - c\epsilon)d \leq \hat{d} \leq (1 + c\epsilon)d$$
$$.96 \cdot d \leq \hat{d} \leq 1.04 \cdot d$$

$$s - \frac{1}{d+1} \leq \epsilon / (d+1)$$
$$\frac{1}{d+1} - s \leq \epsilon / d+1$$

$\{\epsilon = .01\}$

Initial Concentration Bound

So question is how well s concentrates around its mean.

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

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

Initial Concentration Bound

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

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

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

Initial Concentration Bound

So question is how well s concentrates around its mean.

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

Chebyshev's Inequality:

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

Handwritten annotations:
- A bracket above the denominator $(\epsilon \mathbb{E}[s])^2$ indicates the term $\frac{1}{(d+1)^2}$ from the variance bound.
- A bracket below the denominator $(\epsilon \mathbb{E}[s])^2$ indicates the term $\epsilon^2 \cdot \frac{1}{(d+1)^2}$.

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

Initial Concentration Bound

So question is how well s concentrates around its mean.

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

Chebyshev's Inequality:

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

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

Initial Concentration Bound

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:

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

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

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

Initial Concentration Bound

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:

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

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

\mathbf{s} : minimum of d distinct hashes chosen randomly over $[0, 1]$, computed by hashing algorithm. $\hat{\mathbf{d}} = \frac{1}{\mathbf{s}} - 1$: 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 $h : U \rightarrow [0, 1]$ be a random hash function
- $s := 1$
- For $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- Return $\hat{d} = \frac{1}{s} - 1$

Improving Performance

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

Hashing for Distinct Elements (Improved):

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

Improving Performance

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

Hashing for Distinct Elements (Improved):

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

Improving Performance

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

Hashing for Distinct Elements (Improved):

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

Improving Performance

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

Hashing for Distinct Elements (Improved):

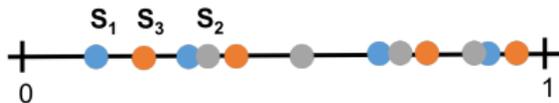
- Let $h_1, h_2, \dots, h_k : U \rightarrow [0, 1]$ be random hash functions
- $s_1, s_2, \dots, s_k := 1$
- For $i = 1, \dots, n$
 - For $j=1, \dots, k, s_j := \min(s_j, h_j(x_i))$
- $s := \frac{1}{k} \sum_{j=1}^k s_j$
- Return $\hat{d} = \frac{1}{s} - 1$

Improving Performance

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

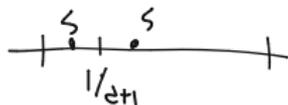
Hashing for Distinct Elements (Improved):

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



Analysis

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



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

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

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

$$\text{Var}[s] = \text{Var}\left[\frac{1}{k} \sum_{j=1}^k s_j\right]$$

$$= \frac{1}{k^2} \sum_{j=1}^k \text{Var}(s_j)$$

$$= \frac{1}{k^2} \cdot k \cdot \frac{1}{(d+1)^2} = \frac{1}{k} \cdot \frac{1}{(d+1)^2}$$

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

$\hat{d} = \frac{1}{s} - 1$: estimate of # distinct elements d .

Analysis

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

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

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

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

Analysis

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

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

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

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

Analysis

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

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

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

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

Analysis

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

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

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

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

Analysis

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

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

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

Chebyshev Inequality:

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

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

Analysis

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

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

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

Chebyshev Inequality:

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

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

Analysis

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

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

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

Chebyshev Inequality:

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

$k \cdot \epsilon^2 = \frac{1}{\delta}$
 $k = \frac{1}{\epsilon^2 \delta}$

How should we set k if we want an error with probability at most δ ?

$$k =$$

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

Analysis

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

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

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

Chebyshev Inequality:

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

How should we set k if we want an error with probability at most δ ?

$$k = \frac{1}{\epsilon^2 \cdot \delta}$$

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

$\hat{d} = \frac{1}{\mathbf{s}} - 1$: estimate of # distinct elements d .

Analysis

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

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

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

Chebyshev Inequality:

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

How should we set k if we want an error with probability at most δ ?

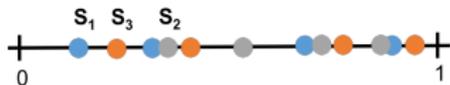
$$k = \frac{1}{\epsilon^2 \cdot \delta}.$$

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

Space Complexity

Hashing for Distinct Elements:

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

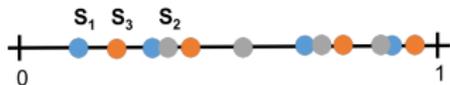


- Setting $k = \frac{1}{\epsilon^2 \cdot \delta}$, algorithm returns \hat{d} with $|d - \hat{d}| \leq 4\epsilon \cdot d$ with probability at least $1 - \delta$.

Space Complexity

Hashing for Distinct Elements:

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

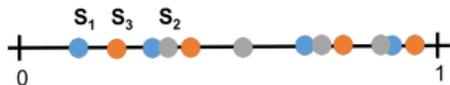


- Setting $k = \frac{1}{\epsilon^2 \cdot \delta}$, algorithm returns \hat{d} with $|d - \hat{d}| \leq 4\epsilon \cdot d$ with probability at least $1 - \delta$.
- Space complexity is $k = \frac{1}{\epsilon^2 \cdot \delta}$ real numbers s_1, \dots, s_k .

Space Complexity

Hashing for Distinct Elements:

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



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

- Space complexity is $k = \frac{1}{\epsilon^2 \cdot \delta}$ real numbers s_1, \dots, s_k .

$$\frac{1}{0.05} = 20$$

- $\delta = 5\%$ failure rate gives a factor 20 overhead in space complexity.

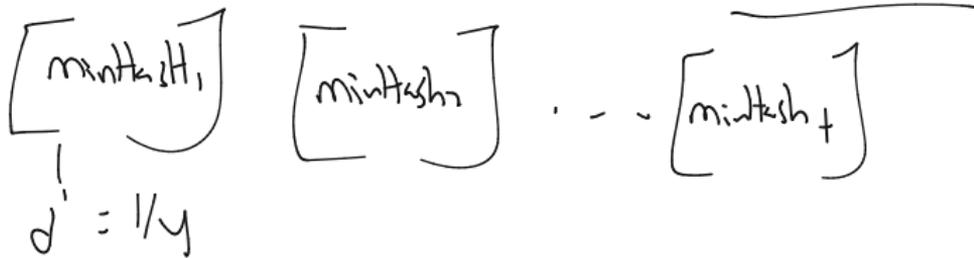
Improved Failure Rate

How can we improve our dependence on the failure rate δ ?

Improved Failure Rate

How can we improve our dependence on the failure rate δ ?

The median trick: Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/4$ - each using $k = \frac{1}{\delta'\epsilon^2} = \frac{4}{\epsilon^2}$ hash functions.



Improved Failure Rate

How can we improve our dependence on the failure rate δ ?

The median trick: Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/4$ – each using $k = \frac{1}{\delta'\epsilon^2} = \frac{4}{\epsilon^2}$ hash functions.

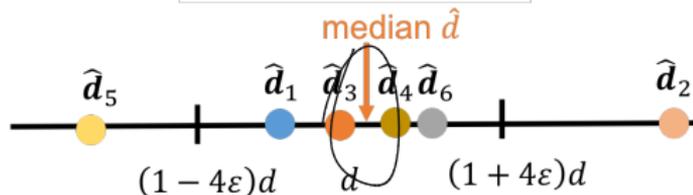
- Letting $\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t$ be the outcomes of the t trials, return $\hat{\mathbf{d}} = \text{median}(\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t)$.

Improved Failure Rate

How can we improve our dependence on the failure rate δ ?

The median trick: Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/4$ – each using $k = \frac{1}{\delta'\epsilon^2} = \frac{4}{\epsilon^2}$ hash functions.

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

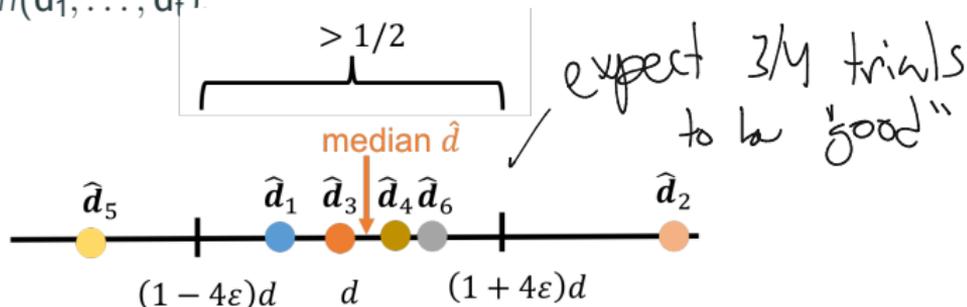


Improved Failure Rate

How can we improve our dependence on the failure rate δ ?

The median trick: Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/4$ – each using $k = \frac{1}{\delta'\epsilon^2} = \frac{4}{\epsilon^2}$ hash functions.

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



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

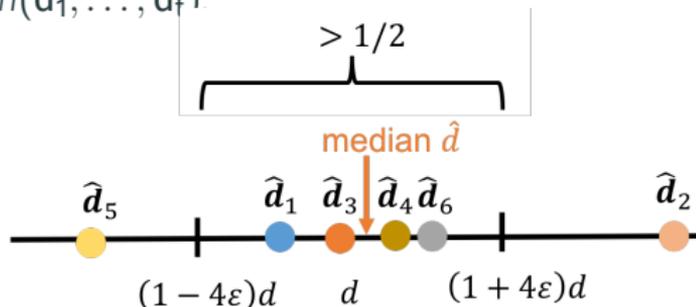
unrelated to δ
comes from conversion between δ and $\hat{\delta}$

Improved Failure Rate

How can we improve our dependence on the failure rate δ ?

The median trick: Run $t = O(\log 1/\delta)$ trials each with failure probability $\delta' = 1/4$ – each using $k = \frac{1}{\delta'\epsilon^2} = \frac{4}{\epsilon^2}$ hash functions.

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



- 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.

The Median Trick

- $\hat{d}_1, \dots, \hat{d}_t$ are the outcomes of the t trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $3/4$.
- $\hat{d} = \text{median}(\hat{d}_1, \dots, \hat{d}_t)$.

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

$\Pr(\hat{d} \text{ good estimate})$
↑
median

\hat{d} good estimate whenever $\geq \frac{1}{2}$ my trials are good estimates
 \hat{d} bad estimate $\Rightarrow < \frac{1}{2}$ trials are good
 $[X \leq t/2 \text{ where } X \# \text{ good trials}]$

The Median Trick

- $\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t$ are the outcomes of the t trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $3/4$.
- $\hat{\mathbf{d}} = \text{median}(\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t)$.

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

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

The Median Trick

- $\hat{d}_1, \dots, \hat{d}_t$ are the outcomes of the t trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $3/4$.
- $\hat{d} = \text{median}(\hat{d}_1, \dots, \hat{d}_t)$.

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

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

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

$\text{if } X \geq \frac{1}{2}t \Rightarrow \hat{d} \in [(1 - 4\epsilon)d, (1 + 4\epsilon)d]$

The Median Trick

- $\hat{d}_1, \dots, \hat{d}_t$ are the outcomes of the t trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $3/4$.
- $\hat{d} = \text{median}(\hat{d}_1, \dots, \hat{d}_t)$.

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

- Let X be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$.
 $\mathbb{E}[X] = 3/4 \cdot t$

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

GOOD

BAD

GOOD

GOOD

GOOD

The Median Trick

- $\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t$ are the outcomes of the t trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $3/4$.
- $\hat{\mathbf{d}} = \text{median}(\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t)$.

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

- Let X be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$.
 $\mathbb{E}[X] = \frac{3}{4} \cdot t$.

$$\begin{aligned} \Pr(\hat{\mathbf{d}} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]) &\leq \Pr\left(X < \frac{1}{2} \cdot t\right) \\ &\leq \frac{2}{3} \cdot \frac{3}{4} \cdot t \\ &\leq \frac{2}{3} \cdot \mathbb{E}[X] \end{aligned}$$

The Median Trick

- $\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t$ are the outcomes of the t trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $3/4$.
- $\hat{\mathbf{d}} = \text{median}(\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t)$.

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

- Let X be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$.
 $\mathbb{E}[X] = \frac{3}{4} \cdot t$.

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

The Median Trick

- $\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t$ are the outcomes of the t trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $3/4$.
- $\hat{\mathbf{d}} = \text{median}(\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t)$.

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

- Let X be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$.
 $\mathbb{E}[X] = \frac{3}{4} \cdot t$.

$$\Pr(\hat{\mathbf{d}} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]) \leq \Pr\left(X < \frac{2}{3} \cdot \mathbb{E}[X]\right) \leq \Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{3}\mathbb{E}[X]\right)$$

$$\Pr\left(X \leq \frac{1}{2}t\right)$$

The Median Trick

- $\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t$ are the outcomes of the t trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $3/4$.
- $\hat{\mathbf{d}} = \text{median}(\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t)$.

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

- Let X be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$.
 $\mathbb{E}[X] = \frac{3}{4} \cdot t$.

$$\Pr(\hat{\mathbf{d}} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]) \leq \Pr\left(X < \frac{2}{3} \cdot \mathbb{E}[X]\right) \leq \Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{3}\mathbb{E}[X]\right)$$

Apply Chernoff bound:

The Median Trick

- $\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t$ are the outcomes of the t trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $3/4$.
- $\hat{\mathbf{d}} = \text{median}(\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t)$.

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

- Let X be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$.
 $\mathbb{E}[X] = \frac{3}{4} \cdot t$.

$$\Pr(\hat{\mathbf{d}} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]) \leq \Pr\left(X < \frac{2}{3} \cdot \mathbb{E}[X]\right) \leq \Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{3}\mathbb{E}[X]\right)$$

Apply Chernoff bound:

$$\Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{3}\mathbb{E}[X]\right) \leq 2 \exp\left(-\frac{\frac{1}{3}^2 \cdot \frac{3}{4}t}{2 + \frac{1}{3}}\right) = O(e^{-ct}).$$

$C = \frac{1}{3^2} \cdot \frac{3}{4} \cdot \frac{3}{2+1/3}$

The Median Trick

- $\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t$ are the outcomes of the t trials, each falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $3/4$.
- $\hat{\mathbf{d}} = \text{median}(\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t)$.

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

- Let X be the # of trials falling in $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$.
 $\mathbb{E}[X] = \frac{3}{4} \cdot t$.

$$\Pr\left(\hat{\mathbf{d}} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]\right) \leq \Pr\left(X < \frac{2}{3} \cdot \mathbb{E}[X]\right) \leq \Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{3}\mathbb{E}[X]\right)$$

Apply Chernoff bound:

$$\Pr\left(|X - \mathbb{E}[X]| \geq \frac{1}{3}\mathbb{E}[X]\right) \leq 2 \exp\left(-\frac{\frac{1}{3}^2 \cdot \frac{3}{4}t}{2 + \frac{1}{3}}\right) = O(e^{-ct}) \leq \delta$$

- 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 $\hat{d} \in [(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ 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 $\hat{d} \in [(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $1 - \delta$.

Total Space Complexity: t trials, each using $k = \frac{1}{\epsilon^2 \delta'}$ hash functions, for $\delta' = 1/4$. Space is $\frac{4t}{\epsilon^2} = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ real numbers (the minimum value of each hash function). $\nearrow O\left(\frac{1}{\epsilon^2 \delta}\right)$

Median Trick

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

Total Space Complexity: t trials, each using $k = \frac{1}{\epsilon^2 \delta'}$ hash functions, for $\delta' = 1/4$. Space is $\frac{4t}{\epsilon^2} = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ 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 $\hat{d} \in [(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ with probability at least $1 - \delta$.

Total Space Complexity: t trials, each using $k = \frac{1}{\epsilon^2 \delta'}$ hash functions, for $\delta' = 1/4$. Space is $\frac{4t}{\epsilon^2} = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ real numbers (the minimum value of each hash function). $\frac{1}{k(d+1)^2}$ $k = \frac{1}{\delta \cdot \epsilon^2}$

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).

$$\tilde{d} = \frac{1}{5} - 1$$



Distinct Elements in Practice

Our algorithm uses continuous valued fully random hash functions.

Distinct Elements in Practice

Our algorithm uses continuous valued fully random hash functions.
Can't be implemented...

Distinct Elements in Practice

Our algorithm uses continuous valued fully random hash functions.
Can't be implemented...

- The idea of using the minimum hash value of x_1, \dots, x_n to estimate the number of distinct elements naturally extends to when the hash functions map to discrete values.

Distinct Elements in Practice

Our algorithm uses continuous valued fully random hash functions.
Can't be implemented...

- The idea of using the minimum hash value of x_1, \dots, x_n to estimate the number of distinct elements naturally extends to when the hash functions map to discrete values.
- Flajolet-Martin (LogLog) algorithm and [HyperLogLog](#).

Distinct Elements in Practice

Our algorithm uses continuous valued fully random hash functions.
Can't be implemented...

- The idea of using the minimum hash value of x_1, \dots, x_n to estimate the number of distinct elements naturally extends to when the hash functions map to discrete values.
- Flajolet-Martin (LogLog) algorithm and [HyperLogLog](#).

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
⋮	
$h(x_n)$	1011000

Distinct Elements in Practice

Our algorithm uses continuous valued fully random hash functions.
Can't be implemented...

- The idea of using the minimum hash value of x_1, \dots, x_n to estimate the number of distinct elements naturally extends to when the hash functions map to discrete values.
- Flajolet-Martin (LogLog) algorithm and **HyperLogLog**.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
⋮	
$h(x_n)$	1011000

Estimate # distinct elements
based on maximum number of
trailing zeros m .

$m=3$

Distinct Elements in Practice

Our algorithm uses continuous valued fully random hash functions.
Can't be implemented...

- The idea of using the minimum hash value of x_1, \dots, x_n to estimate the number of distinct elements naturally extends to when the hash functions map to discrete values.
- Flajolet-Martin (LogLog) algorithm and [HyperLogLog](#).

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
⋮	
$h(x_n)$	1011000

Estimate # distinct elements
based on maximum number of
trailing zeros m .

The more distinct hashes we see,
the higher we expect this
maximum to be.

LogLog Counting of Distinct Elements

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
⋮	
$h(x_n)$	1011000

Estimate # distinct elements
based on maximum number of
trailing zeros m .

LogLog Counting of Distinct Elements

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
⋮	
$h(x_n)$	1011000

Estimate # distinct elements based on maximum number of trailing zeros m .

With d distinct elements, roughly what do we expect m to be?

- a) $O(1)$ b) $O(\log d)$ c) $O(\sqrt{d})$ d) $O(d)$

LogLog Counting of Distinct Elements

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
⋮	
$h(x_n)$	1011000

Estimate # distinct elements based on maximum number of trailing zeros m .

With d distinct elements, roughly what do we expect m to be?

$$\Pr(h(x_i) \text{ has } x \text{ trailing zeros}) =$$

LogLog Counting of Distinct Elements

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
⋮	
$h(x_n)$	1011000

Estimate # distinct elements based on maximum number of trailing zeros m .

With d distinct elements, roughly what do we expect m to be?

$$\Pr(h(x_i) \text{ has } x \text{ trailing zeros}) = \frac{1}{2^x}$$

LogLog Counting of Distinct Elements

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
⋮	
$h(x_n)$	1011000

Estimate # distinct elements based on maximum number of trailing zeros m .

With d distinct elements, roughly what do we expect m to be?

$$\Pr(h(x_i) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}}$$

LogLog Counting of Distinct Elements

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
	⋮
$h(x_n)$	1011000

Estimate # distinct elements based on maximum number of trailing zeros m .

With d distinct elements, roughly what do we expect m to be?

$$\Pr(h(x_i) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}} = \frac{1}{d}.$$

LogLog Counting of Distinct Elements

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
⋮	
$h(x_n)$	1011000

Estimate # distinct elements based on maximum number of trailing zeros m .

With d distinct elements, roughly what do we expect m to be?

$$\Pr(h(x_i) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}} = \frac{1}{d}.$$

So with d distinct hashes, expect to see 1 with $\log d$ trailing zeros.
Expect $m \approx \log d$.

LogLog Counting of Distinct Elements

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
	⋮
$h(x_n)$	1011000

Estimate # distinct elements based on maximum number of trailing zeros m .

With d distinct elements, roughly what do we expect m to be?

$$\Pr(h(x_i) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}} = \frac{1}{d}.$$

So with d distinct hashes, expect to see 1 with $\log d$ trailing zeros. Expect $m \approx \log d$. m takes $\log \log d$ bits to store.

LogLog Counting of Distinct Elements

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
\vdots	
$h(x_n)$	1011000

Estimate # distinct elements based on maximum number of trailing zeros m .

With d distinct elements, roughly what do we expect m to be?

$$\Pr(h(x_j) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}} = \frac{1}{d}.$$

So with d distinct hashes, expect to see 1 with $\log d$ trailing zeros. Expect $m \approx \log d$. m takes $\log \log d$ bits to store.

Total Space: $O\left(\frac{\log \log d}{\epsilon^2}\right)$ for an ϵ approximate count.

LogLog Counting of Distinct Elements

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
\vdots	
$h(x_n)$	1011000

Estimate # distinct elements based on maximum number of trailing zeros m .

With d distinct elements, roughly what do we expect m to be?

$$\Pr(h(x_j) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}} = \frac{1}{d}.$$

So with d distinct hashes, expect to see 1 with $\log d$ trailing zeros. Expect $m \approx \log d$. m takes $\log \log d$ bits to store.

Total Space: $O\left(\frac{\log \log d}{\epsilon^2}\right)$ for an ϵ approximate count.

Note: Careful averaging of estimates from multiple hash functions.

LogLog Space Guarantees

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\text{space used} = O\left(\frac{\log \log d}{\epsilon^2}\right)$$

LogLog Space Guarantees

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\begin{aligned}\text{space used} &= O\left(\frac{\log \log d}{\epsilon^2}\right) \\ &= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} \text{ bits}^1\end{aligned}$$

1. 1.04 is the constant in the HyperLogLog analysis. Not important!

LogLog Space Guarantees

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\begin{aligned}\text{space used} &= O\left(\frac{\log \log d}{\epsilon^2}\right) \\ &= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} \text{ bits}^1 \\ &= \frac{1.04 \cdot 5}{.02^2} = 13000 \text{ bits} \approx 1.6 \text{ kB!}\end{aligned}$$

End here.

1. 1.04 is the constant in the HyperLogLog analysis. Not important!

LogLog Space Guarantees

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\begin{aligned}\text{space used} &= O\left(\frac{\log \log d}{\epsilon^2}\right) \\ &= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} \text{ bits}^1 \\ &= \frac{1.04 \cdot 5}{.02^2} = 13000 \text{ bits} \approx 1.6 \text{ kB!}\end{aligned}$$

Mergeable Sketch: Consider the case (essentially always in practice) that the items are processed on different machines.

1. 1.04 is the constant in the HyperLogLog analysis. Not important!

LogLog Space Guarantees

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\begin{aligned}\text{space used} &= O\left(\frac{\log \log d}{\epsilon^2}\right) \\ &= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} \text{ bits}^1 \\ &= \frac{1.04 \cdot 5}{.02^2} = 13000 \text{ bits} \approx 1.6 \text{ kB!}\end{aligned}$$

Mergeable Sketch: Consider the case (essentially always in practice) that the items are processed on different machines.

- Given data structures (sketches) $HLL(x_1, \dots, x_n)$, $HLL(y_1, \dots, y_n)$ is easy to merge them to give $HLL(x_1, \dots, x_n, y_1, \dots, y_n)$.

1. 1.04 is the constant in the HyperLogLog analysis. Not important!

LogLog Space Guarantees

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\begin{aligned}\text{space used} &= O\left(\frac{\log \log d}{\epsilon^2}\right) \\ &= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} \text{ bits}^1 \\ &= \frac{1.04 \cdot 5}{.02^2} = 13000 \text{ bits} \approx 1.6 \text{ kB!}\end{aligned}$$

Mergeable Sketch: Consider the case (essentially always in practice) that the items are processed on different machines.

- Given data structures (sketches) $HLL(x_1, \dots, x_n)$, $HLL(y_1, \dots, y_n)$ is easy to merge them to give $HLL(x_1, \dots, x_n, y_1, \dots, y_n)$. **How?**

1. 1.04 is the constant in the HyperLogLog analysis. Not important!

LogLog Space Guarantees

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\begin{aligned}\text{space used} &= O\left(\frac{\log \log d}{\epsilon^2}\right) \\ &= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} \text{ bits}^1 \\ &= \frac{1.04 \cdot 5}{.02^2} = 13000 \text{ bits} \approx 1.6 \text{ kB!}\end{aligned}$$

Mergeable Sketch: Consider the case (essentially always in practice) that the items are processed on different machines.

- Given data structures (sketches) $HLL(x_1, \dots, x_n)$, $HLL(y_1, \dots, y_n)$ is easy to merge them to give $HLL(x_1, \dots, x_n, y_1, \dots, y_n)$. **How?**
- Set the maximum # of trailing zeros to the maximum in the two sketches.

1. 1.04 is the constant in the HyperLogLog analysis. Not important!

HyperLogLog In Practice

Implementations: Google PowerDrill, Facebook Presto, Twitter Algebird, Amazon Redshift.

HyperLogLog In Practice

Implementations: Google PowerDrill, Facebook Presto, Twitter Algebird, Amazon Redshift.

HyperLogLog In Practice

Implementations: Google PowerDrill, Facebook Presto, Twitter Algebird, Amazon Redshift.

Use Case: Exploratory SQL-like queries on tables with 100s billions of rows. \sim 5 million count distinct queries per day.

HyperLogLog In Practice

Implementations: Google PowerDrill, Facebook Presto, Twitter Algebird, Amazon Redshift.

Use Case: Exploratory SQL-like queries on tables with 100s billions of rows. \sim 5 million count distinct queries per day. E.g.,

- **Count** number of **distinct** users in Germany that made at least one search containing the word 'auto' in the last month.
- **Count** number of **distinct** subject lines in emails sent by users that have registered in the last week, in comparison to number of emails sent overall (to estimate rates of spam accounts).

HyperLogLog In Practice

Implementations: Google PowerDrill, Facebook Presto, Twitter Algebird, Amazon Redshift.

Use Case: Exploratory SQL-like queries on tables with 100s billions of rows. ~ 5 million count distinct queries per day. E.g.,

- **Count** number of **distinct** users in Germany that made at least one search containing the word 'auto' in the last month.
- **Count** number of **distinct** subject lines in emails sent by users that have registered in the last week, in comparison to number of emails sent overall (to estimate rates of spam accounts).

Traditional *COUNT*, *DISTINCT* SQL calls are far too slow, especially when the data is distributed across many servers.

Questions?