

Frequency Estimation

CS 5968/6968: Data Str & Alg for Scalable Computing Spring
2023

Benwei Shi

University of Utah

2023-02-13

Table of contents

- 1 Problems
 - Frequency
 - Frequency Estimation in Stream
 - Heavy Hitters in Stream
- 2 Misra-Gries Sketch
 - Majority
 - Misra-Gries Sketch
- 3 Count Sketch
 - Count Sketch Algorithm
 - Count Sketch Analysis
- 4 Count-Min Sketch
- 5 Summary
 - Biased vs Unbiased

Table of Contents

- 1 Problems
 - Frequency
 - Frequency Estimation in Stream
 - Heavy Hitters in Stream
- 2 Misra-Gries Sketch
 - Majority
 - Misra-Gries Sketch
- 3 Count Sketch
 - Count Sketch Algorithm
 - Count Sketch Analysis
- 4 Count-Min Sketch
- 5 Summary
 - Biased vs Unbiased

Frequency

Given

- Metadata: a universe $[u] := \{1, 2, \dots, u\}$,
 - u is bigger than memory, $O(\lg u)$ is constant.

Frequency

Given

- Metadata: a universe $[u] := \{1, 2, \dots, u\}$,
 - u is bigger than memory, $O(\lg u)$ is constant.
- Data: a sequence $X := [x_i \in [u]]_{i=1}^n$,
 - i.e. $[x_1, x_2, \dots, x_n]$, $x_i \in [u]$ for all $i \in [n]$.

Frequency

Given

- Metadata: a universe $[u] := \{1, 2, \dots, u\}$,
 - u is bigger than memory, $O(\lg u)$ is constant.
- Data: a sequence $X := [x_i \in [u]]_{i=1}^n$,
 - i.e. $[x_1, x_2, \dots, x_n]$, $x_i \in [u]$ for all $i \in [n]$.

Goal

For any query $q \in [u]$, return the its count or frequency:

$$c(q) := \sum_{x \in X} \mathbb{1}(x, q), \quad f(q) := \frac{c(q)}{n}.$$

Frequency

Given

- Metadata: a universe $[u] := \{1, 2, \dots, u\}$,
 - u is bigger than memory, $O(\lg u)$ is constant.
- Data: a sequence $X := [x_i \in [u]]_{i=1}^n$,
 - i.e. $[x_1, x_2, \dots, x_n]$, $x_i \in [u]$ for all $i \in [n]$.

Goal

For any query $q \in [u]$, return the its count or frequency:

$$c(q) := \sum_{x \in X} \mathbb{1}(x, q), \quad f(q) := \frac{c(q)}{n}.$$

What if n is also too big for memory? Even bigger than external memory?

DDoS Attack Detection at Router

Detect high frequency IP addresses with limited memory.

Frequency Estimation in Stream

Given

- Metadata: a universe $[u] := \{1, 2, \dots, u\}$,
 - u is bigger than memory, $O(\lg u)$ is constant.
- Data: a sequence $X := [x_i \in [u]]_{i=1}^n$,
 - n is bigger than memory.

Goal

For any query $q \in [u]$, return the ε -approximation of its frequency, $\hat{f}_\varepsilon(q)$, s.t.

$$f(q) - \varepsilon \leq \hat{f}_\varepsilon(q) \leq f(q) + \varepsilon$$

Heavy Hitters in Stream

ϕ -Heavy Hitter

$y \in [u]$ is a ϕ -heavy hitter iff $f(y) > \phi$.

Heavy Hitters in Stream

ϕ -Heavy Hitter

$y \in [u]$ is a ϕ -heavy hitter iff $f(y) > \phi$.

ε -approximation of ϕ -heavy hitters, \hat{H}_ε^ϕ

- $y \in \hat{H}_\varepsilon^\phi$ if $f(y) > \phi$.
- $y \notin \hat{H}_\varepsilon^\phi$ if $f(y) < \phi - \varepsilon$.

Heavy Hitters in Stream

ϕ -Heavy Hitter

$y \in [u]$ is a ϕ -heavy hitter iff $f(y) > \phi$.

ε -approximation of ϕ -heavy hitters, \hat{H}_ε^ϕ

- $y \in \hat{H}_\varepsilon^\phi$ if $f(y) > \phi$.
- $y \notin \hat{H}_\varepsilon^\phi$ if $f(y) < \phi - \varepsilon$.
- Given f , you can find all ϕ -Heavy Hitters.

Heavy Hitters in Stream

ϕ -Heavy Hitter

$y \in [u]$ is a ϕ -heavy hitter iff $f(y) > \phi$.

ε -approximation of ϕ -heavy hitters, \hat{H}_ε^ϕ

- $y \in \hat{H}_\varepsilon^\phi$ if $f(y) > \phi$.
- $y \notin \hat{H}_\varepsilon^\phi$ if $f(y) < \phi - \varepsilon$.
- Given f , you can find all ϕ -Heavy Hitters.
- Given a \hat{f}_ε , for any $\phi \geq \varepsilon$, $\{y \in [u] \mid \hat{f}_\varepsilon(y) > \phi - \varepsilon\}$ is a \hat{H}_ε^ϕ .

Table of Contents

- 1 Problems
 - Frequency
 - Frequency Estimation in Stream
 - Heavy Hitters in Stream
- 2 Misra-Gries Sketch
 - Majority
 - Misra-Gries Sketch
- 3 Count Sketch
 - Count Sketch Algorithm
 - Count Sketch Analysis
- 4 Count-Min Sketch
- 5 Summary
 - Biased vs Unbiased

Majority

Goal

Find y if $f(y) > \frac{1}{2}$.

Algorithm: Majority(X)

```
1  $y \leftarrow \text{NaN}, c \leftarrow 0$   
2 forall  $x \in X$  do  
3   if  $y = x$  then  $c \leftarrow c + 1$   
4   else if  $c = 0$  then  $y \leftarrow x, c \leftarrow 1$   
5   else  $c \leftarrow c - 1$   
6 return  $y$ 
```

Boyer and Moore [1981].

Majority

Goal

Find y if $f(y) > \frac{1}{2}$.

Algorithm: Majority(X)

```

1  $y \leftarrow \text{NaN}, c \leftarrow 0$ 
2 forall  $x \in X$  do
3     if  $y = x$  then  $c \leftarrow c + 1$ 
4     else if  $c = 0$  then  $y \leftarrow x, c \leftarrow 1$ 
5     else  $c \leftarrow c - 1$ 
6 return  $y$ 
    
```

Boyer and Moore [1981].

- If there is no m s.t. $f(m) > \frac{1}{2}$, then whatever y is correct.

Majority

Goal

Find y if $f(y) > \frac{1}{2}$.

Algorithm: Majority(X)

```

1  $y \leftarrow \text{NaN}, c \leftarrow 0$ 
2 forall  $x \in X$  do
3   if  $y = x$  then  $c \leftarrow c + 1$ 
4   else if  $c = 0$  then  $y \leftarrow x, c \leftarrow 1$ 
5   else  $c \leftarrow c - 1$ 
6 return  $y$ 

```

Boyer and Moore [1981].

- If there is no m s.t. $f(m) > \frac{1}{2}$, then whatever y is correct.
- Assume $f(m) > \frac{1}{2}$. Whenever c reaches 0:
 - The algorithm goes back to initial state and starts to process the rest of the sequence.
 - m must be the majority of the rest sequence as well.

Misra-Gries Sketch

Algorithm: Majority(X)

```

1  $y \leftarrow \text{NaN}, c \leftarrow 0$ 
2 forall  $x \in X$  do
3   if  $y = x$  then  $c \leftarrow c + 1$ 
4   else if  $c = 0$  then  $y \leftarrow x, c \leftarrow 1$ 
5   else
6      $c \leftarrow c - 1$ 
7 return  $y$ 

```

Algorithm: Misra-Gries(X, k)

```

1  $Y \leftarrow [\text{NaN}] * k, C \leftarrow [0] * k$ 
2 forall  $x \in X$  do
3   if  $\exists i (Y[i] = x)$  then  $C[i] \leftarrow C[i] + 1$ 
4   else if  $\exists i (C[i] = 0)$  then  $Y[i] \leftarrow x, C[i] \leftarrow C[i] + 1$ 
5   else
6     forall  $i$  do  $C[i] \leftarrow C[i] - 1$ 
7 return  $Y, C$ 

```

Misra and Gries [1982].

Extends the majority algorithm by increasing the number of keys and counters from 1 to k .

Misra-Gries Sketch

Algorithm: Majority(X)

```

1  $y \leftarrow \text{NaN}, c \leftarrow 0$ 
2 forall  $x \in X$  do
3   if  $y = x$  then  $c \leftarrow c + 1$ 
4   else if  $c = 0$  then  $y \leftarrow x, c \leftarrow 1$ 
5   else
6      $c \leftarrow c - 1$ 
7 return  $y$ 

```

Algorithm: Misra-Gries(X, k)

```

1  $Y \leftarrow [\text{NaN}] * k, C \leftarrow [0] * k$ 
2 forall  $x \in X$  do
3   if  $\exists i(Y[i] = x)$  then  $C[i] \leftarrow C[i] + 1$ 
4   else if  $\exists i(C[i] = 0)$  then  $Y[i] \leftarrow x, C[i] \leftarrow C[i] + 1$ 
5   else
6     forall  $i$  do  $C[i] \leftarrow C[i] - 1$ 
7 return  $Y, C$ 

```

Misra and Gries [1982].

Extends the majority algorithm by increasing the number of keys and counters from 1 to k .

To approximate $f(q)$ for any $q \in [u]$

$$\hat{f}_{MG}(q) := \begin{cases} \frac{C[i]}{n} & \exists i(Y[i] = q) \\ 0 & \text{otherwise} \end{cases}$$

Misra-Gries Sketch Analysis

Lemma

for all $q \in [u]$,

$$f(q) - \frac{1}{k+1} \leq \hat{f}_{MG(k)}(q) \leq f(q)$$

Misra-Gries Sketch Analysis

Lemma

for all $q \in [u]$,

$$f(q) - \frac{1}{k+1} \leq \hat{f}_{MG(k)}(q) \leq f(q)$$

Proof.

The upper bound is obvious.

When Line 6 executes, there must be $k + 1$ distinct item are decremented. It can happen at most $n/(k + 1)$ times. □

Misra-Gries Sketch Analysis

Lemma

for all $q \in [u]$,

$$f(q) - \frac{1}{k+1} \leq \hat{f}_{MG(k)}(q) \leq f(q)$$

Setting $\frac{1}{k+1} = \varepsilon$, or $k = \frac{1}{\varepsilon} - 1$,

$$f(q) - \varepsilon \leq \hat{f}_{MG(k)}(q) \leq f(q) \leq f(q) + \varepsilon$$

Misra-Gries Sketch Analysis

Lemma

for all $q \in [u]$,

$$f(q) - \frac{1}{k+1} \leq \hat{f}_{MG(k)}(q) \leq f(q)$$

Setting $\frac{1}{k+1} = \varepsilon$, or $k = \frac{1}{\varepsilon} - 1$,

$$f(q) - \varepsilon \leq \hat{f}_{MG(k)}(q) \leq f(q) \leq f(q) + \varepsilon$$

- $\hat{f}_{MG(k)}(q)$ is an ε -approximation of $f(q)$.

Misra-Gries Sketch Analysis

Lemma

for all $q \in [u]$,

$$f(q) - \frac{1}{k+1} \leq \hat{f}_{MG(k)}(q) \leq f(q)$$

Setting $\frac{1}{k+1} = \varepsilon$, or $k = \frac{1}{\varepsilon} - 1$,

$$f(q) - \varepsilon \leq \hat{f}_{MG(k)}(q) \leq f(q) \leq f(q) + \varepsilon$$

- $\hat{f}_{MG(k)}(q)$ is an ε -approximation of $f(q)$.
- Y is a $\hat{H}_\varepsilon^\varepsilon$

Table of Contents

- 1 Problems
 - Frequency
 - Frequency Estimation in Stream
 - Heavy Hitters in Stream
- 2 Misra-Gries Sketch
 - Majority
 - Misra-Gries Sketch
- 3 **Count Sketch**
 - **Count Sketch Algorithm**
 - **Count Sketch Analysis**
- 4 Count-Min Sketch
- 5 Summary
 - Biased vs Unbiased

Count Sketch Algorithm

Algorithm: Count-Sketch(X, t, k)

- 1 $C \leftarrow 0^{t \times k}, H \leftarrow (H_i : [u] \rightarrow [k])_{i=1}^t, S \leftarrow (S_i : [u] \rightarrow [\pm 1])_{i=1}^t$
 - 2 **forall** $x \in X$ **do**
 - 3 **forall** i **in** $[t]$ **do**
 - 4 $C_{i, H_i(x)} \leftarrow C_{i, H_i(x)} + S_i(x)$
 - 5 **return** C, H, S
-

Charikar et al. [2002].

Count Sketch Algorithm

Algorithm: Count-Sketch(X, t, k)

- 1 $C \leftarrow 0^{t \times k}, H \leftarrow (H_i : [u] \rightarrow [k])_{i=1}^t, S \leftarrow (S_i : [u] \rightarrow [\pm 1])_{i=1}^t$
 - 2 **forall** $x \in X$ **do**
 - 3 **forall** i **in** $[t]$ **do**
 - 4 $C_{i, H_i(x)} \leftarrow C_{i, H_i(x)} + S_i(x)$
 - 5 **return** C, H, S
-

Charikar et al. [2002].

- H_i, S_i are independent hash functions.

Count Sketch Algorithm

Algorithm: Count-Sketch(X, t, k)

- 1 $C \leftarrow 0^{t \times k}, H \leftarrow (H_i : [u] \rightarrow [k])_{i=1}^t, S \leftarrow (S_i : [u] \rightarrow [\pm 1])_{i=1}^t$
 - 2 **forall** $x \in X$ **do**
 - 3 **forall** i **in** $[t]$ **do**
 - 4 $C_{i, H_i(x)} \leftarrow C_{i, H_i(x)} + S_i(x)$
 - 5 **return** C, H, S
-

Charikar et al. [2002].

- H_i, S_i are independent hash functions.
- S_i are chosen from a **pairwise** independent family.

Count Sketch Algorithm

Algorithm: Count-Sketch(X, t, k)

- 1 $C \leftarrow 0^{t \times k}, H \leftarrow (H_i : [u] \rightarrow [k])_{i=1}^t, S \leftarrow (S_i : [u] \rightarrow [\pm 1])_{i=1}^t$
 - 2 **forall** $x \in X$ **do**
 - 3 **forall** i in $[t]$ **do**
 - 4 $C_{i, H_i(x)} \leftarrow C_{i, H_i(x)} + S_i(x)$
 - 5 **return** C, H, S
-

$$H = \begin{bmatrix} H_1 \\ H_2 \\ \vdots \\ H_t \end{bmatrix} \quad S = \begin{bmatrix} S_1 \\ S_2 \\ \vdots \\ S_t \end{bmatrix} \quad C = \begin{bmatrix} C_{1,1} & C_{1,2} & \dots & C_{1,k} \\ C_{2,1} & C_{2,2} & \dots & C_{2,k} \\ \vdots & \vdots & \ddots & \vdots \\ C_{t,1} & C_{t,2} & \dots & C_{t,k} \end{bmatrix}$$

Count Sketch Query

To approximate $f(q)$ for any $q \in [u]$

$$\hat{f}_{CS}(q) := \operatorname{median}_{i \in [t]} \hat{f}_i(q), \quad \text{where } \hat{f}_i(q) := \frac{1}{n} S_i(q) C_{i, H_i(q)}.$$

$$H = \begin{bmatrix} H_1 \\ H_2 \\ \vdots \\ H_t \end{bmatrix} \quad S = \begin{bmatrix} S_1 \\ S_2 \\ \vdots \\ S_t \end{bmatrix} \quad C = \begin{bmatrix} C_{1,1} & C_{1,2} & \dots & C_{1,k} \\ C_{2,1} & C_{2,2} & \dots & C_{2,k} \\ \vdots & \vdots & \ddots & \vdots \\ C_{t,1} & C_{t,2} & \dots & C_{t,k} \end{bmatrix}$$

Count Sketch Randomness

Since the algorithm is not deterministic, it is randomized. We will analyze it in a probabilistic way.

Count Sketch Randomness

Since the algorithm is not deterministic, it is randomized. We will analyze it in a probabilistic way.

Question

Where is the randomness come from? Or what are the random variables?

Count Sketch Randomness

Since the algorithm is not deterministic, it is randomized. We will analyze it in a probabilistic way.

Question

Where is the randomness come from? Or what are the random variables?

Answer

The random events are the choices of hash functions in H and S .

Count Sketch Randomness

Since the algorithm is not deterministic, it is randomized. We will analyze it in a probabilistic way.

Question

Where is the randomness come from? Or what are the random variables?

Answer

The random events are the choices of hash functions in H and S . The random variables are H_i and S_i , or $H_i(q)$ and $S_i(q)$ for all $q \in [u]$.

Count Sketch Notations

$$\begin{aligned}C_{i,j} &:= \sum_{x \in X} S_i(x) \mathbb{1}(H_i(x), j) \\ &= \sum_{x \in [u]} nf(x) S_i(x) \mathbb{1}(H_i(x), j)\end{aligned}$$

Count Sketch Notations

$$\begin{aligned}C_{i,j} &:= \sum_{x \in X} S_i(x) \mathbb{1}(H_i(x), j) \\ &= \sum_{x \in [u]} nf(x) S_i(x) \mathbb{1}(H_i(x), j)\end{aligned}$$

- $\mathbb{1}(i, j)$: equal to 1 if $i = j$ and 0 otherwise.

Count Sketch Notations

$$\begin{aligned}
 C_{i,j} &:= \sum_{x \in X} S_i(x) \mathbb{1}(H_i(x), j) \\
 &= \sum_{x \in [u]} nf(x) S_i(x) \mathbb{1}(H_i(x), j)
 \end{aligned}$$

- $\mathbb{1}(i, j)$: equal to 1 if $i = j$ and 0 otherwise.
- $C_{i,j}^x := nf(x) S_i(x) \mathbb{1}(H_i(x), j)$: the part of $C_{i,j}$ caused by $x \in [u]$.

Count Sketch Notations

$$\begin{aligned} C_{i,j} &:= \sum_{x \in X} S_i(x) \mathbb{1}(H_i(x), j) \\ &= \sum_{x \in [u]} nf(x) S_i(x) \mathbb{1}(H_i(x), j) \end{aligned}$$

- $\mathbb{1}(i, j)$: equal to 1 if $i = j$ and 0 otherwise.
- $C_{i,j}^x := nf(x) S_i(x) \mathbb{1}(H_i(x), j)$: the part of $C_{i,j}$ caused by $x \in [u]$.

With these notations, we can simply write each $C_{i,j}$ as

$$C_{i,j} = \sum_{x \in [u]} C_{i,j}^x$$

Count Sketch Analysis - Mean

Lemma

For any $i \in [t], q \in [u]$,

$$E[\hat{f}_i(q)] = f(q)$$

Count Sketch Analysis - Mean

Lemma

For any $i \in [t], q \in [u]$,

$$E[\hat{f}_i(q)] = f(q)$$

$$\hat{f}_i(q) := \frac{1}{n} S_i(q) C_{i, H_i(q)} = f(q) + \frac{1}{n} \sum_{x \in [u], x \neq q} S_i(q) C_{i, H_i(q)}^x$$

Count Sketch Analysis - Mean

Lemma

For any $i \in [t], q \in [u]$,

$$E[\hat{f}_i(q)] = f(q)$$

$$\hat{f}_i(q) := \frac{1}{n} S_i(q) C_{i, H_i(q)} = f(q) + \frac{1}{n} \sum_{x \in [u], x \neq q} S_i(q) C_{i, H_i(q)}^x$$

$$\begin{aligned} & E [S_i(q) C_{i, H_i(q)}^x] \\ &= nf(x) E [S_i(q) S_i(x) \mathbb{1}(H_i(x), H_i(q))] \\ &= nf(x) E [S_i(q) S_i(x)] E [\mathbb{1}(H_i(x), H_i(q))] \quad S_i \text{ and } H_i \text{ are indep.} \\ &= nf(x) E [S_i(q)] E [S_i(x)] E [\mathbb{1}(H_i(x), H_i(q))] \quad S_i \text{ is pairwise indep.} \\ &= 0 \end{aligned}$$

Count Sketch Analysis - Variance

Lemma

For any $i \in [t], q \in [u]$,

$$V[\hat{f}_i(q)] \leq \frac{1}{k} F_2^2$$

where $F_2^2 = \sum_{x \in [u]} f(x)^2$.

Count Sketch Analysis - Variance

Lemma

For any $i \in [t], q \in [u]$,

$$V[\hat{f}_i(q)] \leq \frac{1}{k} F_2^2$$

where $F_2^2 = \sum_{x \in [u]} f(x)^2$.

$$V[\hat{f}_i(q)] = V \left[\frac{1}{n} S_i(q) C_{i, H_i(q)} \right] = \frac{1}{n^2} V \left[\sum_{x \in [u]} S_i(q) C_{i, H_i(q)}^x \right]$$

Count Sketch Analysis - Variance - 2

$$\mathbb{V}[\hat{f}_i(q)] = \frac{1}{n^2} \mathbb{V} \left[\sum_{x \in [u]} S_i(q) C_{i, H_i(q)}^x \right] = \frac{1}{n^2} \sum_{x \in [u]} \mathbb{V} \left[S_i(q) C_{i, H_i(q)}^x \right]$$

Because

$$\begin{aligned} & \text{cov} \left[S_i(q) C_{i, H_i(q)}^x, S_i(q) C_{i, H_i(q)}^y \right] \\ &= \mathbb{E} \left[\left(S_i(q) C_{i, H_i(q)}^x - \mathbb{E}[S_i(q) C_{i, H_i(q)}^x] \right) \left(S_i(q) C_{i, H_i(q)}^y - \mathbb{E}[S_i(q) C_{i, H_i(q)}^y] \right) \right] \\ &= \mathbb{E} \left[\left(S_i(q) C_{i, H_i(q)}^x \right) \left(S_i(q) C_{i, H_i(q)}^y \right) \right] = 0 \end{aligned}$$

for all $x \neq y$, if S_i and H_i are indep., S_i is pairwise indep.,

Count Sketch Analysis - Variance - 3

$$\begin{aligned}
 V[\hat{f}_i(q)] &= \frac{1}{n^2} \sum_{x \in [u]} V \left[S_i(q) C_{i, H_i(q)}^x \right] \\
 &\leq \frac{1}{n^2} \sum_{x \in [u]} E \left[(S_i(q) C_{i, H_i(q)}^x)^2 \right] \\
 &= \frac{1}{n^2} \sum_{x \in [u]} E \left[(nf(x) S_i(x) \mathbb{1}(H_i(x), H_i(q)))^2 \right] \\
 &= \sum_{x \in [u]} f^2(x) E \left[(\mathbb{1}(H_i(x), H_i(q)))^2 \right] \\
 &= \sum_{x \in [u]} f^2(x) \frac{1}{k} = \frac{1}{k} F_2^2
 \end{aligned}$$

Count Sketch Analysis - Failure Probability

Lemma

For any $q \in [u], i \in [t]$,

$$\Pr \left[\left| \hat{f}_i(q) - f(q) \right| \geq \varepsilon \right] \leq \frac{F_2^2}{k\varepsilon^2}$$

Count Sketch Analysis - Failure Probability

Lemma

For any $q \in [u], i \in [t]$,

$$\Pr \left[\left| \hat{f}_i(q) - f(q) \right| \geq \varepsilon \right] \leq \frac{F_2^2}{k\varepsilon^2}$$

The Chebyshev's inequality: $\Pr [|R - E[R]| \geq \varepsilon] \leq \frac{V[R]}{\varepsilon^2}$

Count Sketch Analysis - Failure Probability

Lemma

For any $q \in [u], i \in [t]$,

$$\Pr \left[\left| \hat{f}_i(q) - f(q) \right| \geq \varepsilon \right] \leq \frac{F_2^2}{k\varepsilon^2}$$

The Chebyshev's inequality: $\Pr [|R - E[R]| \geq \varepsilon] \leq \frac{V[R]}{\varepsilon^2}$

$$\Pr \left[\left| \hat{f}_i(q) - f(q) \right| \geq \varepsilon \right] = \Pr \left[\left| \hat{f}_i(q) - E[\hat{f}_i(q)] \right| \geq \varepsilon \right] \leq \frac{V[\hat{f}_i(q)]}{\varepsilon^2} \leq \frac{F_2^2}{k\varepsilon^2}$$

Count Sketch Analysis - Confidence Boosting

Now we know $\Pr \left[\left| \hat{f}_i(q) - f(q) \right| \geq \varepsilon \right] \leq \frac{F_2^2}{k\varepsilon^2}$ for each $i \in [t]$.

Count Sketch Analysis - Confidence Boosting

Now we know $\Pr \left[\left| \hat{f}_i(q) - f(q) \right| \geq \varepsilon \right] \leq \frac{F_2^2}{k\varepsilon^2}$ for each $i \in [t]$.

At the end, we will return

$$\hat{f}_{CS}(q) := \operatorname{median}_{i \in [t]} \hat{f}_i(q)$$

Count Sketch Analysis - Confidence Boosting

Now we know $\Pr \left[\left| \hat{f}_i(q) - f(q) \right| \geq \varepsilon \right] \leq \frac{F_2^2}{k\varepsilon^2}$ for each $i \in [t]$.

At the end, we will return

$$\hat{f}_{CS}(q) := \operatorname{median}_{i \in [t]} \hat{f}_i(q)$$

Why the median?

Count Sketch Analysis - Confidence Boosting

Now we know $\Pr \left[\left| \hat{f}_i(q) - f(q) \right| \geq \varepsilon \right] \leq \frac{F_2^2}{k\varepsilon^2}$ for each $i \in [t]$.

At the end, we will return

$$\hat{f}_{CS}(q) := \operatorname{median}_{i \in [t]} \hat{f}_i(q)$$

Why the median?

If the median has error $\geq \varepsilon$, then at least half of the $\hat{f}_i(q)$ have error $\geq \varepsilon$.

Count Sketch Analysis - Confidence Boosting 2

Let the failure probability of each $\hat{f}_i(q)$ is $p = \frac{F_2^2}{k\epsilon^2}$.
Repeat it t times independently, what is the probability of at least half failures?

Count Sketch Analysis - Confidence Boosting 2

Let the failure probability of each $\hat{f}_i(q)$ is $p = \frac{F_2^2}{k\epsilon^2}$.

Repeat it t times independently, what is the probability of at least half failures?

Binomial distribution! The number of failure is $B(t, p)$.

Count Sketch Analysis - Confidence Boosting 2

Let the failure probability of each $\hat{f}_i(q)$ is $p = \frac{F_2^2}{k\epsilon^2}$.

Repeat it t times independently, what is the probability of at least half failures?

Binomial distribution! The number of failure is $B(t, p)$.

Chernoff bound:

$$\Pr \left[B(t, p) \geq \frac{t}{2} \right] \leq \exp \left(- t(1/2 - p)^2 / (2p) \right)$$

Count Sketch Analysis - Confidence Boosting 3

Set $p = \frac{F_2^2}{k\epsilon^2} = \frac{1}{4}$, or $k = \frac{F_2^2}{4\epsilon^2}$:

$$\Pr \left[B(t, p) \geq \frac{t}{2} \right] \leq \exp \left(- t(1/2 - p)^2 / (2p) \right) \leq \exp(-t/8)$$

Count Sketch Analysis - Confidence Boosting 3

Set $p = \frac{F_2^2}{k\varepsilon^2} = \frac{1}{4}$, or $k = \frac{F_2^2}{4\varepsilon^2}$:

$$\Pr \left[B(t, p) \geq \frac{t}{2} \right] \leq \exp \left(-t(1/2 - p)^2 / (2p) \right) \leq \exp(-t/8)$$

Set $\exp(-t/8) = \delta$, or $t = 8 \log \frac{1}{\delta}$:

$$\Pr \left[\left| \hat{f}_{CS}(q) - f(q) \right| \geq \varepsilon \right] \leq \Pr[B(t, 1/4) \geq t/2] \leq \delta$$

Count Sketch Analysis - Confidence Boosting 3

Set $p = \frac{F_2^2}{k\epsilon^2} = \frac{1}{4}$, or $k = \frac{F_2^2}{4\epsilon^2}$:

$$\Pr \left[B(t, p) \geq \frac{t}{2} \right] \leq \exp \left(-t(1/2 - p)^2 / (2p) \right) \leq \exp(-t/8)$$

Set $\exp(-t/8) = \delta$, or $t = 8 \log \frac{1}{\delta}$:

$$\Pr \left[\left| \hat{f}_{CS}(q) - f(q) \right| \geq \epsilon \right] \leq \Pr[B(t, 1/4) \geq t/2] \leq \delta$$

Theorem

If $k = \frac{F_2^2}{4\epsilon^2}$ and $t = 8 \log \frac{1}{\delta}$, $\hat{f}_{CS}(q)$ is an $\hat{f}_\epsilon(q)$ with probability at least $1 - \delta$ for any $q \in [u]$.

Table of Contents

- 1 Problems
 - Frequency
 - Frequency Estimation in Stream
 - Heavy Hitters in Stream
- 2 Misra-Gries Sketch
 - Majority
 - Misra-Gries Sketch
- 3 Count Sketch
 - Count Sketch Algorithm
 - Count Sketch Analysis
- 4 **Count-Min Sketch**
- 5 Summary
 - Biased vs Unbiased

Count-Min Sketch - Algorithm

Algorithm: Count-Min(X, t, k)

- 1 $C \leftarrow 0^{t \times k}, H \leftarrow (H_i : [u] \rightarrow [k])_{i=1}^t$
 - 2 **forall** $x \in X$ **do**
 - 3 **forall** i **in** $[t]$ **do**
 - 4 $C_{i, H_i(x)} \leftarrow C_{i, H_i(x)} + 1$
 - 5 **return** C, H
-

Cormode and Muthukrishnan [2005]

To approximate $f(q)$ for any $q \in [u]$

$$\hat{f}_{CMS}(q) := \min_{i \in [t]} \hat{f}_i(q), \quad \text{where } \hat{f}_i(q) := \frac{1}{n} C_{i, H_i(q)}.$$

Count-Min Sketch - Bounds

Lemma

For any $q \in [u], i \in [t]$,

$$f(q) \leq \hat{f}_i(q).$$

If H_i is drawn from a pairwise independent hash family, then

$$\mathbb{E}[\hat{f}_i(q) - f(q)] \leq \frac{1}{k}.$$

If $k = \frac{1}{\delta\varepsilon}$, then

$$\Pr[\hat{f}_i(q) - f(q) \geq \varepsilon] \leq \delta$$

Count-Min Sketch - Confidence Boosting

Lemma

For any $q \in [u]$, $i \in [t]$, if H_i is drawn from a pairwise independent hash family, and $k = \frac{2}{\epsilon}$, then $\hat{f}_i(q)$ is a $\hat{f}_\epsilon(q)$ with probability at least $1/2$.

Theorem

If $t = \lg \frac{1}{\delta}$, $k = \frac{2}{\epsilon}$, H_i s are independently drawn from a pairwise independent hash family, then for any $q \in [u]$,
 $\hat{f}_{CMS}(q) := \min_{i \in [t]} \hat{f}_i(q)$ is a $\hat{f}_\epsilon(q)$ with probability at least $1 - \delta$.

Table of Contents

- 1 Problems
 - Frequency
 - Frequency Estimation in Stream
 - Heavy Hitters in Stream
- 2 Misra-Gries Sketch
 - Majority
 - Misra-Gries Sketch
- 3 Count Sketch
 - Count Sketch Algorithm
 - Count Sketch Analysis
- 4 Count-Min Sketch
- 5 Summary
 - Biased vs Unbiased

Summary

Sketch	Space	Technique	Deterministic
Misra-Gries	$O(1/\epsilon)$	Counter	Yes
Count Sketch	$O\left(\frac{F_2^2}{\epsilon^2} \log \frac{1}{\delta}\right)$	Hashing	No
Count-Min Sketch	$O\left(\frac{1}{\epsilon} \log \frac{1}{\delta}\right)$	Hashing	No

Table: Studied Sketches to obtain \hat{f}_ϵ (with probability at least $1 - \delta$ if applicable).

Summary

Sketch	Space	Technique	Deterministic
Misra-Gries	$O(1/\epsilon)$	Counter	Yes
Count Sketch	$O\left(\frac{F_2^2}{\epsilon^2} \log \frac{1}{\delta}\right)$	Hashing	No
Count-Min Sketch	$O\left(\frac{1}{\epsilon} \log \frac{1}{\delta}\right)$	Hashing	No

Table: Studied Sketches to obtain \hat{f}_ϵ (with probability at least $1 - \delta$ if applicable).

It seems like the Count-Min sketch is better than the Count sketch in the error-space tradeoff, but the bound is based on F_2^2 which is usually much smaller than 1. The Count sketch is also more versatile than Count-Min sketch and works very well in practice.

Biased vs Unbiased

Definition (biased, unbiased, under-estimated, over-estimated approximation)

To estimate a ground truth value f , a random variable \hat{f} (the output of any estimation method) is

- unbiased approximation if $E[\hat{f}] = f$;
- biased approximation if $E[\hat{f}] \neq f$;
- under-estimated approximation if $E[\hat{f}] < f$;
- over-estimated approximation if $E[\hat{f}] > f$;

question

Does unbiased approximation always better than biased approximation?

Questions

question

Can Count/Count-Min sketch solve heavy hitters? What is the query time?

Questions

question

Can Count/Count-Min sketch solve heavy hitters? What is the query time?

question

What is the failure probability of Count/Count-Min sketch actually is? For one q or for all $q \in [u]$?

Questions

question

Can Count/Count-Min sketch solve heavy hitters? What is the query time?

question

What is the failure probability of Count/Count-Min sketch actually is? For one q or for all $q \in [u]$?

question

What about weighted data? Real value weights? Negative weights?