

Understanding $|\text{PDy}|_\infty$

- $E[(\text{PDy})_i] = 0$ for each hash bucket i , and $E[(\text{PDy})_i^2] = O\left(\frac{1}{s}\right) (n^{1-\frac{2}{p}} |y|_p^2)$
- Bernstein's bound: Suppose R_1, \dots, R_n are independent, and for all j , $|R_j| \leq K$, and $\text{Var}[\sum_j R_j] = \sigma^2$. There are constants C, c , so that for all $t > 0$,

$$\Pr\left[\left|\sum_j R_j - E\left[\sum_j R_j\right]\right| > t\right] \leq C \left(e^{-\frac{ct^2}{\sigma^2}} + e^{-\frac{ct}{K}}\right)$$
- Recall $(\text{PDy})_i = \sum_j \delta(h(j) = i) \cdot \sigma_j \cdot (\text{Dy})_j$, and set $R_j = \delta(h(j) = i) \cdot \sigma_j \cdot (\text{Dy})_j$
- Want $|\text{PDy}|_\infty \approx |\text{Dy}|_\infty$, where $|\text{Dy}|_\infty \in \left[\frac{|y|_p}{10^{1/p}}, 10^{1/p} |y|_p\right]$ with probability $> 4/5$
- Set $t = \frac{|y|_p}{100}$ and $s = \Theta\left(n^{1-\frac{2}{p}} \log n\right)$, to get $\frac{1}{n^2}$ error probability in Bernstein's bound
- But what is $K = \max_j |R_j|$?

Understanding the Large Elements

- Recall $(PDy)_i = \sum_j \delta(h(j) = i) \cdot \sigma_j \cdot (Dy)_j$, and set $R_j = \delta(h(j) = i) \cdot \sigma_j \cdot (Dy)_j$
- We will separately handle those R_j for which $|R_j| > \frac{\alpha|y|_p}{\log n}$, for a sufficiently small constant $\alpha > 0$. If $|R_j| > \frac{\alpha|y|_p}{\log n}$, then necessarily $|(Dy)_j| \geq \frac{\alpha|y|_p}{\log n}$
 - We call such a j **large** if $|(Dy)_j| \geq \frac{\alpha|y|_p}{\log n}$, otherwise j is **small**. How many indices j are large?
- Recall: $|(Dy)_j| = |y_j|/E_j^{1/p}$
- $\Pr_D[|(Dy)_j| \text{ is large}] = \Pr\left[\frac{|y_j|}{E_j^{1/p}} \geq \frac{\alpha|y|_p}{\log n}\right] = \Pr\left[\frac{|y_j|^p}{\alpha^p|y|_p^p} (\log^p n) \geq E_j\right]$
 $= 1 - e^{-\frac{|y_j|^p (\log^p n)}{\alpha^p|y|_p^p}} \leq \frac{|y_j|^p (\log^p n)}{\alpha^p|y|_p^p}$, so the expected number of large j is $O(\log^p n)$

Understanding the Large Elements

- Recall $(PDy)_i = \sum_j \delta(h(j) = i) \cdot \sigma_j \cdot (Dy)_j$, and set $R_j = \delta(h(j) = i) \cdot \sigma_j \cdot (Dy)_j$
- We have shown the expected number of large j is $O(\log^p n)$, so by a Markov bound we have $O(\log^p n)$ large j with constant probability and we condition on D satisfying this
- We also condition on $|Dy|_\infty \in \left[\frac{|y|_p}{10^{1/p}}, 10^{1/p} |y|_p \right]$, which held with probability $> 4/5$
- All the large j get perfectly hashed into separate hash buckets by P
 - We are throwing $O(\log^p n)$ balls into $s \geq n^{1-2/p}$ bins
- We apply Bernstein for each hash bucket separately
 - We apply Bernstein on the small indices j inside a hash bucket!

Understanding the Large Elements

- $E[(PDy)_i] = 0$ for each hash bucket i , and $E[(PDy)_i^2] = O\left(\frac{1}{s}\right) (n^{1-\frac{2}{p}} |y|_p^2)$
- Bernstein's bound: Suppose R_1, \dots, R_n are independent, and for all j , $|R_j| \leq K$, and $\text{Var}[\sum_j R_j] = \sigma^2$. There are constants C, c , so that for all $t > 0$,
 - $\Pr[|\sum_j R_j - E[\sum_j R_j]| > t] \leq C (e^{-\frac{ct^2}{\sigma^2}} + e^{-\frac{ct}{K}})$
- $(PDy)_i = \sum_j \delta(h(j) = i) \cdot \sigma_j \cdot (Dy)_j$, and $R_j = \delta(h(j) = i) \cdot \sigma_j \cdot (Dy)_j$
- Can assume $K = \max_j |R_j| \leq \frac{\alpha |y|_p}{\log n}$, since there is at most one large j in any hash bucket $(PDy)_i$
- Set $t = \frac{|y|_p}{100}$, and $s = \Theta(n^{1-\frac{2}{p}} \log n)$ in Bernstein's bound, to get for a bucket $(PDy)_i$:

$$\Pr \left[\left| \sum_{\text{small } j} \delta(h(j) = i) \sigma_j (Dy)_j \right| > \frac{|y|_p}{100} \right] \leq C \left(e^{-\Theta(\log n)} + e^{-c \frac{(\log n)}{100\alpha}} \right) \leq \frac{1}{n^2}$$
- By a union bound over all the s buckets, the "signed sum" of small j in every bucket will be at most $\frac{|y|_p}{100}$ 28

Wrapping Up

- For all i ,
 - $|(PDy)_i| \leq \frac{|y|_p}{100}$ if no large indices in i -th bucket
 - $|(PDy)_i| = |\sigma_j(Dy)_j| \pm \frac{|y|_p}{100}$ if exactly one large index j in i -th bucket
 - No bucket contains more than 1 large index j
- We conditioned on $|Dy|_\infty \in \left[\frac{|y|_p}{10^{\frac{1}{p}}}, 10^{\frac{1}{p}}|y|_p \right]$
- What is $|PDy|_\infty$?
 - $|PDy|_\infty \leq 10^{\frac{1}{p}}|y|_p + \frac{|y|_p}{100}$ and $|PDy|_\infty \geq \frac{|y|_p}{10^{\frac{1}{p}}} - \frac{|y|_p}{100}$
- So just output $|PDy|_\infty$ as your estimate to $|y|_p$
- Total space is $s = O(n^{1-\frac{2}{p}} \log n)$ words, which is $O(n^{1-\frac{2}{p}} \log^2 n)$ bits

Outline

- Quick recap of ℓ_1 -regression, and how to speed it up
- Introduction to the Streaming Model
- Estimating Norms in the Streaming Model
- **Heavy Hitters in a Stream**
- Estimating Number of Non-Zero Entries (ℓ_0)

Heavy Hitter Guarantees

- l_1 – guarantee
 - output a set containing all items j for which $|x_j| \geq \phi |x|_1$
 - the set should not contain any j with $|x_j| \leq (\phi - \epsilon) |x|_1$
- l_2 – guarantee
 - output a set containing all items j for which $x_j^2 \geq \phi |x|_2^2$
 - the set should not contain any j with $x_j^2 \leq (\phi - \epsilon) |x|_2^2$
- l_2 – guarantee can be much stronger than the l_1 – guarantee
 - Suppose $x = (\sqrt{n}, 1, 1, 1, \dots, 1)$
 - Item 1 is an l_2 -heavy hitter for constant ϕ, ϵ , but not an l_1 -heavy hitter
 - If $|x_j| \geq \phi |x|_1$, then $x_j^2 \geq \phi^2 |x|_1^2 \geq \phi^2 |x|_2^2$

Heavy Hitter Intuition

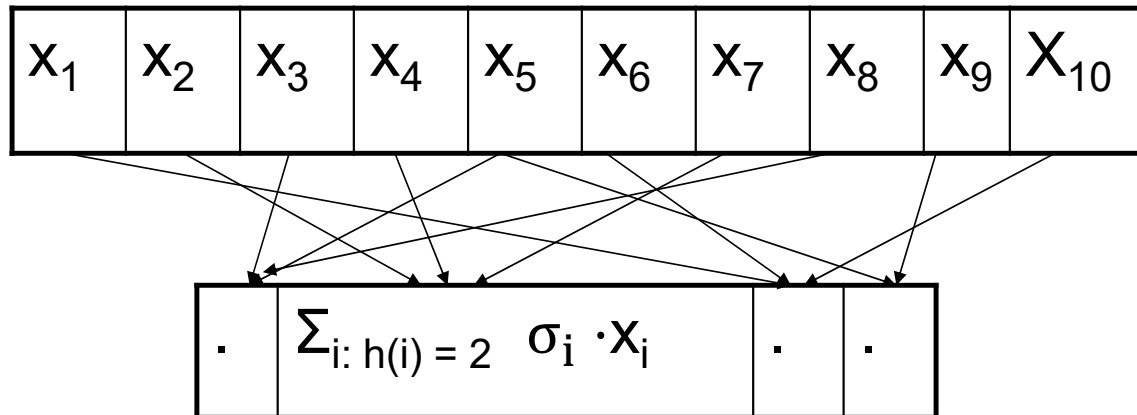
- Suppose you are promised at the end of the stream, $x_i = n$, and $x_j \in \{0,1\}$ for $j \in \{1, 2, \dots, n\}$ with $j \neq i$
- How would you find the identity i ?
- For each j in $\{1, 2, 3, \dots, \log n\}$, let $A_j \subset [n]$ be the set of indices with j -th bit in their binary representation equal to 0, and B_j be the set with j -th bit equal to 1
- Compute $a_j = \sum_{i \in A_j} x_i$ and $b_j = \sum_{i \in B_j} x_i$ for each j in $\{1, 2, \dots, \log n\}$
- Read off the identity of item i

Heavy Hitter Intuition Continued

- Suppose you are promised at the end of the stream, $x_i = 100\sqrt{n \log(n)}$, and $x_j \in \{0,1\}$ for $j \in \{1, 2, \dots, n\}$ with $j \neq i$
- How would you find the identity i ?
- For each j in $\{1, 2, 3, \dots, \log n\}$, let $A_j \subset [n]$ be the set of indices with j -th bit in their binary representation equal to 0, and B_j be the set with j -th bit equal to 1
- Compute $a_j = \sum_{i \in A_j} \sigma_i \cdot x_i$ and $b_j = \sum_{i \in B_j} \sigma_i \cdot x_i$ for each j in $\{1, 2, \dots, \log n\}$
- Read off the identity of item i ?
- Additive Chernoff bound implies magnitude of “noise” in a count is at most $\sqrt{n \log(n)}$ w.h.p.
- Remove assumptions: (1) $x_i = 100\sqrt{n \log(n)}$ and (2) and $x_j \in \{0,1\}$ for $j \in \{1, 2, \dots, n\}$ with $j \neq i$ ³

CountSketch achieves the l_2 -guarantee

- Assign each coordinate i a random sign $\sigma_i \in \{-1,1\}$
- Randomly partition coordinates into B buckets, maintain $c_j = \sum_{i: h(i)=j} x_i \cdot \sigma_i$ in the j -th bucket



- Estimate x_i as $\sigma_i \cdot c_{h(i)}$

Why Does CountSketch Work?

- $E[\sigma_i c_{h(i)}] = \sigma_i \sum_{i': h(i)=h(i')} \sigma_{i'} x_{i'} = x_i$
- Suppose we independently repeat this hashing scheme $O(\log n)$ times
- Output the median of the estimates across the $\log n$ repetitions
- “Noise” in a bucket is $\sigma_i \cdot \sum_{i' \neq i, h(i')=h(i)} \sigma_{i'} \cdot x_{i'}$
- What is the variance of the noise?
- $E\left[\left(\sigma_i \cdot \sum_{i' \neq i, h(i')=h(i)} \sigma_{i'} \cdot x_{i'}\right)^2\right] = \frac{|x|_2^2}{B}$
- So with constant probability, the noise in a bucket is $O\left(\frac{|x|_2}{\sqrt{B}}\right)$ in magnitude
- Since the $\log n$ repetitions are independent, this ensures that our estimate $\sigma_i c_{h(i)}$ will equal $x_i \pm O\left(\frac{|x|_2}{\sqrt{B}}\right)$ with probability $1 - 1/\text{poly}(n)$
- Hence, we approximate *every* x_i simultaneously up to additive error $O\left(\frac{|x|_2}{\sqrt{B}}\right)$

Tail Guarantee

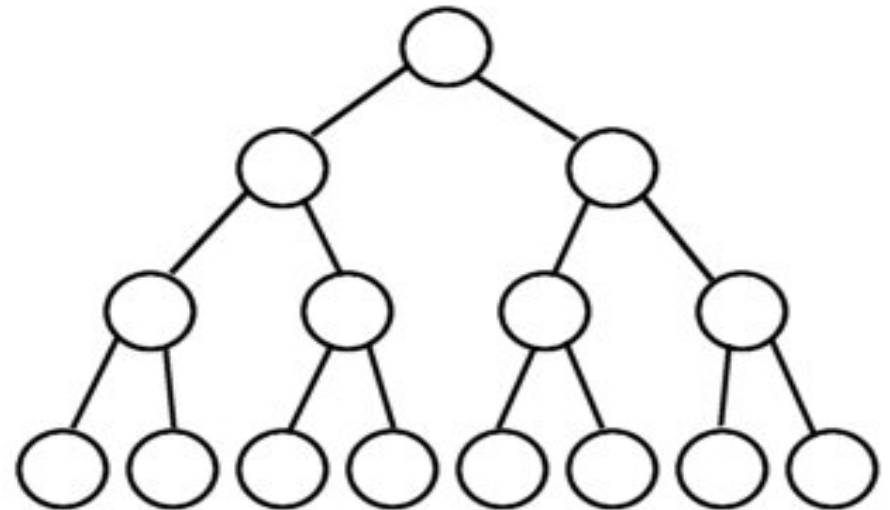
- CountSketch approximates every x_i simultaneously up to additive error $O(\frac{\|x\|_2}{\sqrt{B}})$
- But what if x_1 is a super large poly(n), and $x_2 = n$ and $x_3 = \dots = x_n = 1$?
- We get a pretty bad approximation to x_2
- **Tail Guarantee:** CountSketch approximates *every* x_i simultaneously up to additive error $O(\frac{\|x_{-B/4}\|_2}{\sqrt{B}})$, where $x_{-B/4}$ is x after zero-ing out its top $B/4$ coordinates in magnitude
- Proof: with probability at least $3/4$, in each repetition the top $B/4$ coordinates of x in magnitude do not land in the same hash bucket as x_i
 - Do we need a lot of independence for this?
- What happens if x is $B/4$ -sparse?

How to Find the Top k Heavy Hitters Quickly

- There are 2^i nodes in i -th level of tree
 - Start at the level with $2k$ nodes
- Each node corresponds to a subset of $[n]$ of size $n/2^i$ with the same i -bit prefix
- In i -th level, for each i , hash to $O(k)$ buckets repeat $O(\log k)$ times. Like CountSketch, but in each bucket we run an approximation algorithm to the 2-norm
- In top level our universe has only $2k$ nodes, so we find top k just by computing estimate for all of them

Main idea: in next level, we only need to consider the left and right child of each of the k nodes we found at the previous level. So only $2k \ll n$ nodes to consider.

Full Binary Tree



Why Care About the ℓ_1 -Guarantee?

- l_1 – guarantee
 - output a set containing all items j for which $|x_j| \geq \phi |x|_1$
 - the set should not contain any j with $|x_j| \leq (\phi - \epsilon) |x|_1$
- l_2 – guarantee
 - output a set containing all items j for which $x_j^2 \geq \phi |x|_2^2$
 - the set should not contain any j with $x_j^2 \leq (\phi - \epsilon) |x|_2^2$
- l_2 – guarantee implies the l_1 – guarantee
- So why care about the l_1 – guarantee?
- A nice thing about the l_1 -guarantee is that it can be solved deterministically!

Deterministic ℓ_1 Heavy Hitters

- An $s \times n$ matrix S is ϵ -incoherent if
 - for all columns S_i , $\|S_i\|_2 = 1$
 - for all pairs of columns S_i and S_j , $|\langle S_i, S_j \rangle| \leq \epsilon$
 - entries can be specified with $O(\log n)$ bits of space
- Compute $S \cdot x$ in a stream using $O(s \log n)$ bits of space
- Estimate $\hat{x}_i = S_i^T Sx$
 - $\hat{x}_i = \sum_{j=1, \dots, n} \langle S_i, S_j \rangle x_j = \|S_i\|_2^2 x_i \pm \max_{i,j} |\langle S_i, S_j \rangle| |x|_1 = x_i \pm \epsilon |x|_1$
 - Can figure out which $|x_i| \geq \phi |x|_1$ and which $|x_i| \leq (\phi - \epsilon) |x|_1$
- But do ϵ -incoherent matrices exist?

ϵ -Incoherent Matrices

- Consider a prime $q = \Theta((\log n)/\epsilon)$. Let $d = \epsilon \cdot q = O(\log n)$
- Consider n distinct non-zero polynomials p_1, \dots, p_n each of degree less than d .
 - $q^d - 1 > n$
- Associate p_i with i -th column of S
- Let $s = q^2$ and group the rows of S into q groups of size q
 - In j -th group, the i -th column has a single non-zero on the $p_i(j)$ -th entry
 - $p_i(j)$ -th entry is equal to $1/q^{1/2}$
- Each column S_i has $\|S_i\|_2 = 1$
- S_i and S_j each have the same non-zero in the k -th group iff $p_i(k) = p_j(k)$
- Number of such groups k is at most $d \leq \epsilon q$, so $|\langle S_i, S_j \rangle| \leq \epsilon$

Outline

- Quick recap of ℓ_1 -regression, and how to speed it up
- Introduction to the Streaming Model
- Estimating Norms in the Streaming Model
- Heavy Hitters in a Stream
- Estimating Number of Non-Zero Entries (ℓ_0)

Estimating the Number of Non-Zero Entries

- $|x|_0 = |\{i \text{ such that } x_i \neq 0\}|$
- How can we output a number Z with $(1 - \epsilon)Z \leq |x|_0 \leq (1 + \epsilon)Z$ with prob. 9/10?
 - Want $O((\log n)/\epsilon^2)$ bits of space
- Suppose $|x|_0 = O(\frac{1}{\epsilon^2})$. What can we do in this case?
- Use our algorithm for recovering a k -sparse vector from last time, $k = O(\frac{1}{\epsilon^2})$
 - What is another way?
- But what if $|x|_0 \gg \frac{1}{\epsilon^2}$?

Estimating the Number of Non-Zero Entries

- Suppose we somehow had an estimate Z with $Z \leq |x|_0 \leq 2Z$, what could we do?
- Independently sample each coordinate i with probability $p = 100/(Z \epsilon^2)$
- Let Y_i be an indicator random variable if coordinate i is sampled
- Let y be the vector restricted to coordinates i for which $Y_i = 1$
- $E[|y|_0] = \sum_{i \text{ such that } x_i \neq 0} E[Y_i] = p|x|_0 \geq \frac{100}{\epsilon^2}$
- $\text{Var}[|y|_0] = \sum_{i \text{ such that } x_i \neq 0} \text{Var}[Y_i] \leq \frac{200}{\epsilon^2}$
- $\Pr \left[\left| |y|_0 - E[|y|_0] \right| > \frac{100}{\epsilon} \right] \leq \frac{\text{Var}[|y|_0] \epsilon^2}{100^2} \leq \frac{1}{50}$
- Use sparse recovery or CountSketch to compute $|y|_0$ exactly
- Output $\frac{|y|_0}{p}$

But we don't
know Z ...

Estimating the Number of Non-Zero Entries

- Guess Z in powers of 2
- Since $0 \leq |x|_0 \leq n$, there are $O(\log n)$ guesses
- The i -th guess $Z = 2^i$ corresponds to sampling each coordinate with probability $p = \min(1, \frac{100}{2^i \epsilon^2})$
- Sample the coordinates as nested subsets $[n] = S_0 \supseteq S_1 \supseteq S_2 \supseteq \dots \supseteq S_{\log n}$
- Run previous algorithm for each guess
- One of our guesses Z satisfies $Z \leq |x|_0 \leq 2Z$ and we should use that guess
- *But how do we know which one?*

Estimating the Number of Non-Zero Entries

- Use the largest guess $Z = 2^i$ for which $\frac{400}{\epsilon^2} \leq |y|_0 \leq \frac{3200}{\epsilon^2}$
- If $\frac{800}{\epsilon^2} \leq E[|y|_0] \leq \frac{1600}{\epsilon^2}$, then $\frac{400}{\epsilon^2} \leq |y|_0 \leq \frac{3200}{\epsilon^2}$ with probability at least $49/50$
- If $\frac{100}{\epsilon^2} \leq E[|y|_0] \leq \frac{200}{\epsilon^2}$, then $|y|_0 < \frac{400}{\epsilon^2}$ with probability at least $49/50$
- So with probability $48/50$, we choose an i for which $\frac{200}{\epsilon^2} \leq E[|y|_0] \leq \frac{1600}{\epsilon^2}$
- There are only 4 such indices i , and all 4 of them satisfy $|y|_0 = (1 \pm \epsilon)E[|y|_0]$ simultaneously with probability $1-4/50$. So doesn't matter which i we choose
- Overall, our success probability is $1-2/50-4/50 > 4/5$

What is Our Overall Space Complexity?

- If we use our k -sparse recovery algorithm for $k = O\left(\frac{1}{\epsilon^2}\right)$, then it takes $O\left(\frac{\log n}{\epsilon^2}\right)$ bits of space in each of $\log n$ levels, so $O\left(\frac{\log^2 n}{\epsilon^2}\right)$ total bits of space ignoring random bits
 - How much randomness do we need?
 - Pairwise independence is enough for Chebyshev's inequality
 - Implement nested sampling by choosing a hash function $h: [n] \rightarrow [n]$, checking if first i bits of $h(j) = 0$
 - $O(\log n)$ bits of space for the randomness
- Can improve to $O\left(\frac{\log n (\log\left(\frac{1}{\epsilon}\right) + \log \log n)}{\epsilon^2}\right)$ bits. How?
- Just need to know number of non-zero counters, so reduce counters from $\log n$ bits to $O\left(\log\left(\frac{1}{\epsilon}\right) + \log \log n\right)$ to bits

Reducing Counter Size

- In sampling levels that we care about, we have $O\left(\frac{1}{\epsilon^2}\right)$ counters, each of $O(\log n)$ bits
- At most $O\left(\frac{\log n}{\epsilon^2}\right)$ prime numbers dividing any of these counters
- Choose a random prime $q = O\left(\frac{\log n \log \log n}{\epsilon^2}\right)$. Unlikely that q divides any counters
- Just maintain our sparse recovery structure mod q , so $O\left(\frac{(\log \log n + \log\left(\frac{1}{\epsilon}\right))}{\epsilon^2}\right)$ bits per each of $O(\log n)$ sparse recovery instances