

# Outline

- Wrapup of Lower Bounds
- Streaming Algorithms
  - Heavy Hitters
  - L0 estimation

# Aspects of 1-Way Communication of Index

- Alice has  $x \in \{0,1\}^n$
- Bob has  $i \in [n]$
- Alice sends a (randomized) message  $M$  to Bob
- $I(M ; X | R) = \sum_i I(M ; X_i | X_{<i}, R)$   
 $\geq \sum_i I(M ; X_i | R)$   
 $= n - \sum_i H(X_i | M, R)$
- Fano:  $H(X_i | M, R) \leq H(\delta)$  if Bob can guess  $X_i$  with probability  $> 1 - \delta$
- $CC_\delta(\text{Index}) \geq I(M ; X | R) \geq n(1-H(\delta))$

*The same lower bound applies if the protocol is only correct on average over  $x$  and  $i$  drawn independently from a uniform distribution*

# Distributional Communication Complexity



X

$f(X, Y)$ ?



Y

- $(X, Y) \sim \mu$
- $\mu$ -distributional complexity  $D_\mu(f)$ : the minimum communication cost of a protocol which outputs  $f(X, Y)$  with probability  $2/3$  for  $(X, Y) \sim \mu$ 
  - Yao's minimax principle:  $R(f) = \max_{\mu} D_\mu(f)$
- 1-way communication: Alice sends a single message  $M(X)$  to Bob

# Indexing is Universal for Product Distributions [Kremer, Nisan, Ron]

- Communication matrix  $A_f$  of a Boolean function  $f: X \times Y \rightarrow \{0,1\}$  has  $(x,y)$ -th entry equal to  $f(x,y)$
- $\max_{\text{product } \mu} D_\mu(f) = \Theta(\text{VC - dimension})$  of  $A_f$
- Implies a reduction from Index is optimal for product distributions

$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

# Indexing with Low Error

- Index Problem with  $1/3$  error probability and  $0$  error probability both have  $\Omega(n)$  communication
- Sometimes, want lower bounds in terms of error probability
- Indexing on Large Alphabets:
  - Alice has  $x \in \{0, 1\}^{n/\delta}$  with  $\text{wt}(x) = n$ , Bob has  $i \in [n/\delta]$
  - Bob wants to decide if  $x_i = 1$  with error probability  $\delta$
  - [Jayram, W] 1-way communication is  $\Omega(n \log(1/\delta))$
  - Can be used to get an  $\Omega(\log(\frac{1}{\delta}))$  bound for norm estimation
  - We've seen an  $\Omega(\log n + \epsilon^{-2} + \log(\frac{1}{\delta}))$  lower bound for norm estimation
  - There is an  $\Omega(\epsilon^{-2} \log \frac{1}{\delta} \log n)$  bit lower bound

# Beyond Product Distributions

*Although  $R(f) = \max_{\mu} D_{\mu}(f)$ , it may be that  $\max_{\mu} D_{\mu}(f) \gg \max_{\text{product } \mu} D_{\mu}(f)$ , so one often can't get good lower bounds by looking at product distributions...*

# Non-Product Distributions

- Needed for stronger lower bounds
- Example: approximate  $|x|_\infty$  up to a multiplicative factor of  $B$  in a stream
  - Lower bounds for  $p$ -norms

Gap $_\infty$  (x,y)  
Problem



$$x \in \{0, \dots, B\}^n$$



$$y \in \{0, \dots, B\}^n$$

- Promise:  $|x - y|_\infty \leq 1$  or  $|x - y|_\infty \geq B$
- Hard distribution non-product
- $\Omega(n/B^2)$  lower bound [Saks, Sun] [Bar-Yossef, Jayram, Kumar, Sivakumar]

# Outline

- Wrapup of Lower Bounds
- Streaming Algorithms
  - Heavy Hitters
  - L0 estimation



# 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  $\phi, \epsilon$ , 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

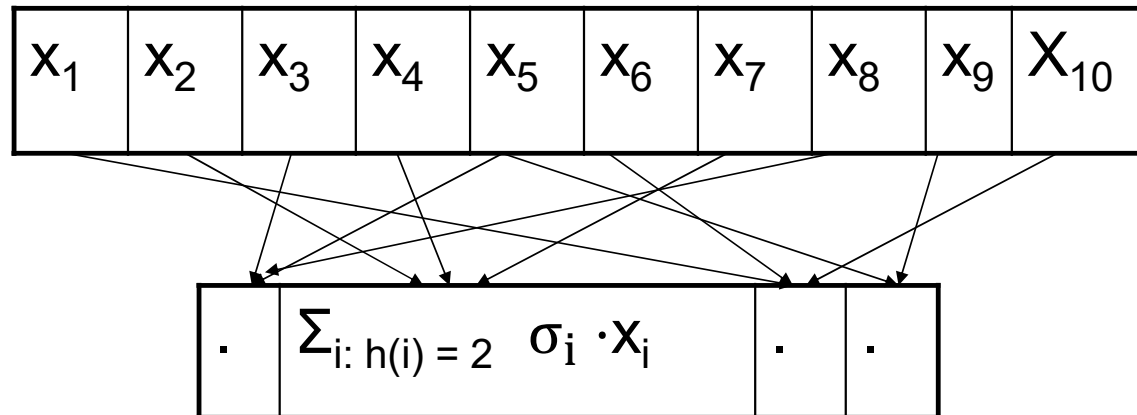
- Suppose you are promised at the end of the stream,  $x_i = n$ , and  $x_j \in \{0,1\}$  for  $j \in \{1, 2, \dots, n\}$  with  $j \neq i$
- How would you find the identity  $i$ ?
- For each  $j$  in  $\{1, 2, 3, \dots, \log n\}$ , let  $A_j \subset [n]$  be the set of indices with  $j$ -th bit in their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

# 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

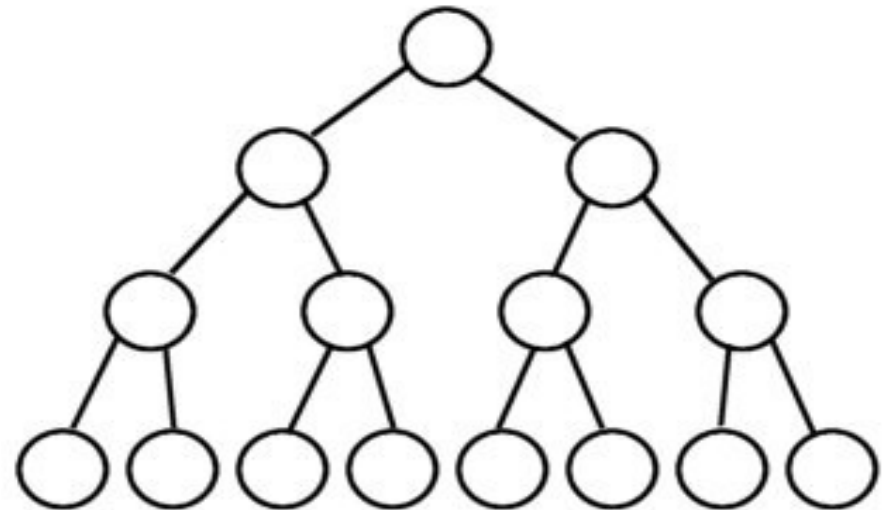
- CountSketch approximates every  $x_i$  simultaneously up to additive error  $O\left(\frac{\|x\|_2}{\sqrt{B}}\right)$
- But what if  $x_1$  is a super large poly(n), and  $x_2 = n$  and  $x_3 = \dots = x_n = 1$ ?
- We get a pretty bad approximation to  $x_2$
- **Tail Guarantee:** CountSketch approximates every  $x_i$  simultaneously up to additive error  $O\left(\frac{\|x_{-B/4}\|_2}{\sqrt{B}}\right)$ , 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 $\ell_1$ -Guarantee?

- $\ell_1$  – guarantee

- output a set containing all items  $j$  for which  $|x_j| \geq \phi |x|_1$
- the set should not contain any  $j$  with  $|x_j| \leq (\phi - \epsilon) |x|_1$

- $\ell_2$  – 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  $x_j^2 \leq (\phi - \epsilon) |x|_2^2$

- $\ell_2$  – guarantee implies the  $\ell_1$  – guarantee

- So why care about the  $\ell_1$  – guarantee?

- A nice thing about the  $\ell_1$ -guarantee is that it can be solved deterministically!



# Deterministic $\ell_1$ Heavy Hitters

- An  $s \times n$  matrix  $S$  is  $\epsilon$ -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  $\epsilon$ -incoherent matrices exist?

# $\epsilon$ -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

- Wrapup of Lower Bounds
- Streaming Algorithms
  - Heavy Hitters
  - **L0 estimation**

# 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  $O\left(\frac{\log\left(\log\left(\frac{1}{\epsilon}\right) + \log \log n\right)}{\epsilon^2}\right)$  bits. How?
- Just need to know number of non-zero counters, so reduce counters from  $\log n$  bits to  $O\left(\log\left(\frac{1}{\epsilon}\right) + \log \log n\right)$  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