# Sublinear Algorithms

## Lecture 03: Streaming III

Karl Bringmann

May 18, 2020

# Recap: Distinct Elements

approximate the number of distinct items among $x_1, \ldots, x_m \in [n]$

**Data structure problem:**

maintain set $D$ and its size $t$

**update($x$):** add $x$ to $D$

**query():** output $\tilde{t}$ with

$$\mathbb{P}\big[|\tilde{t} - t| \geq \varepsilon t\big] \leq \delta$$

**Idealized Setting:** FM (Flajolet, Martin 1985)

1) Pick random function $h: [n] \to [0,1]$

2) On update(): Maintain $X = \min_i h(x_i)$

3) On query(): Output $\tilde{t} = 1/X - 1$

Let $y_1, \ldots, y_t$ be the distinct items in the stream

Suppose that $y_1, \ldots, y_t$ are *random* in $[0,1]$

Then we expect $1/\min_i y_i \approx t$

# Recap: Distinct Elements

approximate the number of distinct items among $x_1, \dots, x_m \in [n]$

**Theoretical Variant:** (Bar-Yossef et al. 2002)

1) Pick a prime $m$ with $n^3 \le m \le n^{O(1)}$

2) Pick pairwise independent

hash function $h': [n] \to [m]$

3) Denote $h(x) := h'(x)/m \in (0,1]$

4) Maintain a set $L$ containing the $k := \lceil 36/\varepsilon^2 \rceil$

smallest distinct values among $h(x_1), \dots, h(x_m)$

5) On query():

   If $|L| < k$:   Output $\tilde{t} = |L|$

   Otherwise:   Output $\tilde{t} = k / \max(L)$

**Idealized Setting:** FM (Flajolet, Martin 1985)

1) Pick random function $h: [n] \to [0,1]$

2) On update(): Maintain $X = \min_i h(x_i)$

3) On query(): Output $\tilde{t} = 1/X - 1$

**Space usage:**   $O\left(\frac{1}{\varepsilon^2} \log n\right)$

# Recap: Pairwise Independence

Random variables $X_1, \ldots, X_n$ are **independent** if for any $j_1, \ldots, j_n$ we have

$$\mathbb{P}[X_1 = j_1 \text{ and } \ldots \text{ and } X_n = j_n] = \mathbb{P}[X_1 = j_1] \cdot \ldots \cdot \mathbb{P}[X_n = j_n]$$

$k = n$

Random variables $X_1, \ldots, X_n$ are **pairwise independent** if
for any $i \neq i'$ the random variables $X_i$ and $X_{i'}$ are independent.

$k = 2$

**Lem:** For pairwise independent $X_1, \ldots, X_n$ we have

$$\mathrm{Var}[X_1 + \cdots + X_n] = \mathrm{Var}[X_1] + \cdots + \mathrm{Var}[X_n]$$

Random variables $X_1, \ldots, X_n$ are $k$-**wise independent** if
for any distinct $i_1, \ldots, i_k$ the random variables $X_{i_1}, \ldots, X_{i_k}$ are independent.

# Recap: Pairwise Independent Hash Function

Let $p$ be a prime with $p \geq n$, and pick $h \in \mathcal{H}$ uniformly at random

Let $\mathcal{H}$ be the set of all functions $h: [n] \to [m]$ of the form

$h(i) = (a \cdot i + b) \bmod p$ where $a, b \in [p]$

Each hash value $h(i)$ is **uniformly distributed** in $[m]$

The random variables $h(1), \dots, h(n)$ are **pairwise independent**

A function $h \in \mathcal{H}$ can be represented by the pair $(a, b) \in [p]^2$, using $2\lceil \log p \rceil$ bits

We can sample a function $h \in \mathcal{H}$ in time $O(1)$

# k-wise Independent Hash Function

Let $p$ be a prime with $p \geq n$, and pick $h \in \mathcal{H}$ uniformly at random

Let $\mathcal{H}$ be the set of all functions $h: [n] \to [p]$ of the form

$$h(i) = \left(a_0 + a_1 \cdot i + \cdots + a_{k-1} \cdot i^{k-1}\right) \bmod p \ \text{ where } a_0, \ldots, a_{k-1} \in [p]$$

Each hash value $h(i)$ is **uniformly distributed** in $[m]$

The random variables $h(1), \ldots, h(n)$ are $k$-**wise** **independent**

A function $h \in \mathcal{H}$ can be represented by the tuple $(a_0, \ldots, a_{k-1}) \in [p]^k$, using $k\lceil \log p \rceil$ bits

We can sample a function $h \in \mathcal{H}$ in time $O(k)$

# k-wise Independent Hash Function

Let $p$ be a prime with $p \geq n$, and pick $h \in \mathcal{H}$ uniformly at random

Let $\mathcal{H}$ be the set of all functions $h \colon [n] \to [p]$ of the form

$$h(i) = \left(a_0 + a_1 \cdot i + \cdots + a_{k-1} \cdot i^{k-1}\right) \bmod p \quad \text{where } a_0, \dots, a_{k-1} \in [p]$$

$$\left\{ 0, \dots, 4 \frac{1}{5} \right\}$$

$$\{-1, 1\}$$

**Arbitrary Codomain:** Can we get codomain $Y = \{y_1, \dots, y_t\}$ for $t \ll p$?

$\sigma(i) \coloneqq y_{h(i) \bmod t}$ is *almost uniform*, that is, $\mathbb{P}[\sigma(i) = y_j] = 1/t \pm O(1/p)$

**Lem:** Fix $k$. For any $n$ and $Y$, there is a family $\mathcal{H}$ of functions from $[n]$ to $Y$ such that

- $h \in \mathcal{H}$ can be stored using $O(\log(n + |Y|))$ bits, and sampled in time $O(1)$,

- for random $h \in \mathcal{H}$ and fixed $i$, the value $h(i)$ is (almost) **uniformly distributed** in $Y$,

- for random $h \in \mathcal{H}$, the random variables $h(1), \dots, h(n)$ are $k$-**wise independent.**

# Outline

1) **Turnstile Model + Moment Estimation**

2) Point Query + Heavy Hitters

# Generalized Streaming Model

maintain a vector $x \in \mathbb{Z}^n$ under updates of the form $x_i = x_i + \Delta$

**General data structure problem:**

maintain vector $x \in \mathbb{Z}^n$

**update($i, \Delta$):** $x_i = x_i + \Delta$

**query():** approximate $f(x)$

**Insertion-only:**     Each update has $\Delta = 1$

This is what we studied so far

(E.g. distinct elements: #non-zero $x_i$'s)

$z_1, \ldots, z_m \in [n]$

$x_i = |\{j \mid z_j = i\}|$

# Generalized Streaming Model

maintain a vector $x \in \mathbb{Z}^n$ under updates of the form $x_i = x_i + \Delta$

**General data structure problem:**

maintain vector $x \in \mathbb{Z}^n$

**update($i, \Delta$):** $x_i = x_i + \Delta$

**query():** approximate $f(x)$

**Insertion-only:** Each update has $\Delta = 1$

This is what we studied so far

(E.g. distinct elements: #non-zero $x_i$'s)

One-way

**Strict Turnstile:** $\Delta \in \mathbb{Z}$ (may be negative!)

Promise: $x_i \geq 0$ for all $i$ at all times

two-way
closed room

**General Turnstile:** $\Delta \in \mathbb{Z}, x_i \in \mathbb{Z}$

two-way
open

$x_1$ $x_2$ $x_3$ $x_4$

# Second Moment Estimation

approximate $F_2 = \|x\|_2^2 = \sum_{i=1}^n x_i^2$ for a vector $x \in \mathbb{Z}^n$ given in turnstile model

**Data structure problem:**

maintain vector $x \in \mathbb{Z}^n$

**update($i, \Delta$):** $x_i = x_i + \Delta$

**query():** output $\tilde{F}_2$ with

$$\mathbb{P}\left[\left|\tilde{F}_2 - F_2\right| \geq \varepsilon F_2\right] \leq \delta$$

**AMS Sketch:** (Alon, Matias, Szegedy 1999)

1) Pick 4-wise independent
   hash function $\sigma: [n] \rightarrow \{-1, 1\}$

2) Maintain $y = \sum_{i=1}^n \sigma(i) \cdot x_i$

3) Output $y^2$

initially $x = (0, \ldots, 0)$

assume $|\Delta| = n^{O(1)}$

so entries of $x$ are $O(\log n)$-bit integers

update$(i, \Delta)$: $y = y + \sigma(i) \cdot \Delta$

# Second Moment Estimation

approximate $F_2 = \|x\|_2^2 = \sum_{i=1}^n x_i^2$ for a vector $x \in \mathbb{Z}^n$ given in turnstile model

**Standard Approach:**

> **Lem:** AMS Sketch is an unbiased estimator,
> that is, $\mathbb{E}[y^2] = \|x\|_2^2$

**AMS Sketch:** (Alon, Matias, Szegedy 1999)

1) Pick 4-wise independent $\sigma: [n] \rightarrow \{-1, 1\}$

2) Maintain $y = \sum_{i=1}^n \sigma(i) \cdot x_i$

3) Output $y^2$

**Proof:**
$$\mathbb{E}[y^2] = \mathbb{E}[(\sum_{i=1}^n \sigma(i)x_i)^2]$$

$$= \mathbb{E}[\sum_{i,j} \sigma(i)\sigma(j)x_i x_j]$$

$$= \sum_i \mathbb{E}[\sigma(i)^2]x_i^2 + \sum_{i \neq j} \mathbb{E}[\sigma(i)\sigma(j)]x_i x_j$$

$$= \boxed{\sum_i x_i^2}$$

*independent*

$$y^2 = \sum_i x_i^2 + \sum_{i \neq j} \sigma(i)\sigma(j)x_i x_j$$

$$= \mathbb{E}[\sigma(i)] \cdot \mathbb{E}[\sigma(j)]$$

$$b$$

$$= 0$$

...since $\sigma$ is pairwise independent

# Second Moment Estimation

approximate $F_2 = \|x\|_2^2 = \sum_{i=1}^n x_i^2$ for a vector $x \in \mathbb{Z}^n$ given in turnstile model

**Standard Approach:**

> **Lem:** AMS Sketch is an unbiased estimator, that is, $\mathbb{E}[y^2] = \|x\|_2^2$

> **Lem:** We have $\mathbb{E}[y^4] \leq 3\mathbb{E}[y^2]^2$.

**AMS Sketch:** (Alon, Matias, Szegedy 1999)

1) Pick 4-wise independent $\sigma: [n] \to \{-1,1\}$

2) Maintain $y = \sum_{i=1}^n \sigma(i) \cdot x_i$

3) Output $y^2$

**Proof:** $\mathbb{E}[y^4] = \mathbb{E}[(\sum_{i=1}^n \sigma(i)x_i)^4] \;\dot{=}\; \sum_{i=1}^n \mathbb{E}[(\sigma(i)x_i)^4] + 6\sum_{i<j} \mathbb{E}[(\sigma(i)x_i)^2] \cdot \mathbb{E}\left[(\sigma(j)x_j)^2\right]$

$\underbrace{\phantom{\mathbb{E}[(\sigma(i)x_i)^2]}}_{x_i^2} \quad \underbrace{\phantom{\mathbb{E}[(\sigma(j)x_j)^2]}}_{x_j^2}$

...since $\sigma$ is 4-wise independent and has expectation 0, see Exercise Sheet 1

$= \sum_{i=1}^n x_i^4 + 6\sum_{i<j} x_i^2 \cdot x_j^2 \;\leq\; 3\left(\sum_{i=1}^n x_i^4 + 2\sum_{i<j} x_i^2 \cdot x_j^2\right) \;=\; 3\left(\sum_{i=1}^n x_i^2\right)^2 = 3\mathbb{E}[y^2]^2$

# Second Moment Estimation

approximate $F_2 = \|x\|_2^2 = \sum_{i=1}^n x_i^2$ for a vector $x \in \mathbb{Z}^n$ given in turnstile model

**Standard Approach:**

> **Lem:** AMS Sketch is an unbiased estimator, that is, $\mathbb{E}[y^2] = \|x\|_2^2$

> **Lem:** We have $\mathbb{E}[y^4] \leq 3\mathbb{E}[y^2]^2$.

**AMS Sketch:** (Alon, Matias, Szegedy 1999)

1) Pick 4-wise independent $\sigma : [n] \to \{-1,1\}$

2) Maintain $y = \sum_{i=1}^n \sigma(i) \cdot x_i$

3) Output $y^2$

By Chebyshev: $\mathbb{P}[|y^2 - F_2| \geq \varepsilon F_2] \leq \dfrac{2}{\varepsilon^2}$

> **Boosting via Chebyshev:** AMS+ = average over $\dfrac{8}{\varepsilon^2}$ runs of AMS has error prob. $1/4$

> **Boosting via Chernoff:** AMS++ = median of $8\log(2/\delta)$ runs of AMS has error prob. $\delta$

# Second Moment Estimation

approximate $F_2 = \|x\|_2^2 = \sum_{i=1}^n x_i^2$ for a vector $x \in \mathbb{Z}^n$ given in turnstile model

**Data structure problem:**

maintain vector $x \in \mathbb{Z}^n$

**update($i, \Delta$):** $x_i = x_i + \Delta$

**query():** output $\tilde{F}_2$ with

$$\mathbb{P}\left[\left|\tilde{F}_2 - F_2\right| \geq \varepsilon F_2\right] \leq \delta$$

**Space usage:** $O\left(\frac{1}{\varepsilon^2} \log\left(\frac{1}{\delta}\right) \log n\right)$ bits

**AMS Sketch:** (Alon, Matias, Szegedy 1999)

1) Pick 4-wise independent $\sigma : [n] \to \{-1, 1\}$

2) Maintain $y = \sum_{i=1}^n \sigma(i) \cdot x_i$

3) Output $y^2$

**AMS+:** average over $O(1/\varepsilon^2)$ runs of AMS

**AMS++:** median of $O(\log(1/\delta))$ runs of AMS+

# Remark 1: Moment Estimation

approximate $F_p = \|x\|_p^p = \sum_{i=1}^{n} |x_i|^p$ for a vector $x \in \mathbb{Z}^n$ given in turnstile model

For $0 \leq p \leq 2$ and constant $\delta$:  There is streaming algorithm using $\mathrm{poly}(\varepsilon^{-1} \log n)$ space.

For $p > 2$ and constant $\varepsilon, \delta$:  Space complexity is $n^{1-2/p}$ up to logfactors.

[Alon,Matias,Szegedy'99, Bar-Yossef,Jayram,Kumar,Sivakumar'04, Indyk,Woodruff'05, Indyk'06]

# Remark 2: Linear Sketch

... is an algorithm that maintains $\Pi x$, for some matrix $\Pi \in \mathbb{R}^{m \times n}$

Want to maintain vector $x \in \mathbb{Z}^n$

Pick a suitable matrix $\Pi$
and instead maintain $y \in \mathbb{Z}^m$ with $y = \Pi x$

$$y = \Pi \cdot x$$

We cannot store $\Pi$ explicitly!

**Implicit representation** of $\Pi$: Given $i, j$ we can efficiently compute the entry $\Pi_{i,j}$

On update$(i, \Delta)$: $y = y + \Delta \cdot \Pi_i$, where $\Pi_i$ is the $i$-th column of $\Pi$ $\quad \rightarrow$ time $O(m)$

$\Pi$ may be **randomized** (that is, chosen from some probability distribution)

# Remark 2: Linear Sketch

... is an algorithm that maintains $\Pi x$, for some matrix $\Pi \in \mathbb{R}^{m \times n}$

Want to maintain vector $x \in \mathbb{Z}^n$

Pick a suitable matrix $\Pi$
and instead maintain $y \in \mathbb{Z}^m$ with $y = \Pi x$

$$y = \Pi \cdot x$$

Any algorithm in strict/general turnstile model can be converted into a linear sketch, at the cost of at most a logarithmic factor in the space bound.

[Li,Nguyen,Woodruff'14]

# Remark 2: Linear Sketch

approximate $F_2 = \|x\|_2^2 = \sum_{i=1}^{n} x_i^2$ for a vector $x \in \mathbb{Z}^n$ given in turnstile model

**Data structure problem:**

maintain vector $x \in \mathbb{Z}^n$

**update($i, \Delta$):** $x_i = x_i + \Delta$

**query():** output $\tilde{F}$ with

$$\mathbb{P}\big[\big|\tilde{F} - F_2\big| \geq \varepsilon F_2\big] \leq \delta$$

**AMS+ Sketch:** (Alon, Matias, Szegedy 1999)

1) Pick 4-wise independent

   hash functions $\sigma_1, \dots, \sigma_s : [n] \to \{-1,1\}$

2) Maintain $y_j = \sum_{i=1}^{n} \sigma_j(i) \cdot x_i$

3) Output $\frac{1}{s} \sum_{j=1}^{s} y_j^2$

This is a linear sketch!

$(y_1, \dots, y_s) = y = \Pi x$ with $\Pi_{j,i} = \sigma_j(i)$

$\Pi$ is implicitly represented

# Outline

1) Turnstile Model + Moment Estimation

2) **Point Query + Heavy Hitters**

$\#\{i \mid x_i > \varepsilon\|x\|_1\} \leq 1/\varepsilon$

# Point Query

approximate $x_i \pm \varepsilon\|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in turnstile model

**Data structure problem:**

maintain vector $x \in \mathbb{Z}^n$

**update$(i, \Delta)$:** $x_i = x_i + \Delta$

**query$(i)$:** output $\tilde{x}_i$ with

$\mathbb{P}[|\tilde{x}_i - x_i| \geq \varepsilon\|x\|_1] \leq \delta$

**CountMin Sketch:** (Cormode,Muthukrishnan'05)

1) Pick 2-wise independent

hash function $h: [n] \to [t]$ for $t = \lceil 4/\varepsilon \rceil$

2) Maintain counters $C_j = \sum_{i \text{ s.t. } h(i)=j} x_i$

That is, initially $C_1, \ldots, C_t = 0$

On update$(i, \Delta)$: Add $\Delta$ to $C_{h(i)}$

3) On query$(i)$: output $C_{h(i)}$

$\left( x_i \pm 1/3 \longrightarrow x_i \begin{smallmatrix} \nearrow 0 \\ \searrow 1 \end{smallmatrix} \longrightarrow n \text{ bits} \right)$

$\|x\|_1 = \sum_i |x_i|$

# Point Query

approximate $x_i \pm \varepsilon \|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in turnstile model

**Lem:** $\mathbb{P}\left[\left|C_{h(i)} - x_i\right| > \varepsilon\|x\|_1\right] \leq 1/4$

**CountMin Sketch:** (Cormode,Muthukrishnan'05)

1) Pick 2-wise independent
   hash function $h: [n] \to [t]$ for $t = \lceil 4/\varepsilon \rceil$

2) Maintain counters $C_j = \sum_{i \text{ s.t. } h(i)=j} x_i$

   That is, initially $C_1, \dots, C_t = 0$

   On update$(i, \Delta)$: Add $\Delta$ to $C_{h(i)}$

3) On query$(i)$: output $C_{h(i)}$

**Proof:** Fix $i$. For $j \neq i$:

$$Z_j = \begin{cases} 1, & \text{if } h(j) = h(i) \\ 0, & \text{otherwise} \end{cases}$$

$$\mathbb{E}[Z_j] = \mathbb{P}[Z_j = 1] = 1/t$$

$$C_{h(i)} = x_i + \sum_{j \neq i} x_j Z_j$$

$$\left|C_{h(i)} - x_i\right| \leq \sum_{j \neq i} |x_j| \cdot Z_j$$

(Markov)

$$\mathbb{P}\left[\left|C_{h(i)} - x_i\right| > \varepsilon\|x\|_1\right] \leq \mathbb{P}\left[\sum_{j \neq i} |x_j| \cdot Z_j > \varepsilon\|x\|_1\right] \leq \frac{\sum_{j \neq i} |x_j|/t}{\varepsilon\|x\|_1} \leq \frac{1}{\varepsilon t} \leq \frac{1}{4}$$

$\geq 0$

# Point Query

approximate $x_i \pm \varepsilon\|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in turnstile model

**Data structure problem:**

maintain vector $x \in \mathbb{Z}^n$

**update($i, \Delta$):** $x_i = x_i + \Delta$

**query($i$):** output $\tilde{x}_i$ with

$\mathbb{P}[|\tilde{x}_i - x_i| \geq \varepsilon\|x\|_1] \leq \delta$

**Space usage:** $O\left(\frac{1}{\varepsilon}\log\left(\frac{1}{\delta}\right)\log n\right)$

In **strict** turnstile model we can use *minimum* instead of *median*

**CountMin Sketch:** (Cormode, Muthukrishnan'05)

1) Pick 2-wise independent
   hash function $h: [n] \to [t]$ for $t = \lceil 4/\varepsilon\rceil$

2) Maintain counters $C_j = \sum_{i \text{ s.t. } h(i)=j} x_i$
   That is, initially $C_1, \ldots, C_t = 0$
   On update($i, \Delta$): Add $\Delta$ to $C_{h(i)}$

3) On query($i$): output $C_{h(i)}$

**CM++:** median of $O(\log(1/\delta))$ runs of CM

$\|x\|_1 = \sum_i |x_i|$

$\leq 1/\varepsilon$ many

# Heavy Hitters

Goal: $O(\mathrm{polylog}\, n)$ space

compute all $i$ with $x_i \geq \varepsilon\|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in **strict** turnstile model

**Data structure problem:**

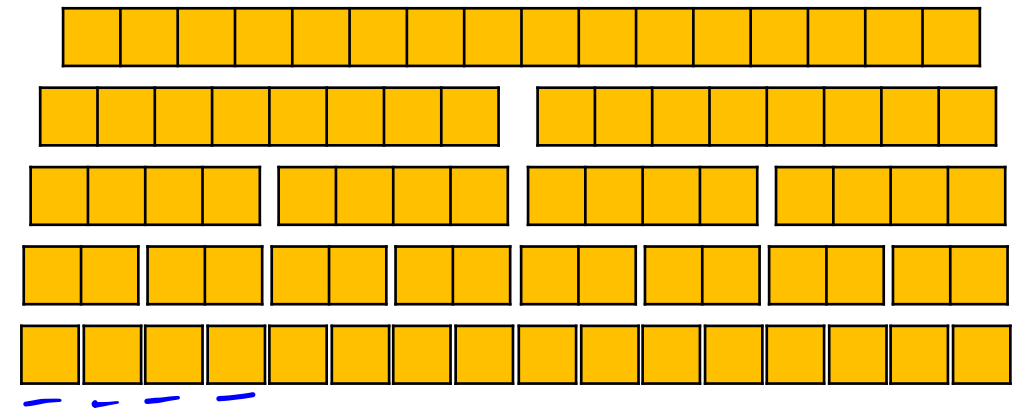maintain vector $x \in \mathbb{Z}^n$

**update($i, \Delta$):** $x_i = x_i + \Delta$

**query($x$):** output a set $R$ s.t.

$R$ contains **all** $i$ with $x_i \geq \varepsilon\|x\|_1$ .

$R$ contains **no** $i$ with $x_i < \frac{\varepsilon}{2}\|x\|_1$ .

with failure probability $\delta$

# Heavy Hitters: Dyadic Trick

compute all $i$ with $x_i \geq \varepsilon \|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in **strict** turnstile model

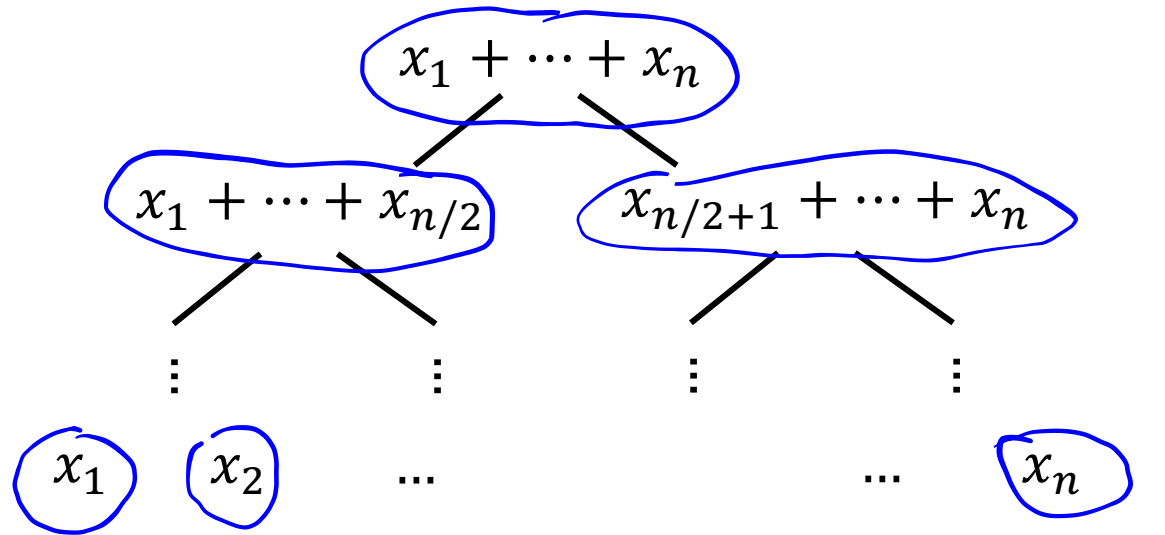**Data structure problem:**

maintain vector $x \in \mathbb{Z}^n$

**update($i, \Delta$):** $x_i = x_i + \Delta$

**query($i$):** output a set $R$ s.t.

$R$ contains **all** $i$ with $x_i \geq \varepsilon \|x\|_1$

$R$ contains **no** $i$ with $x_i < \frac{\varepsilon}{2}\|x\|_1$

with failure probability $\delta$

$$x_1 + \cdots + x_n$$

$$x_1 + \cdots + x_{n/2} \qquad x_{n/2+1} + \cdots + x_n$$

$$\vdots \qquad \vdots \qquad \vdots \qquad \vdots$$

$$x_1 \qquad x_2 \qquad \cdots \qquad \cdots \qquad x_n$$
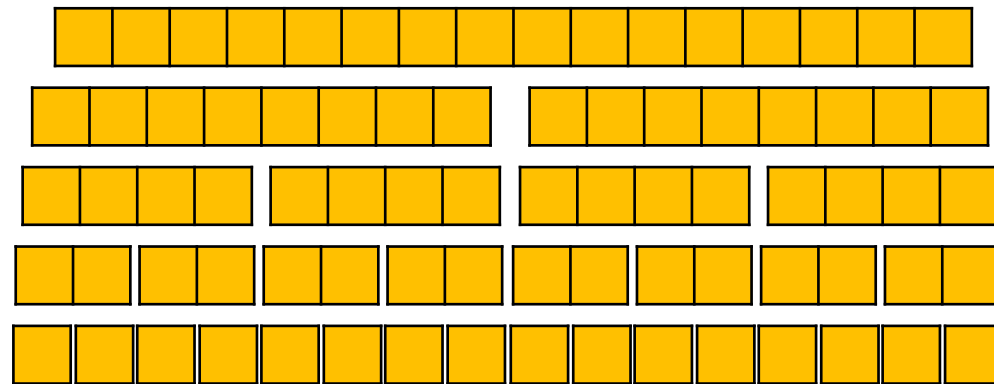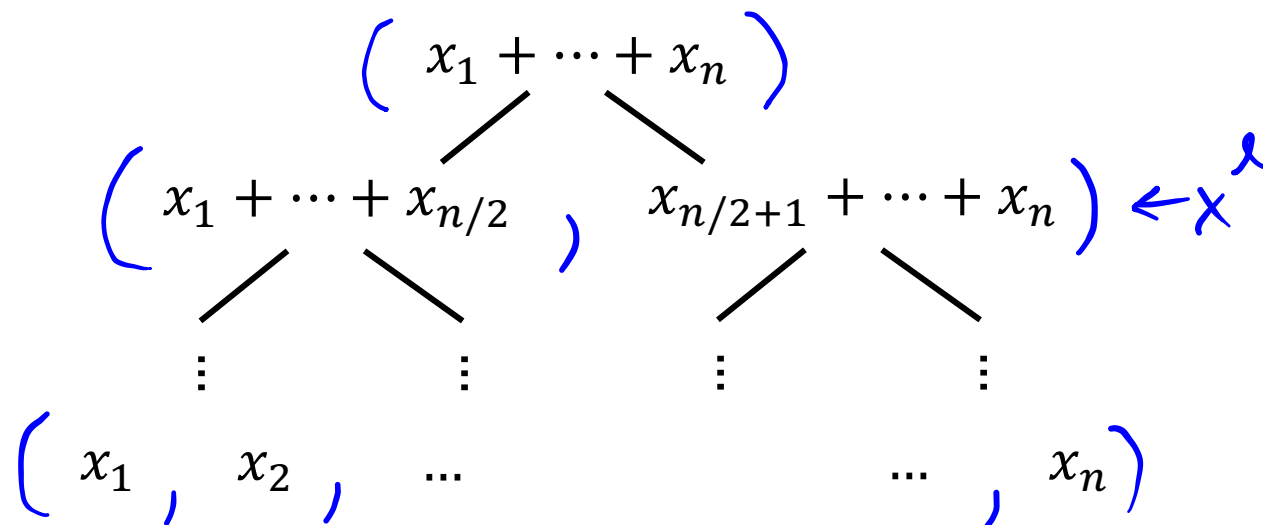
# Heavy Hitters: Dyadic Trick

compute all $i$ with $x_i \geq \varepsilon \|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in **strict** turnstile model

Level $\ell$ corresponds to vector $\boxed{x^\ell} \in \mathbb{Z}^{2^\ell}$ with

$$x_i^\ell = x_{(i-1)n/2^\ell + 1} + \cdots + x_{in/2^\ell} \leftarrow$$

**Store a CM++ sketch for every vector $x^\ell$**

with $\mathbb{P}\left[ \left| \tilde{x}_i^\ell - x_i^\ell \right| \geq \frac{\varepsilon}{4} \|x^\ell\|_1 \right] \leq \frac{\delta \varepsilon}{4 \log n}$

$$\left( x_1 + \cdots + x_n \right)$$

$$\left( x_1 + \cdots + x_{n/2} \right) \quad \left. x_{n/2+1} + \cdots + x_n \right) \leftarrow x^\ell$$

$$\left( x_1, \quad x_2, \quad \ldots \qquad \ldots, \quad x_n \right)$$

# Heavy Hitters: Dyadic Trick

compute all $i$ with $x_i \geq \varepsilon \|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in **strict** turnstile model

Level $\ell$ corresponds to vector $x^\ell \in \mathbb{Z}^{2^\ell}$ with

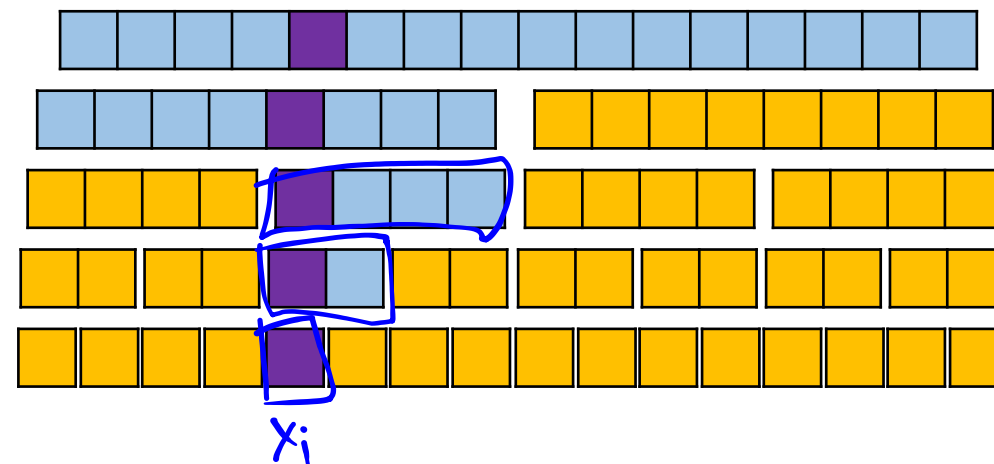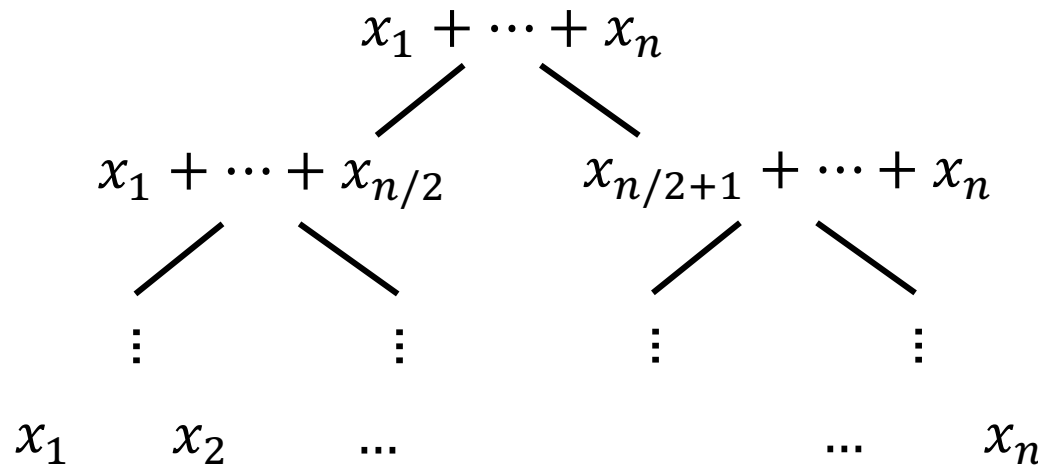$$x_i^\ell = x_{(i-1)n/2^\ell + 1} + \cdots + x_{in/2^\ell}$$

**Store a CM++ sketch for every vector $x^\ell$**

with $\mathbb{P}\left[ \left| \tilde{x}_i^\ell - x_i^\ell \right| \geq \frac{\varepsilon}{4} \|x^\ell\|_1 \right] \leq \frac{\delta \varepsilon}{4 \log n}$



$$x_1 + \cdots + x_n$$

$$x_1 + \cdots + x_{n/2} \qquad x_{n/2+1} + \cdots + x_n$$

$$x_1 \qquad x_2 \qquad \ldots \qquad \ldots \qquad x_n$$

**On update($i, \Delta$):**

Affects one entry of every vector $x^\ell$

So perform $O(\log n)$ updates of CM++

# Heavy Hitters: Dyadic Trick

compute all $i$ with $x_i \geq \varepsilon \|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in **strict** turnstile model

Level $\ell$ corresponds to vector $x^\ell \in \mathbb{Z}^{2^\ell}$ with

$$x_i^\ell = x_{(i-1)n/2^\ell+1} + \cdots + x_{in/2^\ell}$$

**Store a CM++ sketch for every vector $x^\ell$**

with $\mathbb{P}\left[\left|\tilde{x}_i^\ell - x_i^\ell\right| \geq \frac{\varepsilon}{4}\|x^\ell\|_1\right] \leq \frac{\delta\varepsilon}{4\log n}$

DFS($\boxed{x_i^\ell}$): *node $(\ell, i)$*

    Use CM++ sketch to decide if $x_i^\ell \geq \frac{3}{4}\varepsilon\|x\|_1$

    If CM++ says „larger": *// $x_i^\ell$ is heavy hitter*

        If $x_i^\ell$ is a leaf:  add $i$ to result $R$

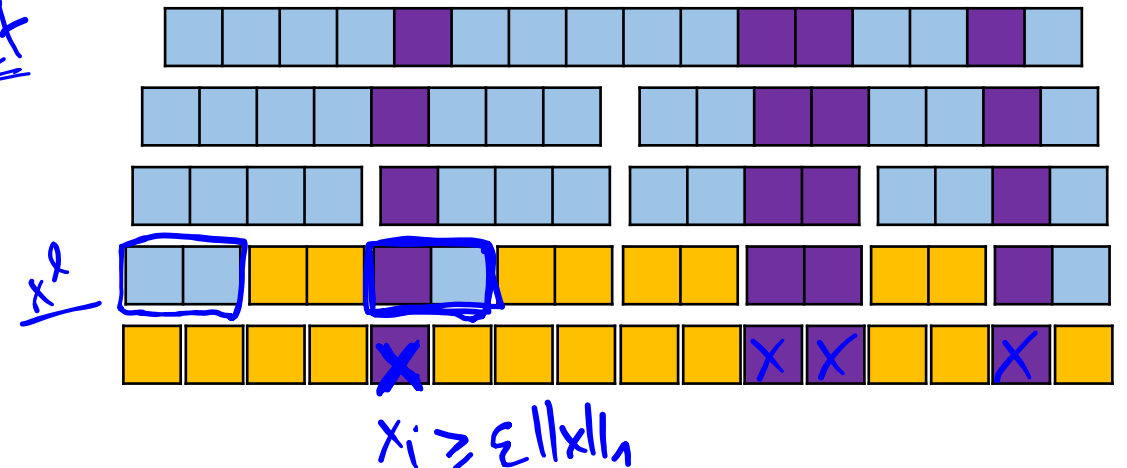        Else:  run DFS($x_{2i-1}^{\ell+1}$) and DFS($x_{2i}^{\ell+1}$)

**On query():**

Same norm on every level: $\left\|x^\ell\right\|_1 = \|x\|_1$ → *strict*

Every node is less than or equal to its parent

So all heavy hitters on all levels form a **subtree**

There are $O(1/\varepsilon)$ heavy hitters on each level



$x^\ell$

$x_i \geq \varepsilon\|x\|_1$

# Heavy Hitters: Dyadic Trick

compute all $i$ with $x_i \geq \varepsilon \|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in **strict** turnstile model

Level $\ell$ corresponds to vector $x^\ell \in \mathbb{Z}^{2^\ell}$ with

$$x_i^\ell = x_{(i-1)n/2^\ell+1} + \cdots + x_{in/2^\ell}$$

**Store a CM++ sketch for every vector $x^\ell$**

with $\mathbb{P}\left[\left|\tilde{x}_i^\ell - x_i^\ell\right| \geq \dfrac{\varepsilon}{4}\|x^\ell\|_1\right] \leq \dfrac{\delta\varepsilon}{4\log n}$

DFS($x_i^\ell$):

    Use CM++ sketch to decide if $x_i^\ell \geq \dfrac{3}{4}\varepsilon\|x\|_1$

    If CM++ says „larger":

        If $x_i^\ell$ is a leaf:  add $i$ to result $R$

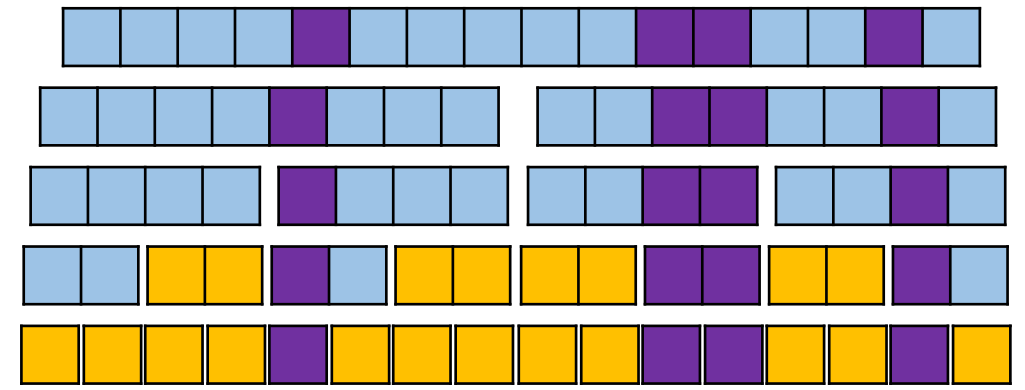        Else:  run DFS($x_{2i-1}^{\ell+1}$) and DFS($x_{2i}^{\ell+1}$)

**On query():**

Assuming correctness of CM++:

    CM++ says „larger" only if $x_i^\ell \geq \dfrac{\varepsilon}{2}\|x\|_1$

    On each level we explore the children of $\leq 2/\varepsilon$ nodes

So $\leq \dfrac{4}{\varepsilon}\log n$ calls to CM++, so error probability $\delta$
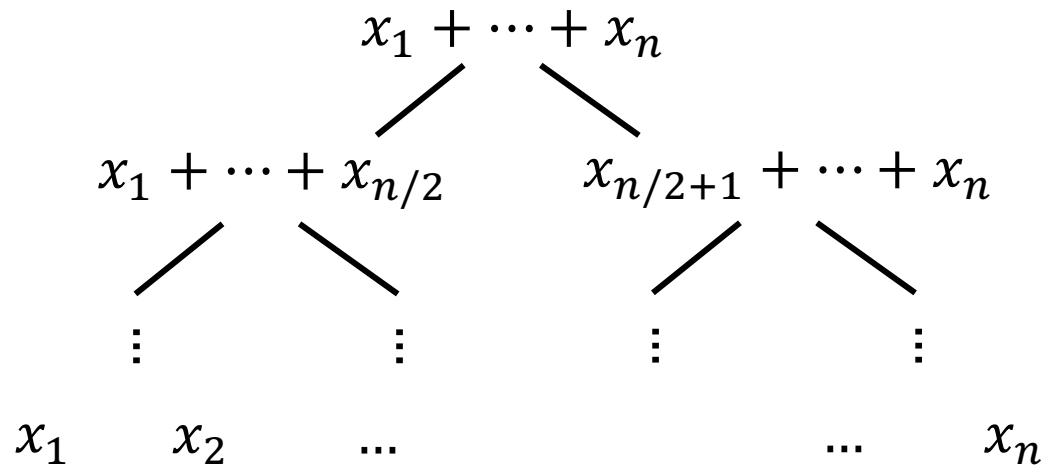
# Heavy Hitters: Dyadic Trick

compute all $i$ with $x_i \geq \varepsilon \|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in **strict** turnstile model

Level $\ell$ corresponds to vector $x^\ell \in \mathbb{Z}^{2^\ell}$ with

$$x_i^\ell = x_{(i-1)n/2^\ell+1} + \cdots + x_{in/2^\ell}$$

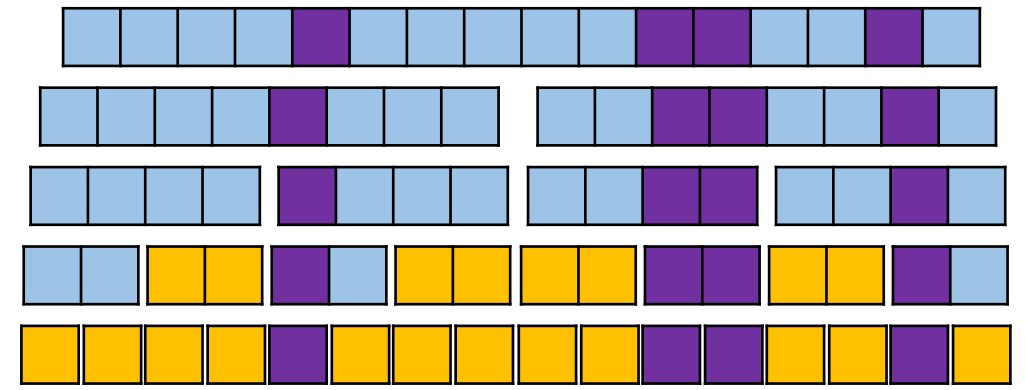**Store a CM++ sketch for every vector $x^\ell$**

with $\mathbb{P}\left[\left|\tilde{x}_i^\ell - x_i^\ell\right| \geq \frac{\varepsilon}{4}\|x^\ell\|_1\right] \leq \frac{\delta\varepsilon}{4\log n}$

CM++: $O\left(\frac{1}{\varepsilon} \log \frac{1}{\delta} \cdot \log n\right)$

**Space:** $O(\log n) \cdot O\left(\frac{1}{\varepsilon} \log\left(\frac{\log n}{\delta\varepsilon}\right) \log n\right)$ bits

#levels

$$= O\left(\frac{1}{\varepsilon} \log^2 n \left(\log\log n + \log\left(\frac{1}{\delta\varepsilon}\right)\right)\right)$$

$x_1 + \cdots + x_n$

$x_1 + \cdots + x_{n/2}$     $x_{n/2+1} + \cdots + x_n$

$\vdots$    $\vdots$    $\vdots$    $\vdots$

$x_1$     $x_2$     $\ldots$     $\ldots$     $x_n$

# Heavy Hitters

Goal: $O(\text{polylog}\, n)$ space

compute all $i$ with $x_i \geq \varepsilon\|x\|_1$ for a vector $x \in \mathbb{Z}^n$ given in **strict** turnstile model

**Data structure problem:**

maintain vector $x \in \mathbb{Z}^n$

**update**$(i, \Delta)$: $x_i = x_i + \Delta$

**query**(): output a set $R$ s.t.

$R$ contains **all** $i$ with $x_i \geq \varepsilon\|x\|_1$

$R$ contains **no** $i$ with $x_i < \frac{\varepsilon}{2}\|x\|_1$

with failure probability $\delta$

**Space:** $O\left(\frac{1}{\varepsilon}\log^2 n\left(\log\log n + \log\left(\frac{1}{\delta\varepsilon}\right)\right)\right)$ bits

Level $\ell$ corresponds to vector $x^\ell \in \mathbb{Z}^{2^\ell}$ with

$$x_i^\ell = x_{(i-1)n/2^\ell+1} + \cdots + x_{in/2^\ell}$$

**Store a CM++ sketch for every vector $x^\ell$**

with $\mathbb{P}\left[\left|\tilde{x}_i^\ell - x_i^\ell\right| \geq \frac{\varepsilon}{4}\|x^\ell\|_1\right] \leq \frac{\delta\varepsilon}{4\log n}$

CM++: $O(\log 1/\delta)$

**Update time:** $O\left(\log n\left(\log\log n + \log\left(\frac{1}{\delta\varepsilon}\right)\right)\right)$

**Query time:** $O\left(\frac{1}{\varepsilon}\log n\left(\log\log n + \log\left(\frac{1}{\delta\varepsilon}\right)\right)\right)$

# More Material

*Moment estimation, AMS Sketch:*

[Alon, Matias, Szegedy „The space complexity of approximating the frequency moments" 1999]

*Point Query + Heavy Hitters, CountMin Sketch:*

[Cormode, Muthukrishnan „An improved data stream summary: the count-min sketch and its applications" 2005]

— Course Website → Material → Link to Summer School on Streaming by Jelani Nelson

— **Exercise Sheet 1** due on **Friday, May 22**

**Nect lecture by Vasileios Nakos on May 28!**