# Sparse Recovery Application of Count-Sketch

Hu Fu @SHUFE, Oct 20, 2022

# The Sparse Recovery Problem

- Recall: In a streaming algorithm setting, we are sometimes interested in having a sparse vector that approximates the frequency vector  $\mathbf{x} \in \mathbb{R}^d$ 
  - A vector is sparse if it has few non-zero entries
  - We may measure the quality of approximation by  $\ell_2$  distance
- So given  $k \in \mathbb{N}$  and  $\epsilon \in (0, \frac{1}{2})$ , we are interested in finding  $\mathbf{y} \in \mathbb{R}^d$ , with
  - $\|\mathbf{y}\|_0 \le k$
  - $\|\mathbf{y} \mathbf{x}\|_2 \le (1 + \epsilon) E_2^k(\mathbf{x})$ , where  $E_2^k(\mathbf{x})$

Quantifying the error using  $E_2^k(\mathbf{x})$  is necessary. It can be as large as comparable to  $\|\mathbf{x}\|_2$ 

$$\mathbf{x}) := \min_{\mathbf{z} \in \mathbb{R}^d, \|z\|_0 \le k} \|\mathbf{z} - \mathbf{x}\|_2$$

## Sparse Recovery with Count-Sketch

- Offline optimum solution: pick the k largest entries of  $\mathbf{x}$ •
- **Recall Count-Sketch:** •
  - Draw  $\ell = O(\log d)$  hash functions  $h_1, \dots, h_\ell : [d] \to [w]$ , independently from a pairwise independent hash family
  - Draw  $\ell$  hash functions  $g_1, \dots, g_{\ell} : [d] \to \{-1, +1\}$ , independently from a pairwise independent hash family
  - At input  $i_t$ , increase counter  $C_i[h_i(i_t)]$  by  $g_i(i_t)\Delta_t$ , for  $j = 1, ..., \ell$
  - Output: for coordinate *i*, report  $\tilde{x}_i := \text{median} \{g_j(i)C_j[h_j(i)]\}$
- To solve sparse recovery, take  $w = 3k/\epsilon^2$ , take the k largest coordinates of  $\tilde{\mathbf{x}}$



- Main idea:
  - If we chose the k "correct" entries, since total error should be  $\epsilon E_2^k$ , the error If we chose the  $\kappa$  concert entry, and the controlled to  $\frac{\epsilon}{\sqrt{k}}E_2^k$
  - But the k entries we chose may differ from the "correct" ones. We should argue that, when all entries are estimated accurately enough, this doesn't introduce too much error.

### **Ideas** of **Proof**

- Main idea:

**Lemma.** Count-Sketch with  $w = 3k/\epsilon^2$ ,  $\ell = O(\log n)$  guarantees  $|x_j - \tilde{x}_j| \le \frac{\epsilon}{\sqrt{k}} E_2^k$  for

each *j* with high probability.

#### **Ideas** of **Proof**

# • If we chose the *k* "correct" entries, since total error should be $\epsilon E_2^k$ , the error allowed for each entry should be controlled to $\frac{\epsilon}{\sqrt{k}}E_2^k$

- Main idea:
  - If we chose the k "correct" entries, since total error should be  $\epsilon E_2^k$ , the error allowed for each entry should be controlled to  $\frac{\epsilon}{\sqrt{k}}E_2^k$
  - But the k entries we chose may differ from the "correct" ones. We should argue that, when all entries are estimated accurately enough, this doesn't introduce too much error.

**y**, then  $\|\mathbf{x} - \mathbf{z}\| \leq (1 + 5\epsilon)E_2^k(\mathbf{x})$ 

### **Ideas** of **Proof**



### Proof of First Lemma

**Lemma.** Count-Sketch with  $w = 3k/\epsilon^2$ ,  $\ell = O(\log n)$  guarantees  $|x_j - \tilde{x}_j| \le \frac{\epsilon}{\sqrt{k}} E_2^k$  for each j with high

probability.

Recall the analysis of Count-Sketch. For each  $j \in [d]$ ,  $i \in [\ell]$ , the *i*-th estimate is  $z_i := C_i[h_i(j)]g_i(j)$ , then  $\mathbb{E}[z_i] = x_i$ . To apply Chebyshev's inequality, we bound the variance of  $z_i$ . Let  $Y_{j,j'}$  be the indicator variable for the event  $h_i(j) = h_i(j')$ , then by pairwise independence of the hash family,  $\mathbb{P}[Y_{j,j'}] = \frac{1}{w}$ .  $\operatorname{Var}(z_i) = \mathbb{E}[(z_i - x_i)^2] = \mathbb{E}\left[\left(\sum_{j=1}^{n} g_i(j)g_i(j')Y_{j,j'}x_{j'}\right)^2\right] = \sum_{j=1}^{n} x_{j'}^2 \mathbb{E}[Y_{j,j'}^2] \le \frac{||x||^2}{w}$ j'≠j

$$\left\| \sum_{j' \neq j}^{2} \right\| = \sum_{j' \neq j} x_{j'}^{2} \mathbb{E}[Y_{j,j'}^{2}] \le \frac{\|x\|^{2}}{w}$$

### Refining the Analysis

- What is  $E_2^k(\mathbf{x})$ ?
  - Let T be the set of k entries of x with the largest absolute values, then  $E_2^k(\mathbf{x}) = \sqrt{\sum_{i' \notin T} x_{j'}^2}$
- The error introduced by collision with entries not in T is controllable by  $E_2^k$ What about collision with the entries in *T*?  $\bullet$
- - With w growing with k, this can be made to happen with small probability.

### Proof of First Lemma

**Lemma.** Count-Sketch with  $w = 3k/\epsilon^2$ ,  $\ell = O(\log n)$  guarantees  $|x_j - \tilde{x}_j| \le \frac{\epsilon}{\sqrt{k}} E_2^k$  for each j with high

probability.

Recall the analysis of Count-Sketch. For each  $j \in$ then  $\mathbb{E}[z_i] = x_j$ .

**Lemma.** 
$$\mathbb{P}[|z_i - x_j| \ge \frac{\epsilon}{\sqrt{k}} E_2^k(\mathbf{x})] \le \frac{2}{5}$$

Let  $z'_i = \sum_{\substack{j' \notin T, j' \neq j}} g_i(j)g_i(j')x_{j'}$ , then  $\operatorname{Var}[z'_i] \leq \frac{(E_2^k(\mathbf{x}))}{w}$ By Chebyshev inequality,  $\mathbb{P}[|z'_i - x_i| > \frac{\epsilon}{\sqrt{k}}E_2^k(\mathbf{x})]$ 

#### Recall the analysis of Count-Sketch. For each $j \in [d]$ , $i \in [\ell]$ , the *i*-th estimate is $z_i := C_i[h_i(j)]g_i(j)$ ,

$$\frac{(\mathbf{x})^2}{(\mathbf{x})^2} = \frac{3k}{\epsilon^2} (E_2^k(\mathbf{x}))^2.$$
$$\mathbf{x}^2 = \frac{1}{3}.$$

Let *A* be the event that none of entries in *T* collide with *j* under  $h_i$ , then  $\mathbb{P}[A] \ge 1 - \frac{e^2}{3}$ 



#### Proof of First Lemma

**Lemma.** 
$$\mathbb{P}[|z_i - x_j| \ge \frac{\epsilon}{\sqrt{k}} E_2^k(\mathbf{x})] \le \frac{2}{5}$$

Let 
$$z'_i = \sum_{j' \notin T, j' \neq j} g_i(j)g_i(j')x_{j'}$$
, then  $\operatorname{Var}[z'_i] \leq \frac{(E_2^k(\mathbf{x}))^2}{w} = \frac{3k}{\epsilon^2}(E_2^k(\mathbf{x}))^2$ .  
By Chebyshev inequality,  $\mathbb{P}[|z'_i - x_i| > \frac{\epsilon}{\sqrt{k}}E_2^k(\mathbf{x})] \leq \frac{1}{3}$ .

If  $|z_i - x_j| \ge \frac{\epsilon}{\sqrt{k}} E_2^k(\mathbf{x})$ , then either A happens or (or both happen) By union bound, the probability is  $|z'_i - x_j| \ge \frac{\epsilon}{\sqrt{t_i}} E_i$ 

Let A be the event that some entry in T collides with j under  $h_{i}$ , then  $\mathbb{P}[A] \leq \frac{\epsilon^2}{2}$ 

$$C_2^k(\mathbf{x}) \text{ at most } \frac{\epsilon^2}{3} + \frac{1}{3} \le \frac{2}{5}$$

Recall the proof we gave for the performance of SkipList. We had a similar use of union bound.

### **Proof of Second Lemma**

**Lemma.** If for 
$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$$
 we have  $\|\mathbf{x} - \mathbf{y}\|_{\infty} \le \frac{\epsilon}{\sqrt{k}}$   
 $\|\mathbf{x} - \mathbf{z}\|$ 

Let  $T \subseteq [d]$  be the set of k "big" entries of x, and T' be that for y, then  $||\mathbf{x} - \mathbf{z}||_2^2$  has three parts:

- entries in  $\overline{T} \cap \overline{T'}$ : by definition, these are original components of  $(E_2^k(\mathbf{x}))^2$
- entries in T T' and T' T: in  $(E_2^k(\mathbf{x}))^2$  we should have  $\sum x_j^2$

• note that |T - T'| = |T' - T| since |T| = |T'| =

• Key observation: entries in T - T' and T' - T must all be close (in absolute value)

 $-E_2^k(\mathbf{x})$ , let  $\mathbf{z}$  be the k-sparse recovery of  $\mathbf{y}$ , then

 $\leq (1+5\epsilon)E_2^k(\mathbf{x})$ 

 $\|\mathbf{x} - \mathbf{z}\|_{2}^{2} = \sum |x_{j} - z_{j}|^{2} + \sum |x_{j}^{2} + \sum |x_{j}^{2} + \sum |x_{j} - z_{j}|^{2}$  $j \in T \cap T' \qquad j \notin T \cup T' \qquad j \in T \setminus T' \qquad j \in T \setminus T'$ entries in  $T \cap T'$  and  $T' \setminus T$ : by assumption, each entry contributes  $\leq \frac{\epsilon^2}{k} (E_2^k(\mathbf{x}))^2$ , and there are k of them

$$= k.$$



**Lemma.** If for  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$  we have  $\|\mathbf{x} - \mathbf{y}\|_{\infty} \leq \frac{\epsilon}{\sqrt{k}} E_2^k(\mathbf{x})$ , let  $\mathbf{z}$  be the k-sparse recovery of  $\mathbf{y}$ , then

Let  $T \subseteq [d]$  be the set of k "big" entries of x, and T' be that for y, then  $||\mathbf{x} - \mathbf{z}||_2^2$  has three parts:  $j \in T \cap T'$ 

• entries in T - T' and T' - T: in  $(E_2^k(\mathbf{x}))^2$  we should have  $\sum x_j^2$ 

 $j \in T \setminus T'$ 

• note that |T - T'| = |T' - T| since |T| = |T'| = k.

• Key observation: entries in T - T' and T' - T must all be close (in absolute value)

**<u>Claim.</u>** If  $j \in T \setminus T'$  and  $j' \in T' \setminus T$ , then  $x_j \leq x_{j'} + \frac{2\pi}{T}$ 

#### $\|\mathbf{x} - \mathbf{z}\| \le (1 + 5\epsilon) E_2^k(\mathbf{x})$

 $\|\mathbf{x} - \mathbf{z}\|_{2}^{2} = \sum \|x_{j} - z_{j}\|^{2} + \sum |x_{j}^{2} + \sum |x_{j}^{2} + \sum |x_{j} - z_{j}|^{2}$  $j \notin T \cup T' \qquad j \in T \setminus T' \qquad j \in T' \setminus T$ 

$$\frac{\epsilon}{k} E_2^k(\mathbf{x})$$

$$\sum_{j \in T \setminus T'} x_j^2 \leq \sum_{j \in T' \setminus T} (|x_j| + \frac{2\epsilon}{\sqrt{k}} E_2^k(\mathbf{x}))^2 \leq \sum_{j \in T' \setminus T} x_j^2 + \frac{4\epsilon^2}{k} (E_2^k(\mathbf{x}))^2 + \frac{4\epsilon |x_j|}{\sqrt{k}} E_2^k(\mathbf{x}) \leq \sum_{j \in T' \setminus T} x_j^2 + 8\epsilon (E_2^k(\mathbf{x}))^2$$

$$By Cauchy-Schwartz, \sum_{j \in T' \setminus T} |x_j| \leq \sum_{j \notin T} |x_j| \leq \sqrt{k \sum_{j \notin T} x_j^2} = \sqrt{k} E_2^k(\mathbf{x})$$



### Proof of Second Lemma

**Lemma.** If for 
$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$$
 we have  $\|\mathbf{x} - \mathbf{y}\|_{\infty} \le \frac{\epsilon}{\sqrt{k}}$   
 $\|\mathbf{x} - \mathbf{z}\|$ 

 $\|\mathbf{x} - \mathbf{z}\|_2^2 = \sum |x_j - z_j|^2 + \sum |x_j^2 + \sum |x_j^2 + \sum |x_j - z_j|^2$ entries in T - T' and T' - T:  $\sum x_i^2 \le \sum x_j^2 + 8\epsilon (E_2^k(\mathbf{x}))^2$ 

Let  $T \subseteq [d]$  be the set of k "big" entries of x, and T' be that for y, then  $||\mathbf{x} - \mathbf{z}||_2^2$  has three parts:  $j \in \overline{T \cap T'} \qquad j \notin \overline{T \cup T'} \qquad j \in \overline{T \setminus T'} \qquad j \in \overline{T \setminus T}$ entries in  $T \cap T'$  and  $T' \setminus T$ : by assumption, each entry contributes  $\leq \frac{\epsilon^2}{k} (E_2^k(\mathbf{x}))^2$ , and there are k of them • entries in  $\overline{T} \cap \overline{T'}$ : by definition, these are original components of  $(E_2^k(\mathbf{x}))^2$  $j \in T \setminus T' \qquad j \in T' \setminus T$ Putting everything together,  $\|\mathbf{x} - \mathbf{z}\|_2^2 \le (1 + 9\epsilon)(E_2^k(\mathbf{x}))^2$ , hence  $\|\mathbf{x} - \mathbf{z}\| \le \sqrt{1 + 9\epsilon}E_2^k(\mathbf{x}) \le (1 + 5\epsilon)E_2^k(\mathbf{x})$ .

 $-E_2^k(\mathbf{x})$ , let  $\mathbf{z}$  be the k-sparse recovery of  $\mathbf{y}$ , then

 $\leq (1+5\epsilon)E_2^k(\mathbf{x})$ 



## Putting Things Together..

**Lemma.** Count-Sketch with  $w = 3k/\epsilon^2$ ,  $\ell = O(\log k)$ 

**Lemma.** If for 
$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$$
 we have  $\|\mathbf{x} - \mathbf{y}\|_{\infty} \le \frac{\epsilon}{\sqrt{k}}$   
 $\|\mathbf{x} - \mathbf{z}\|$ 

Putting the two Lemmas together, we have that with has error  $\leq (1 + 5\epsilon)E_2^k(\mathbf{x})$ 

One last thing: to give the sketch from  $ilde{{f x}}$ , naïvely we need to go through all the coordinates, which takes time O(d).

We can do faster by maintaining a record as the input comes!

n) guarantees 
$$|x_j - \tilde{x}_j| \le \frac{\epsilon}{\sqrt{k}} E_2^k$$
 for each  $j$  w.h.p.

 $-E_2^k(\mathbf{x})$ , let  $\mathbf{z}$  be the k-sparse recovery of  $\mathbf{y}$ , then

#### $|\leq (1+5\epsilon)E_2^k(\mathbf{x})$

Putting the two Lemmas together, we have that with high probability, the sparse recovery yielded by Count-Sketch

