

Linear Sketching

CountMin Sketch

Count Sketch

EECS 477

Seth Pettie

Streaming & Sketching

- You need to maintain a vector $x \in \mathbb{Z}^u$ of integers. Initially $x = 0$.
 - $\text{Update}(i, \Delta)$: Set $x(i) = x(i) + \Delta$.
 - $\text{Query}(i)$: Return $x(i)$.
 - *Look at other types of queries later...*
- “Incremental” : all Δ s are positive.
- “Strict Turnstyle” : positive and negative Δ , but $\forall i. x(i) \geq 0$ at all times.
- $\Theta(u)$ space is necessary and sufficient. Problem solved!

Streaming & Sketching

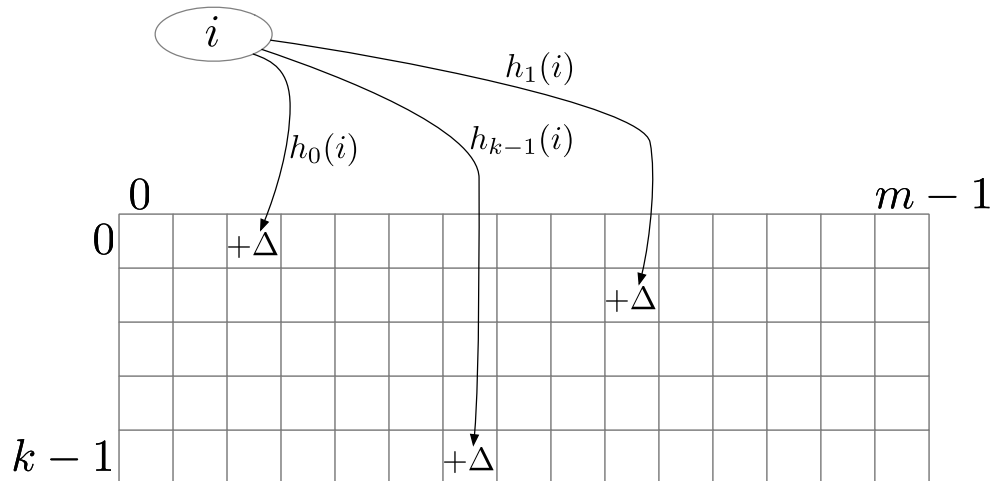
- You need to maintain a vector $x \in \mathbb{Z}^u$ of integers. Initially $x = 0$.
 - $\text{Update}(i, \Delta)$: Set $x(i) = x(i) + \Delta$.
 - $\text{Query}(i)$: Return $\tilde{x}(i) = x(i) \pm \text{Err}(x)$.
 - *Look at other types of queries later...*
- “Incremental” : all Δ s are positive.
- “Strict Turnstyle” : positive and negative Δ , but $\forall i. x(i) \geq 0$ at all times.
- Now how much space do you need? Depends a lot on what “ $\text{Err}(x)$ ” is...

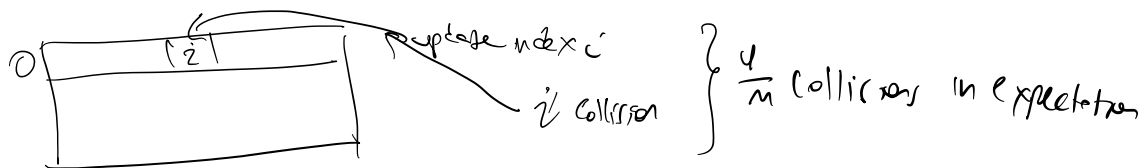
CountMin Sketch

- **Theorem.** A data structure of size $O(\epsilon^{-1} \log u)$ can handle **Update/Query** operations in $O(\log u)$ time. **Query** returns $\tilde{x}(i) = x(i) \pm \|x\|_1$.
- “ ℓ_1 -norm” $\|x\|_1 = \sum_{i=0}^{u-1} x(i)$. *Manhattan distance*
- “ ℓ_2 -norm” (aka Euclidean norm) $\|x\|_2 = \sqrt{\sum_{i=0}^{u-1} x(i)^2}$.
- $F_2 = 2^{\text{nd}}$ moment $= \|x\|_2^2 = \sum_{i=0}^{u-1} x(i)^2$.

CountMin Sketch

- Choose k hash functions $h_1, \dots, h_k: [u] \rightarrow [m]$ from a 2-independent family.
- Allocate an $k \times m$ array A , initially all 0.
- Update(i, Δ) :
 - For($j = 0; j < k; j++$)
 - $A[j, h_j(i)] += \Delta$;





$$I_{i'} : h_0(i) = h_0(i')$$

indicator for every i'

$$E(I_{i'}) = \begin{cases} 1 & i' = i \\ \frac{1}{m} & i' \neq i \end{cases}$$

$$\text{Query}(i) = \sum_j m_j(A[j, h_j(i)])$$

$$= O(k)$$

$$E(A[0, h_0(i)])$$

$$= E\left(\sum_{i'=0}^{m-1} I_{i'} \cdot x(i')\right)$$

constant

$$= \sum_{i'=0}^{m-1} E(I_{i'}) \cdot x(i') \rightarrow x \text{ norm except } x(i)$$

$$= x(i) + \left(\frac{\|x\|_1 - x(i)}{m}\right) \text{ positive noise}$$

$$\downarrow$$

$$\left(\frac{1}{m}\right) (x(0) + \dots + x(m-1)) \rightarrow \text{except } x(i)$$

- What is $E(A[0, h_0(i)]) = ?$
 - Hint: indicators!
- $I_{i'}$: an indicator for the event that $h_0(i) = h_0(i')$.
- $E(A[0, h_0(i)]) = \sum_{i'=0}^{u-1} E(I_{i'}) \cdot x(i')$
- $= x(i) + (\|x\|_1 - x(i))/m$
- How should we implement $\text{Query}(i)$?
- $\text{Query}(i)$
 - Return $\min\{A[0, h_0(i)], \dots, A[k-1, h_{k-1}(i)]\}$.

$$P(A[0, h_0(i)] - x(i) \geq \text{err}) \leq \frac{\mathbb{E}(A[0, h_0(i)] - x(i))}{\text{err}} \quad \text{Markov's}$$

$$\leq \frac{\|x\|_2 / m}{\text{err}} \stackrel{\text{want}}{=} \frac{1}{2}$$

$$\text{err} = \left(\frac{2}{M}\right) \|x\|_2$$

At 12:46

$$P(\forall j A[0, h_0(j)] - x(j) \geq \text{err}) \leq \left(\frac{1}{2}\right)^k = \frac{1}{\mu^2} \quad k = 2 \log \mu$$

- Define $\text{Err} = \left(\frac{2}{m}\right) \|x\|_1$
- What is $\Pr(\text{Query}(i) \notin [x(i), x(i) + \text{Err}])$?
- $\Pr\left(A[j, h_j(i)] - x(i) > 2 \cdot E\left(A[j, h_j(i)] - x(i)\right)\right) < \frac{1}{2}$.
Markov's Ineq.
 $- E\left(A[j, h_j(i)] - x(i)\right) < \frac{\|x\|_1}{m} = \frac{\text{Err}}{2}$.
- Thus, $\Pr(A[j, h_j(i)] > x(i) + \text{Err}) < \frac{1}{2}$.
- Finally, $\Pr(\forall j. A[j, h_j(i)] > x(i) + \text{Err}) < \left(\frac{1}{2}\right)^k$.
- Set $k = 2\log u, m = 2\epsilon^{-1}$, then
 $- \forall i. \text{Query}(i) \in [x(i), x(i) + \epsilon \|x\|_1]$ with prob. $1 - 1/u$.

Heavy Hitters

- Want to implement

$$\epsilon m = \epsilon \|x\|_1 \quad m \geq 2\epsilon^{-1}$$

- HeavyHitters() : return a short list L that includes all (2ϵ) -heavy hitters: i s.t. $x(i) \geq 2\epsilon\|x\|_1$.

- For i from 0 to $u - 1$:

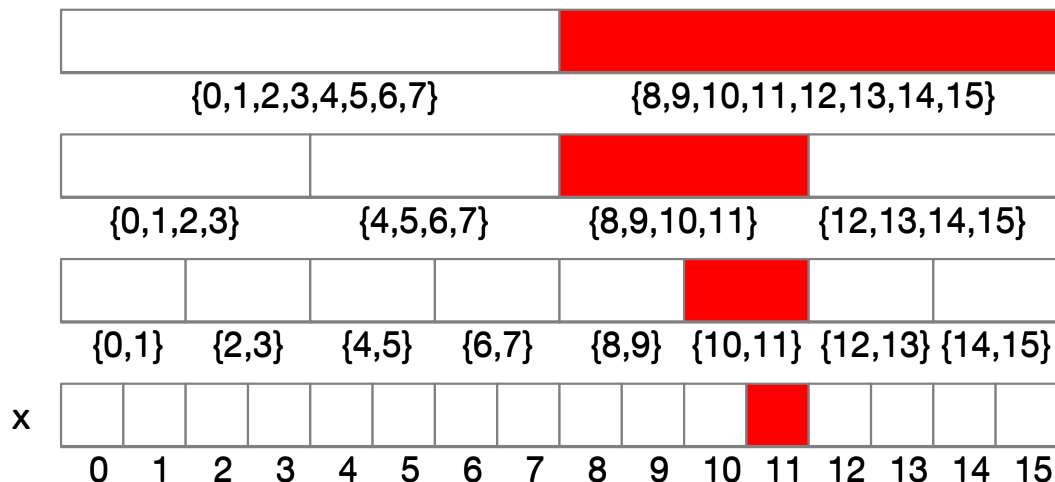
- If $\text{Query}(i) \geq 2\epsilon\|x\|_1$, add i to L .

- **Claim:** $|L| \leq \epsilon^{-1}$.

- **Proof:** Any $i \in L$ has $x(i) \geq \epsilon\|x\|_1$.

Heavy Hitters

- Want to implement
 - HeavyHitters() : return a short list L that includes all (2ϵ) -heavy hitters: i s.t. $x(i) \geq 2\epsilon\|x\|_1$.
- At most ϵ^{-1} indices i at each level with $x(i) \geq 2\epsilon\|x\|_1$.
- $O(\epsilon^{-1} \log^2 u)$ time to find them all.



Build CountMin structure for each of $\log u$ vectors

If an element is a 2ϵ -HH, its parent is also a 2ϵ -HH.

If an element is a 2ϵ -HH, 0, 1, or 2 of its children Could be 2ϵ -HHs.

Count Sketch

- CountMin Sketch: ℓ_1 error; good for finding **2ϵ -heavy hitters**. (i s.t. $x(i) \geq 2\epsilon\|x\|_1$)
- Count Sketch: ℓ_2 error; good for finding ℓ_2 **2ϵ -heavy hitters**: i s.t $x(i) \geq 2\epsilon\|x\|_2$.
- Pick hash functions $h_0, \dots, h_{k-1}: [u] \rightarrow [m]$ as before.
- Pick hash functions $g_0, \dots, g_{k-1}: [u] \rightarrow \{-1, 1\}$ from a 2-wise independent family.
- Update(i, Δ) :
 - For($j = 0$; $j < k$; $j++$)
 - $A[j, h_j(i)] += g_j(i) \cdot \Delta$;

$$x = (\sqrt{n}, \sqrt{n}, 1, 1, 1, \overbrace{\quad\quad\quad}^{\sim 25}, 1, 0, \dots)$$

$\|x\|_2 = \sqrt{2n}$ ℓ_2 norm picks up different heavy hitters $\epsilon = 0.01$

$$\begin{aligned}
 E(g_0(i) \cdot A[0, h_0(i)]) &= E\left(\sum_{i'} I_{i'} \cdot g_0(i') \cdot x(i')\right) \\
 &= \underbrace{g_0(i)^2}_{\substack{\downarrow \\ 0}} x(i) + \sum_{i' \neq i} E(I_{i'} \cdot g_0(i')) \cdot x(i') \\
 &= g_0(i) x(i) \\
 &= x(i)
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(\quad) &= E\left((g_0(i) \cdot A[0, h_0(i)] - x(i))^2\right) \\
 &= E\left(\left(\sum_{i' \neq i} g_0(i) \cdot g_0(i') \cdot x(i') \cdot I_{i'}\right)^2\right)
 \end{aligned}$$

$$= E\left(\sum_{i'=i''} x(i)^2 \cdot I_{i'} + \sum_{i' \neq i''} \overbrace{g_0(i)^2}^{=1} \overbrace{g_0(i) g_0(i'')}^{=1} \overbrace{x(i') x(i'')}^{=1} \overbrace{I_{i'} I_{i''}}^{=1}\right)$$

$$= E\left(\sum_{i'=i''} x(i)^2 \cdot I_{i'}\right) = \sum_i x(i)^2 \cdot E(I_i) \leq \frac{\|x\|_2^2}{m} = \frac{F_2}{m}$$

$$\begin{aligned}
 P(|g_0(i) \cdot A[0, h_0(i)] - x(i)|^2 > t^2) &\stackrel{\text{Markov's}}{\leq} \frac{E(g_0(i) \cdot A[0, h_0(i)] - x(i))^2}{t^2} \\
 &\leq \frac{\|x\|_2^2}{m t^2} = \frac{1}{3}
 \end{aligned}$$

$$t = \varepsilon \|x\|_2 \quad m = 3 \cdot \varepsilon^{-2}$$

Query(i) : $g_0(i) \cdot A[0, h_0(i)]$
 \vdots
 $g_{k-1}(i) \cdot A[k-1, h_{k-1}(i)]$
 get the median of this
 $K = \log n$

Count Sketch

- What is $E(A[0, h_0(i)]) = ?$
- $= E(\sum_{i'} I_{i'} \cdot g(i')x(i'))$
- $= g(i)x(i) + \sum_{i'} E(I_{i'} \cdot g(i'))x(i') = g(i)x(i).$
- $\text{Var}(g(i)A[0, h_0(i)]) =$
 $E\left(\left(\sum_{i' \neq i} g(i)g(i')I_{i'} \cdot x(i')\right)^2\right)$
- $= \sum_{i' \neq i} E(I_{i'}) \cdot x(i')^2 +$
 $\sum_{i' \neq i''} g(i')g(i'')I_{i'}I_{i''}x(i')x(i'')$
- $\leq \|x\|_2^2/m = F_2/m.$

$$\Pr(|g_0(i) A[\Sigma_0, h_0(i)] - x(i)| > t) < \frac{\mathbb{E}[(A[\Sigma_0, h_0(i)] - x(i))^2]}{t^2}$$

square working
chance

$$= t = \epsilon \|x\|_2, \quad m = \epsilon^2 \cdot 3$$

$$\Rightarrow \frac{\|x\|_2^2 / m}{t^2} = \frac{1}{3}$$

(p now we get good approx with prob 2/3

• Query (i)

• Return $\text{median}_{j \leq k} \{g_j(i) A[\Sigma_j, h_j(i)]\}$

Claim: if $k \geq 60 \ln n$, Query returns $x(i) \pm \epsilon \|x\|_2$
with prob $1 - 1/n^2$

Chernoff Bounds

→ pf included

- $X = X_1 + \dots + X_n$, $X_i \in \{0,1\}$, and $\Pr(X_i = 1) = p$.
- $\Pr(X > (1 + \delta)np) < \left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right)^{np} < e^{-\frac{\delta^2}{2+\delta}np}$.
- $\Pr(X < (1 - \delta)np) < \left(\frac{e^{-\delta}}{(1-\delta)^{1-\delta}}\right)^{np} < e^{-\frac{\delta^2}{2}np}$.

Back to Count Sketch

- $X = X_1 + \dots + X_k$ defined so that
 - $X_j = \begin{cases} 0 & \text{if } g_j(i)A[j, h_j(i)] \text{ is good approx of } x(i) \\ 1 & \text{if } g_j(i)A[j, h_j(i)] \text{ is bad approx of } x(i) \end{cases}$
 - $\Pr(X_j = 1) = \frac{1}{3} = p$. $E(X) = k/3$.
- $\text{Query}(i)$ could return a bad approx. if $X \geq k/2$.
- $\Pr(X \geq (1.5)E(X)) \leq \left(\frac{e^{\frac{1}{2}}}{1.5^{1.5}}\right)^{E(X)} < e^{-.1\left(\frac{k}{3}\right)} < \frac{1}{u^2}$
- The last inequality holds for $k = 60 \ln u$.