# Lecture 11    10/1/2025

## Count-min and Count Sketches

In terms of frequency moments $F_\infty$ is the most frequently occuring item in the stream. It is a brittle measure. In most application we want to know the "heavy hitters", items that occur very frequently.

F... a theoretical perspective we

will call an index $i \in [n]$ a
heavy hitter if $\bar{f_i} \geq \alpha F_1 = \alpha m$
for some sufficiently large constant $\alpha \in (0,1)$

Alternatively $\bar{f_i} \geq \frac{m}{k}$ for some
integer $k$.

A classical algorithm shows that
one can identify items $i$ with
$\bar{f_i} \geq \frac{m}{k}$.

Misra-Gries $(k)$

— we have a data structure $D$ that stores $k$ items along with a counter for each. $D$ is initialized to empty set

— $m \leftarrow 0$

— while (stream is not empty) do

    $m \leftarrow m + 1$

    $e_m$ is current item

    If $e_m \in D$ then

        increment counter for $e_m$ in $D$

    Else

        If $D$ has $< k$ elements

            add $e_m$ to $D$ with counter value 1

        Else

            decrease counter value by

1 for all current elements
delete from $D$ any element
with counter set to 0

— end while

— Output values stored in $D$
and the counters values.

Implicitly it defines an estimate
$\tilde{f}_i$ for each $i$
if $i \in D$ at the end then
$\tilde{f}_i$ is the counter value
otherwise it is 0.

**Theorem:** $\forall i \in [n]$

$$\tilde{f}_i - \frac{m}{k+1} \leq \hat{f}_i \leq f_i$$

Hence if $i$ is a heavy hitter it will be in $D$. Space usage is $O(k)$.

Although Misra-Gries is nice it does not allow deletions and also does not provide a sketch.

Count-min and Count sketches are a way to use hashing to

identify heavy hitters and they have led to many applications

Basic idea is simple.
Suppose we use a hash function
$$h: [n] \rightarrow [ck] \text{ for some}$$
sufficiently large constant $k$.
Then $h$ spreads the $n$ items into $ck$ buckets. Suppose the heavy hitters are $i_1, i_2, \ldots, i_k$.
Then we expect that they will

not collide and we can use separate counts in each bucket.

We will use amplification as usual by considering multiple hash functions rather than a single one.

[Cormode-Muthukrishnan]

Count-Min Sketch $(w, d)$
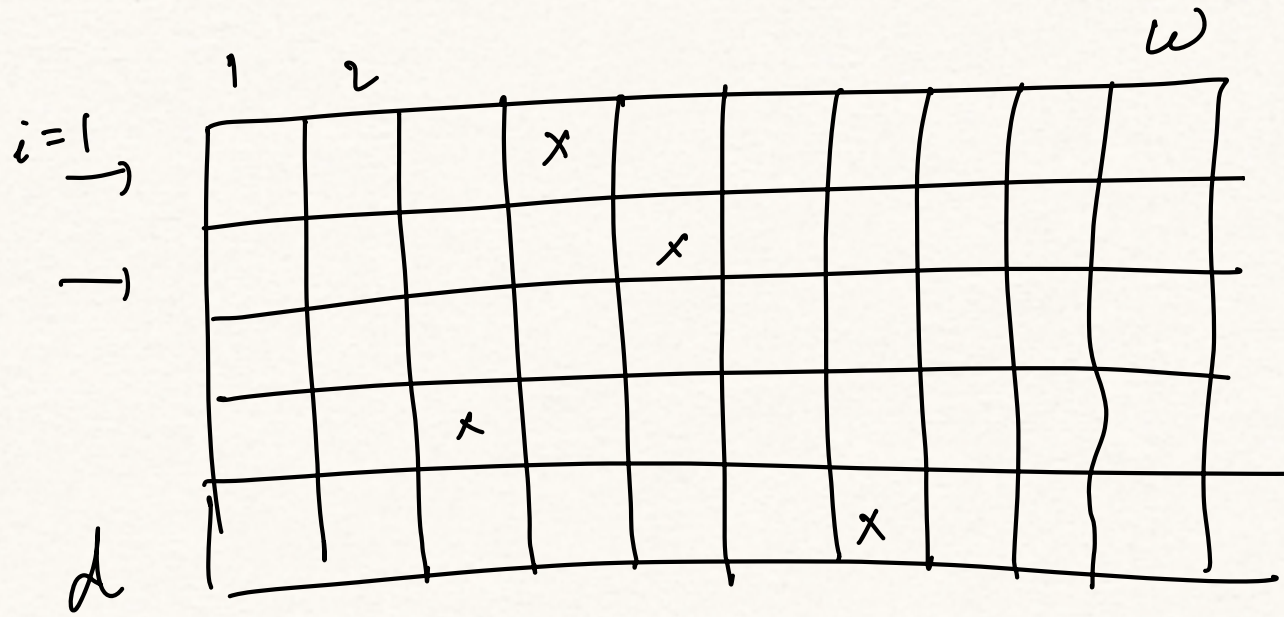— $h_1, h_2, \ldots, h_d$ are $d$ independent pairwise indep hash functions from $[n] \to [w]$.

- While (stream is not empty)

$e_t = (i_t, \Delta_t)$ is current item

for $\ell = 1$ to $d$ do

$$C[\ell, h(i_t)] \leftarrow C[\ell, h(i_t)] + \Delta_t .$$

→ end while

- for $i = 1$ to $[n]$

$$\tilde{X}_i = \min_{\ell=1}^{d} C[\ell, h(i)]$$

$w$ is width of the sketch

d is # of independent... 
we use.



$i=1$ →

→

d

$$C[d,s] \quad = \quad \sum_{i:\, h_d(i)=s} x_i$$

Lemma: Consider strict turnstile model ($\bar{X} \geq 0$). Let $d = \Omega(\ln \frac{1}{\delta})$ and $w \geq \frac{2}{\varepsilon}$. Then, $\forall i \in [n]$

(i) $\tilde{X}_i \geqslant X_i$

(ii) $\Pr\left[\tilde{X}_i \geqslant X_i + \varepsilon \|x\|_1\right] \leq \delta$.

**Proof:** Fix $i$.

For $\ell \in [d]$

$$Z_\ell = C[\ell, h_\ell(i)] = X_i + \sum_{i' \neq i: \, h_\ell(i') = h_\ell(i)} X_{i'}$$

$$Z_\ell - X_i = \sum_{i' \neq i: \, h_\ell(i') = h_\ell(i)} X_{i'}$$

$$E[Z_\ell - X_i] = \frac{1}{\omega} \sum_{i' \neq i} X_{i'} \leq \frac{\|x\|_1 - X_i}{\omega}.$$

By pairwise independence.

$$\mathbb{E}[Z_\ell - X_i] \le \frac{\varepsilon}{2} \|\bar{X}\|_1$$

By Markov
$$\Pr[Z_\ell - X_i] \ge \varepsilon \|\bar{X}\|_1 \le \frac{1}{2}.$$

Thus $\Pr\left[\min_\ell (Z_\ell - X_i)\right] \ge \varepsilon \|X\|_1$

$\uparrow$ by independence. $\le \left(\frac{1}{2}\right)^d \le \delta.$

$\square.$

Choosing $d = \Omega(\log n)$ we have

$$\tilde{X}_i \le X_i + \varepsilon \|X\|_1 \quad \forall i \in [n]$$

with high probability.

Count Min gives over estimates

Total space is $O(dw)$ counters

$$O\left(\frac{1}{\varepsilon}\log n\right).$$

Advantage: $\frac{1}{\varepsilon}$ dependence, simple

Disadvantage: only handles $\bar{x} \geq 0$.

Exercise: Show that Count Min is
linear sketch.

a when

# Count Sketch

Similar to Count Min in using $d$ independent hash functions but uses $F_2$ estimation ideas and median estimator instead of min.

## Count Sketch $(w, d)$

- $h_1, h_2, \ldots, h_d$ independent hash functions from $[n] \to [w]$
  hash functions

$- g_1, g_2, \ldots, g_d$ indep

from $[n] \to \{-1, +1\}$.

$\sigma_{\tau}$ While (stream is not empty) do

$$e_t \leftarrow (i_t, \Delta_t).$$

for $\ell = 1$ to $d$ do

$$C[\ell, h_\ell(i_t)] \leftarrow C[\ell, h_\ell(i_t)]$$
$$+ g_\ell(i_t) \Delta_t$$

end for

end while

$-$ for $i \in [n]$

$$\tilde{x}_i = \operatorname*{median}_{\ell=1}^{d} \left( g_\ell(i) C[\ell, h_\ell(i)] \right)$$

$\tilde{X}_i$ can be negative even of

$\overline{X} \geq 0$.

cancellations happen like ni $F_2$

estimation

Lemma: Let $d \geq 4 \ln \frac{1}{\delta}$ and

$\omega \geq \frac{3}{\varepsilon^2}$ . Then for any $i \in [n]$

(i) $E[\tilde{X}_i] = X_i$ and

(ii) $\Pr\left[ |\tilde{X}_i - X_i| \geq \varepsilon \|\bar{X}\|_2 \right] \leq \delta.$

**Proof:** Fix $i$. For $\ell \in [d]$

To make analysis easier, let

$Y_{i'}$ for $i' \neq i$ be the indicator

for $h_\ell(i) = h_\ell(i')$

$$Z_\ell = g_\ell(i) \, C[\ell, h_\ell(i)]$$

$$= g_\ell(i) \left[ g_\ell(i) x_i + \sum_{i' \neq i} g_\ell(i') Y_{i'} x_{i'} \right]$$

$$= X_i + \sum_{i' \neq i} g_\ell(i) g_\ell(i') Y_{i'} X_{i'}$$

$$E[Z_\ell] = X_i \qquad \text{by pairwise indep}$$
$$\text{of } g_\ell .$$

We note that $E[Y_{i'}] = \frac{1}{\omega}$

and $E[Y_{i'}^2] = \frac{1}{\omega}$ by pairwise indep of $h_\ell$.

$$\text{Var}(Z_\ell) = E[(Z_\ell - X_i)^2]$$

$$= E\left[\left(\sum_{i' \neq i} g_\ell(i) g_\ell(i') Y_{i'} x_{i'}\right)\right]$$

$$= \sum_{i' \neq i} x_{i'}^2 E[Y_{i'}^2]$$

$$= \frac{1}{\omega} \sum_{i' \neq i} x_{i'}^2 \quad .$$

$$\leq \frac{\|\bar{x}\|_2^2}{\omega} \cdot \leq \frac{\varepsilon^2}{3} \|\bar{x}\|_2^2$$

Hence using Chebyshev

$$E\left[|Z_\ell - x_i| \geq \varepsilon \|\bar{x}\|_2\right] \leq \frac{1}{3} .$$

Via Chernoff bounds

$$\Pr_\alpha \left[ \left| \text{med}\left( Z_1, Z_2, \dots Z_d \right) - X_i \right| \geqslant \varepsilon \| \bar{X} \|_2 \right]$$
$$\leq \delta.$$
$$\square$$

Important: Sketches do not store directly the "identity" of the heavy hitters. Given $i \in [n]$ we

can estimate $\tilde{x}_i$ from the sketch.
But outputting all $i$ such that
$\tilde{x}_i$ is high requires a linear
scan through $[n]$. Can maintain
multiple data structures and
use additional information to
find the heavy hitters in $\tilde{O}(k)$
space and time.

# Sparse Recovery

One nice and powerful application of Count Sketch is for sparse recovery. Suppose $\bar{x} \in R^n$ is sparse or close to sparse. Meaning that only $k$ of the coordinates are non-zero. Can we recover $x$ without knowing which of the coordinates are

going to be important? Want to use only $\tilde{O}(k)$ space.

Defn: Given $\bar{x} \in R^n$ let

$$\text{error}_2^k(\bar{x}) = \min_{\bar{z}\,:\,\|\bar{z}\|_0 \leq k} \|\bar{x} - \bar{z}\|_2 .$$

That is, what is the best $k$-space approximation to $\bar{x}$.

Offline, easy to compute.

$$z_i^* = \bar{X}_i \quad \text{if} \quad i \text{ is among the largest}$$

absolute value $k$ coordinates of $\bar{X}$

$$= 0 \quad \text{otherwise.}$$

Can we find $z^*$ in streaming setting?

**Theorem:** Count Sketch with $\omega = \frac{3K}{\varepsilon^2}$ and $d = \Omega(\log n)$ allows us to find a $\bar{Z}$ such that

$$\|\bar{Z}\|_0 \leq K \text{ and with high probability}$$

$$\|\bar{X} - \bar{Z}\|_2 < (1+\varepsilon) \, err_2^k(\bar{X}) .$$

In particular if $\bar{X}$ is $K$-Sparse then exact recovery.

# Compressed Sensing and RIP Matrices

Count Sketch guarantees that we can recover any sparse $\bar{x}$ with high probability. Can we guarantee probability 1 with a linear sketch?

Yes!

There exist $\ell \times n$ matrices $\Pi$ for

$$\ell = O\left(k \log \frac{n}{k}\right) \text{ such that}$$

given any $k$-sparse $\bar{x} \in R^n$ one can recover $\bar{x}$ from $\Pi \bar{x}$.

Note that $\Pi \bar{x}$ takes $O(l)$ space and since $l = O(k \log \frac{n}{k})$ we are not storing much more than what we want to recover.

Such matrices are called RIP matrices for "restricted isoperimetry property".

Turns out that a random $\ell \times n$ matrix with each entry chosen independently from a $N(0,1)$ Gaussian distribution satisfies the RIP. But cannot easily verify that a given matrix is RIP.

This area is called Compressed Sensing and has several applications in

signal processing.