

# Average Sensitivity of Hierarchical $k$ -Median Clustering

Shijie Li<sup>1</sup>   Weiqiang He<sup>1</sup>   Ruobing Bai<sup>1</sup>   Pan Peng<sup>1</sup>

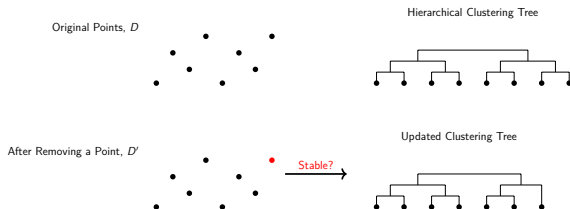
<sup>1</sup> School of Computer Science and Technology,  
University of Science and Technology of China

ICML 2025

# Motivation

- Hierarchical clustering reveals data structure at multiple scales
- But classic methods can be highly sensitive to small input changes [BLG14]<sup>1</sup>
- This instability harms interpretability and reliability

**Our goal:** Analyze and minimize the expected change in clustering output (symmetric difference) under perturbation



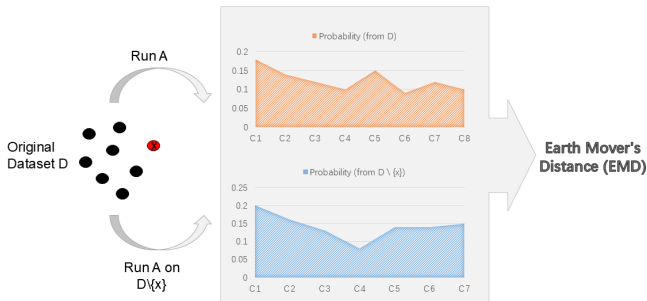
<sup>1</sup>Balcan, Liang, and Gupta. "Robust hierarchical clustering". In JMLR 2014

# Average Sensitivity of Randomized Algorithms

**Setup:** Randomized algorithm  $A$  and dataset  $P \subseteq [0, \Lambda]^d$

**Average sensitivity** [VY21]<sup>2</sup>:

$$\text{avg}_{p \in P} [d_{\text{EM}}(A(P), A(P \setminus \{p\}))]$$

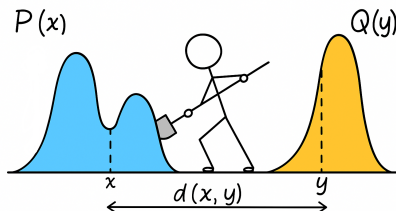


# Average Sensitivity of Randomized Algorithms

## Earth Mover's Distance (EMD):

$$d_{\text{EM}}(A(P), A(P \setminus \{p\})) = \min_D \mathbb{E}_{(S, S') \sim D} [|S \Delta S'|]$$

- $A(P), A(P \setminus \{p\})$ : distributions over clustering outputs
- $D$ : joint distribution with marginals matching each algorithm output



# Problem Formulation

## Hierarchical $k$ -median clustering

- Input: Dataset  $P \subseteq [0, \Lambda]^d$
- Output:
  - Centers:  $c_1, \dots, c_k$
  - Clusters:  $P_1, \dots, P_k$  minimizing the  $k$ -median cost for all  $k \in [n]$ , i.e.

$$\text{COST}(P, \{c_1, \dots, c_k\}) = \min_{P_1, \dots, P_k} \sum_{i \in [k]} \sum_{p \in P_i} \|p - c_i\|$$

- Perturbation model: Uniformly at random delete one point  $p \in P$

**Goal:** Design hierarchical  $k$ -median algorithms with **provably low average sensitivity**

# Our Contributions

## Theorem (Main Theorem, informal)

*Given a point set  $P$  of size  $n$  and a parameter  $\varepsilon > 0$ , our algorithm computes a hierarchical  $k$ -median clustering for all  $k \in \{1, \dots, n\}$  with:*

- *Expected cost:*

$$\mathbb{E}[\text{COST}_T(P, S_k)] \leq O(d \log \Lambda \cdot (1 + \varepsilon)^k) \cdot \text{OPT}(P, k)$$

- *Average sensitivity:*  $O\left(\frac{k \ln n}{\varepsilon}\right)$
- *Success probability:*  $\geq 1 - \frac{k}{n^2}$
- *Running time:*  $O(dn \log \Lambda + n^3)$

Here,  $S_k$  is the level- $k$  center set, and  $\text{COST}_T(P, S_k)$  is the clustering cost on RHST tree  $T$ .

# Our Contributions

We prove lower bounds on the average sensitivity of Single Linkage and deterministic CLNSS [CLN+21]<sup>3</sup>.

## Lemma (Single Linkage)

*The average sensitivity of Single Linkage is at least  $\Omega(n)$ .*

## Lemma (Deterministic CLNSS)

*The average sensitivity of the deterministic CLNSS algorithm is at least  $\Omega(n)$ .*

---

<sup>3</sup>Cohen-Addad et al. “Parallel and efficient hierarchical k-median clustering”. In NeurIPS 2021

# Our Algorithm

## Low-Sensitivity Hierarchical $k$ -Median Algorithm

**Input:** A set of points  $P$

**Output:** Centers  $c_1, \dots, c_n$ , clusterings  $\mathcal{P}_1, \dots, \mathcal{P}_n$

- ① Apply a random shift to each point in  $P$ .
- ② Construct a 2-RHST<sup>a</sup> tree  $T$ .
- ③ Initialize  $S_0 \leftarrow \emptyset$ ,  $\mathcal{P}_0 \leftarrow \{P\}$ .
- ④ Label all internal nodes of the RHST as unlabelled.
- ⑤ **For**  $t = 1$  to  $n$ , do the following:
  - ① Sample  $\lambda$  from a dataset-dependent interval.
  - ② Sample  $c_t$  with probability  $\propto \exp(-\text{COST}_T(P, x \cup S_{t-1})/\lambda)$ .
  - ③ Label the highest unlabelled ancestor of  $c_t$  with  $c_t$ .
  - ④ Update  $S_t \leftarrow c_t \cup S_{t-1}$ .
  - ⑤ Define  $\mathcal{P}_t$  by assigning points to closest labelled ancestor.

---

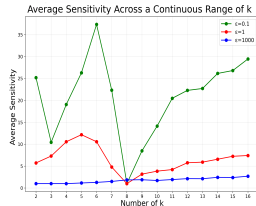
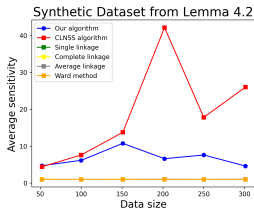
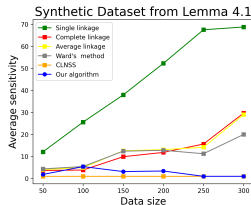
<sup>a</sup>Restricted 2-hierarchically well-separated tree (2-RHST)



# Experimental Setup

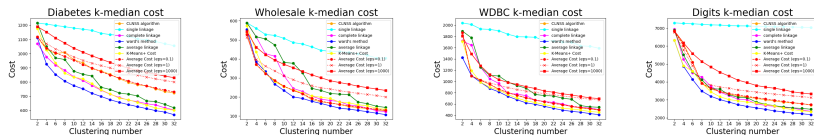
- **Datasets:** Synthetic and real-world (Scikit-learn , UCI Repository )
- **Baselines:**
  - Hierarchical methods: single, complete, average, Ward's
  - CLNSS algorithm
- **Metrics:**
  - Average sensitivity (robustness to deletions)
  - Clustering cost (e.g.,  $k$ -median)
  - Effect of  $\varepsilon$  (randomness impact)

# Experimental Results on Synthetic Datasets

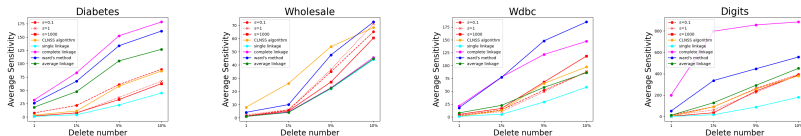


- Left and middle: average sensitivity of Single Linkage and CLN55 (vs others) on synthetic datasets to show the lower bounds.
- Right: results on a synthetic regression dataset with 500 points.

# Experimental Results on Real-World Datasets



- **$k$ -Median Cost:** Comparison across algorithms for varying  $k$  on real datasets.



- **Average Sensitivity ( $k = 4$ ):** Slightly worse than single linkage, but better than all other methods.

**Thank you!**