

# Robust Sparsification via Sensitivity

Chansophea Wathanak In<sup>1</sup>

Yi Li<sup>1</sup>

David Woodruff<sup>2</sup>

Xuan Wu<sup>1</sup>



ICML 2025

## 1. Robust Sparsification Problem

**Def:** For an optimization problem

$$\min_{x \in \mathbb{R}^d} F(x) \text{ s.t.}$$

$$F(x) = F_1(x) + \dots + F_n(x)$$

**Robust Version:**

$$\min_{x \in \mathbb{R}^d} F^{(m)}(x) \text{ s.t.}$$

$F^{(m)}$  = sum of all but the  $m$  largest functions

**Motivation:** Capture Outliers

**Def:** ( $\epsilon, m$ )-coreset for  $F^{(m)}(x)$

$$\tilde{F}(x) = \sum_{j=1}^s w_j F_{i_j}(x) \text{ Sparse Representation}$$

$\tilde{F}$  is  $(\epsilon, m)$ -coreset of  $F$  if

$$\forall x \in \mathbb{R}^d, \forall t \leq m, \tilde{F}^{(t)}(x) \in (1 \pm \epsilon) \cdot F^{(t)}(x)$$

**Sensitivity**

$$\sigma(F_i) = \sup_{x \in \mathbb{R}^d} \frac{F_i(x)}{F(x)} \quad \text{Total Sensitivity}$$

$$T = \sum_{i=1}^n \sigma(F_i)$$

**Applications:** 1. Robust Clustering

2. Robust Subspace Embedding

3. Trimmed Squares Regression

4. Robust PCA

## 2. Our Contributions

### • Algorithms:

- For every  $F$  that has total sensitivity  $T$  and a vanilla coresset of size  $Q$ , we construct an  $(\epsilon, m)$ -coreset of size  $\tilde{O}(mT/\epsilon + Q)$
- As applications, we give  $\exp(O(m/\epsilon + m \log d)) + O(nd)$  time algorithms for robust regression and robust PCA with  $m$  outliers

### • Lower Bound:

- For subspace embedding, we prove  $mT/\epsilon$  is necessary ( $T = d$ )

### • Experiments:

- Our algorithm outperforms uniform sampling in real-world data

## 3. Prior Results

Most existing robust coresets are designed for k-clustering:

- $\exp(k+m)$  [FS12]
- $m + \text{poly}(mk/\epsilon)$  [HJLW23, HLLW25]
- $\min\{mk/\epsilon, m/\epsilon^2\} + Q$  [JL25]

Beyond clustering, existing results only satisfy some weaker coresset property:

- Bi-criteria [FL11, HJLW18]
- Local coreset [WGD21]

## 4. Main Novel Techniques

**Input:** A set  $A$  of functions, parameters  $\epsilon$  and  $m$

**Output:** A subset  $D \subseteq A$

- $B \leftarrow \emptyset$
- for each  $f \in A$ , with probability  $\frac{1}{m}$ , add  $f$  to  $B$
- for each  $f \in B$ , compute the sensitivity  $\sigma_B(f)$
- $D \leftarrow \{f \in B : \sigma_B(f) \geq \frac{\epsilon}{4}\}$
- Return**  $D$

### Technical Novelty:

- We characterize a subset of “contributing functions”, whose removal makes a vanilla coresset for remaining data also robust.
- A sampling based algorithm to find all contributing functions.
- An iterative reduction based analysis to control the number of contributing functions and the number of iterations.

## 5. Experiments

- Optimization problem:  
 $\ell_2$  subspace embedding

- Result: Outperform  
uniform sampling

- Dataset:  
1. Energy Dataset: 19735x28  
2. Emission Dataset: 36733x11

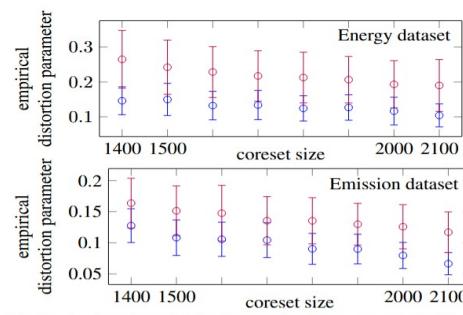


Figure 1: Results of subspace embedding coresets with  $m = 10$  and  $\epsilon = 0.25$ . Blue plots correspond to our coresset algorithm and the red plots to uniform sampling.