

Foundations of Comparison-Based Hierarchical Clustering

Debarghya Ghoshdastidar, Michaël Perrot, and Ulrike von Luxburg

Objective: Study hierarchical clustering when only similarity comparisons are available, that is without features nor explicit similarities.

Comparison-Based Machine Learning

Humans are bad at giving unbiased, quantitative information. Better at giving *relative information*.

Example: The **left vehicles** are *more similar* to each other than the **right vehicles**.



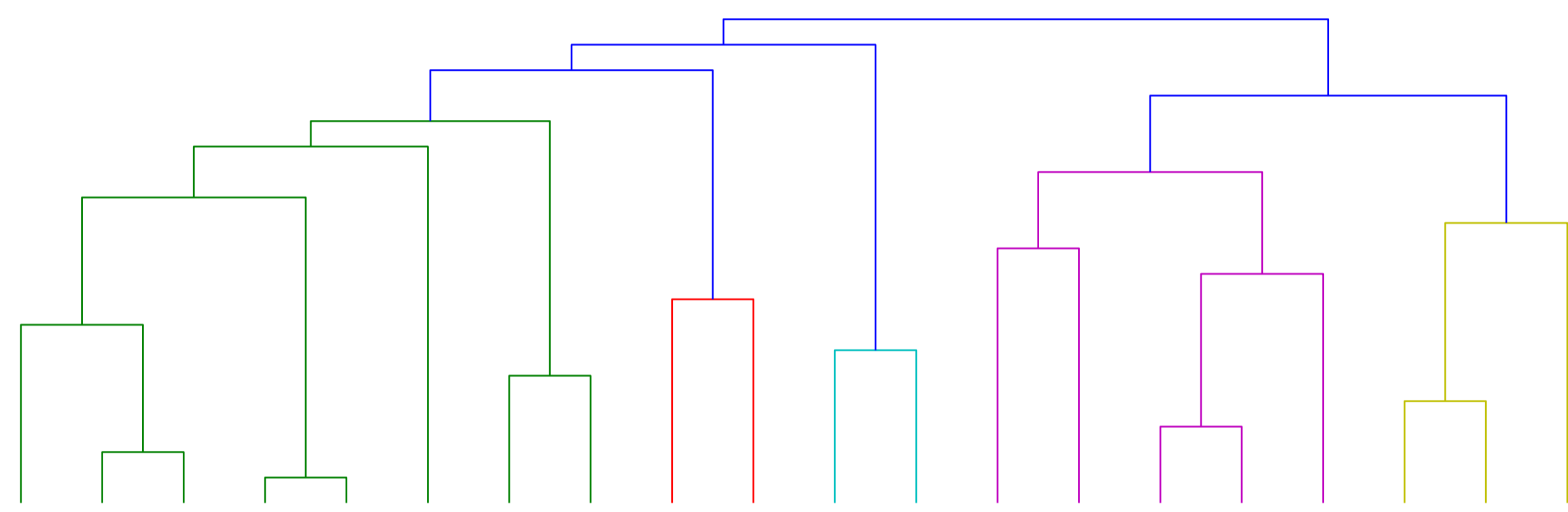
Given an unknown similarity function w , the corresponding quadruplet is

$$w_{\text{SUV left}; \text{SUV right}} \neq w_{\text{Sport car}; \text{Tractor}}$$

Challenging problem: No features (coordinates), not even distances!

Given a list of quadruplets, can we solve standard machine learning tasks such as *clustering*?

Example: Let $X = \{x_1, \dots, x_N\}$ be a set of N cars. Can we build a dendrogram that reflects their similarities using only a limited set of quadruplets \mathcal{Q} ?



Existing solutions:

- *Embedding based methods:* Retrieve a Euclidean representation of the objects that respects the quadruplets, then use standard machine learning methods.
- *Direct methods:* Algorithms that directly handle the quadruplets to solve a specific task.

Obtaining the comparisons:

- *Actively:* quadruplets chosen by the algorithm.
- *Passively:* quadruplets given to the algorithm with no way to make new queries.



```

input : Set of objects  $X = \{x_1, \dots, x_N\}$ ; Cluster-level similarity  $W : 2^X \times 2^X \rightarrow \mathbb{R}$ .
output: Binary tree, or dendrogram, representing a hierarchical clustering of  $X$ .
begin
  Let  $B$  be a collection of  $N$  singleton trees  $G_1, \dots, G_N$  with root nodes  $c_1, \dots, c_N$ .
  while  $|B| > 1$  do
    Let  $C = \{c, c'\}$  be the pair of trees in  $B$  for which  $W_{\text{root}(c), \text{root}(c')}$  is maximum.
    Create  $C^2$  with  $C^2.\text{root} = C.\text{root} \cup C'.\text{root}$ ,  $C^2.\text{left} = c$ , and  $C^2.\text{right} = c'$ .
    Add  $C^2$  to the collection  $B$ , and remove  $C, c'$ .
  end
  return The surviving element in  $B$ .
end
Algorithm 1: Agglomerative Hierarchical Clustering.

```



Contributions:

New algorithms for hierarchical clustering that *directly* use quadruplets. *Sufficient conditions* to guarantee exact recovery of a planted model.

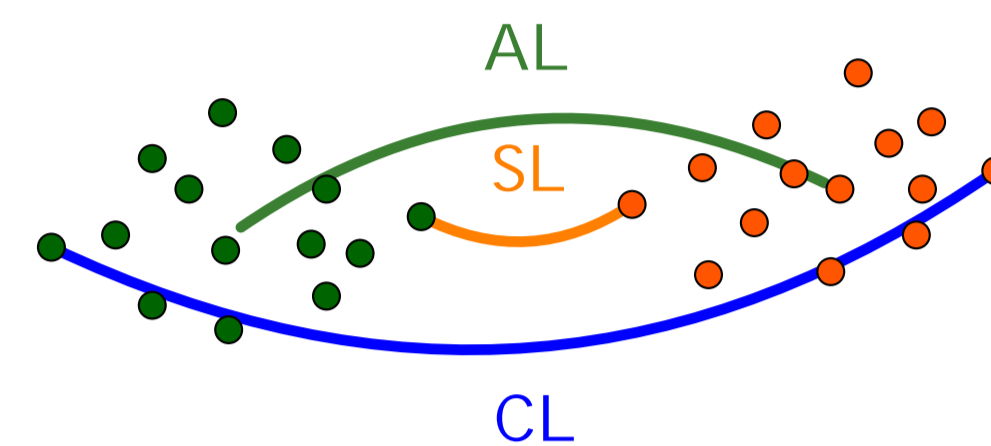
Setting

Hierarchical Clustering: Iteratively group clusters using a linkage function. Given G and G' :

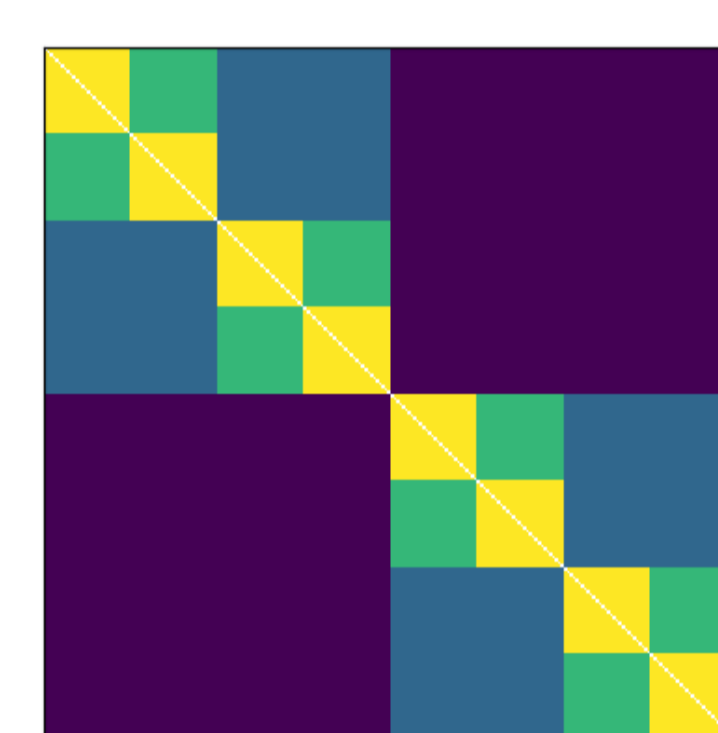
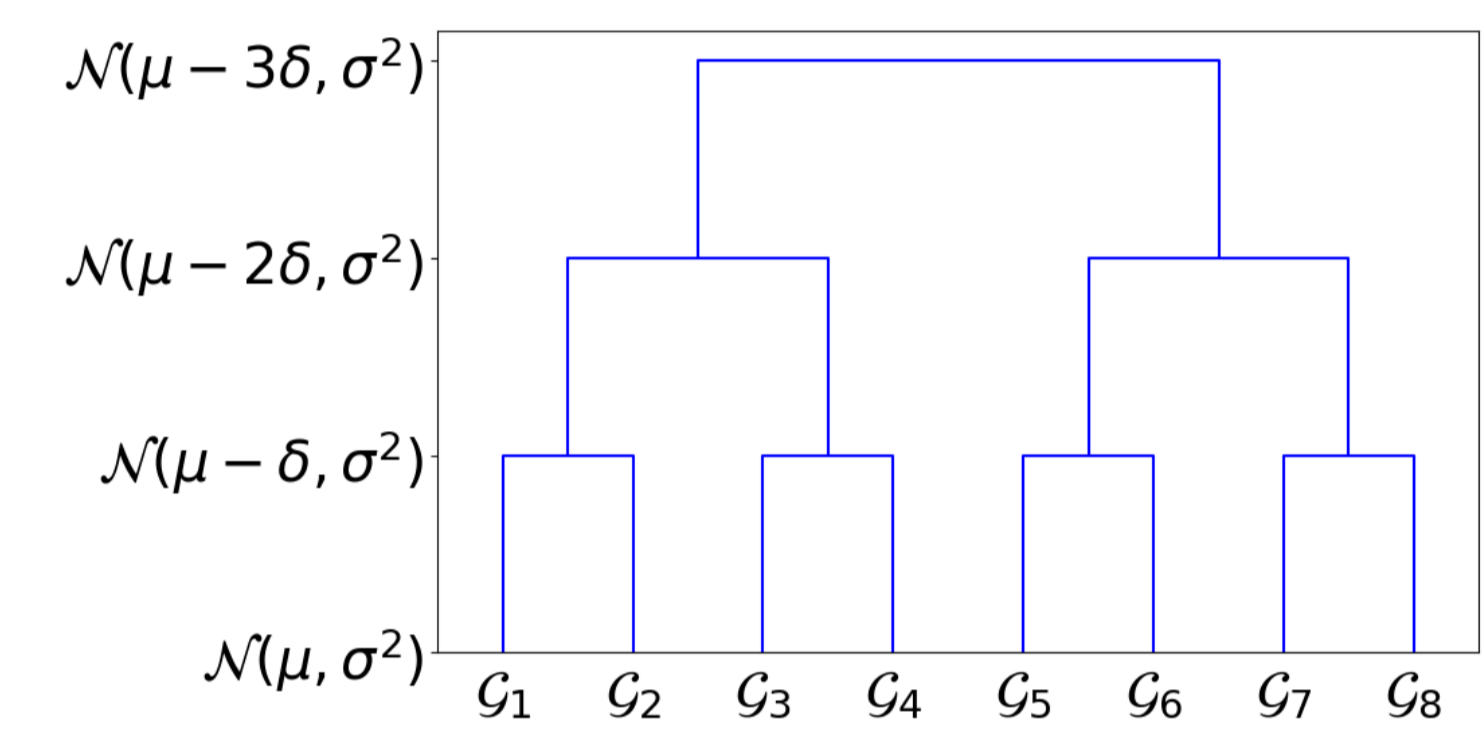
– Single Linkage (SL): $W_{\text{PG}; G'q} = \max_{x \in P, y \in G'} w_{ij}$

– Complete Linkage (CL): $W_{\text{PG}; G'q} = \min_{x \in P, y \in G'} w_{ij}$

– Average Linkage (AL): $W_{\text{PG}; G'q} = \frac{1}{|P||G'|} \sum_{x \in P, y \in G'} w_{ij}$



Planted Model: A noisy hierarchical block matrix with L levels, 2^L pure clusters of size N_0 and signal to noise δ .



Hierarchical structure.

Expected similarities.

Quadruplets Kernel Average Linkage (4K-AL)

Summary: Use the quadruplets to derive a proxy for the similarities between the examples.

Kernel function: Two similar objects should behave similarly with respect to any third object.

– **Active comparisons:** Let $w_{i_0 j_0}$ be a reference similarity and S be a set of landmarks:

$$K_{ij} = \frac{1}{|S|} \sum_{k \in S} \frac{w_{w_{ik} i_0 j_0} w_{w_{jk} i_0 j_0}}{w_{w_{ik} i_0 j_0} w_{w_{jk} i_0 j_0}}$$

– **Passive comparisons:** Use all the similarities as references and all the examples as landmarks:

$$K_{ij} = \frac{1}{N} \sum_{k: i, k} \frac{1}{N} \sum_{l: r, l} \frac{w_{p_i, r, k, l, q} w_{p_k, l, i, j, q}}{w_{p_i, r, k, l, q} w_{p_k, l, i, j, q}}$$

Quadruplets-Based Average Linkage (4-AL)

Summary: Use passive comparisons to define a cluster-level similarity function.

Cluster-level similarity: Clusters $G_1; G_2$ are more similar to each other than $G_3; G_4$ if their objects are, on average, more similar to each other than the objects of G_3 and G_4 :

$$W_{\text{PG}_1; G_2; G_3; G_4} = \frac{1}{|P|} \sum_{x \in P} \frac{1}{|G_3| |G_4|} \sum_{y \in G_3, z \in G_4} w_{xy}$$

Averaging over all cluster pairs gives rise to the following linkage function:

$$W_{\text{PG}; Gq} = \frac{1}{K} \sum_{r, s} \frac{W_{\text{PG}; G_r; G_s}}{K}$$

Theory

Summary: 4K-AL and 4-AL have better guarantees than SL and CL and use less quadruplets.

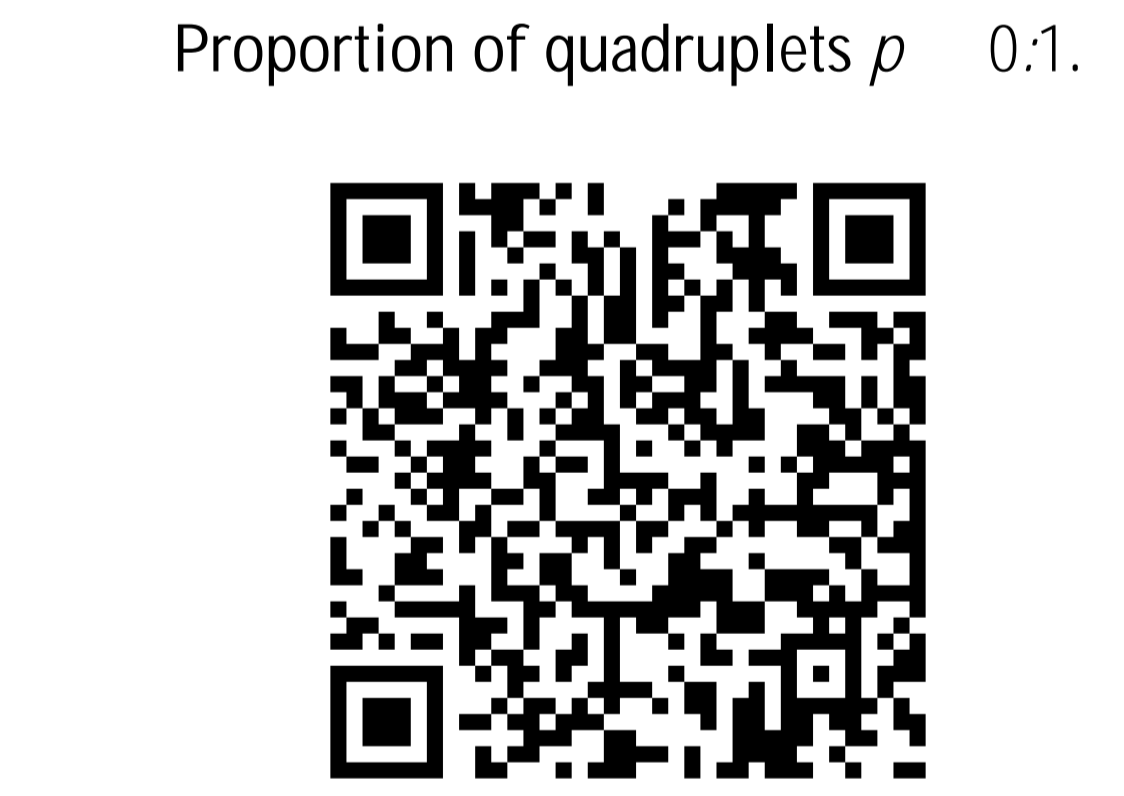
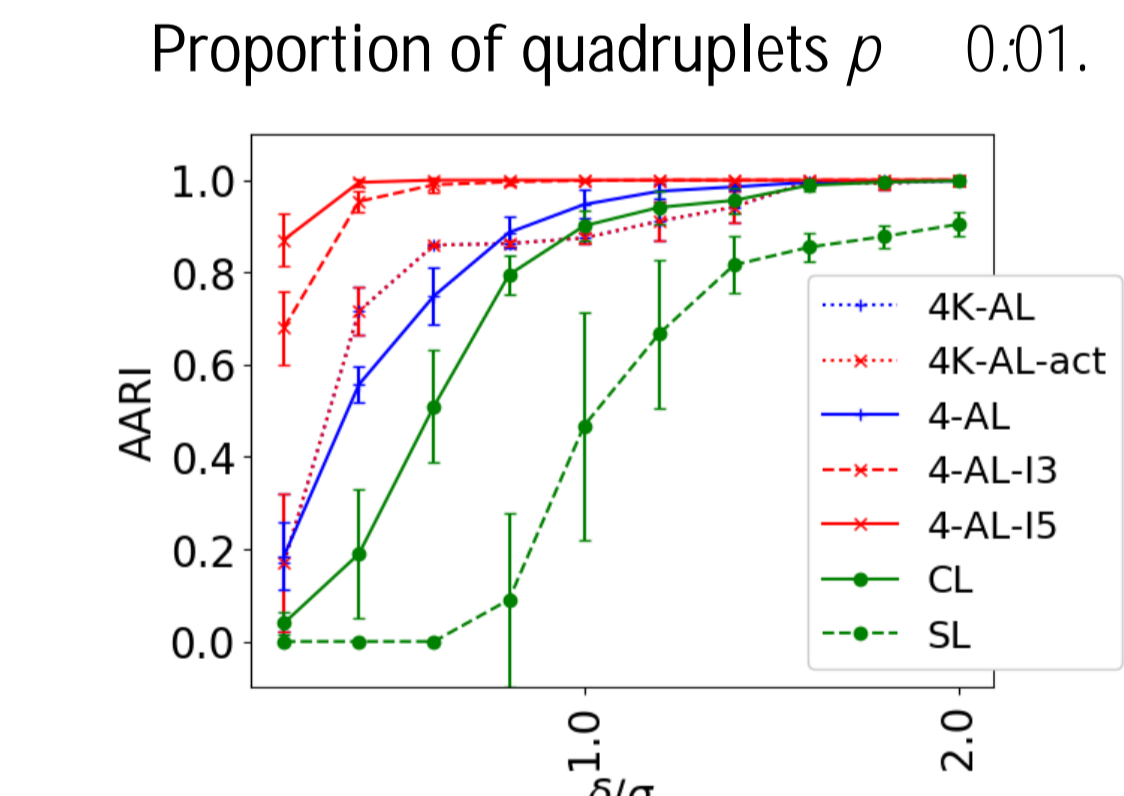
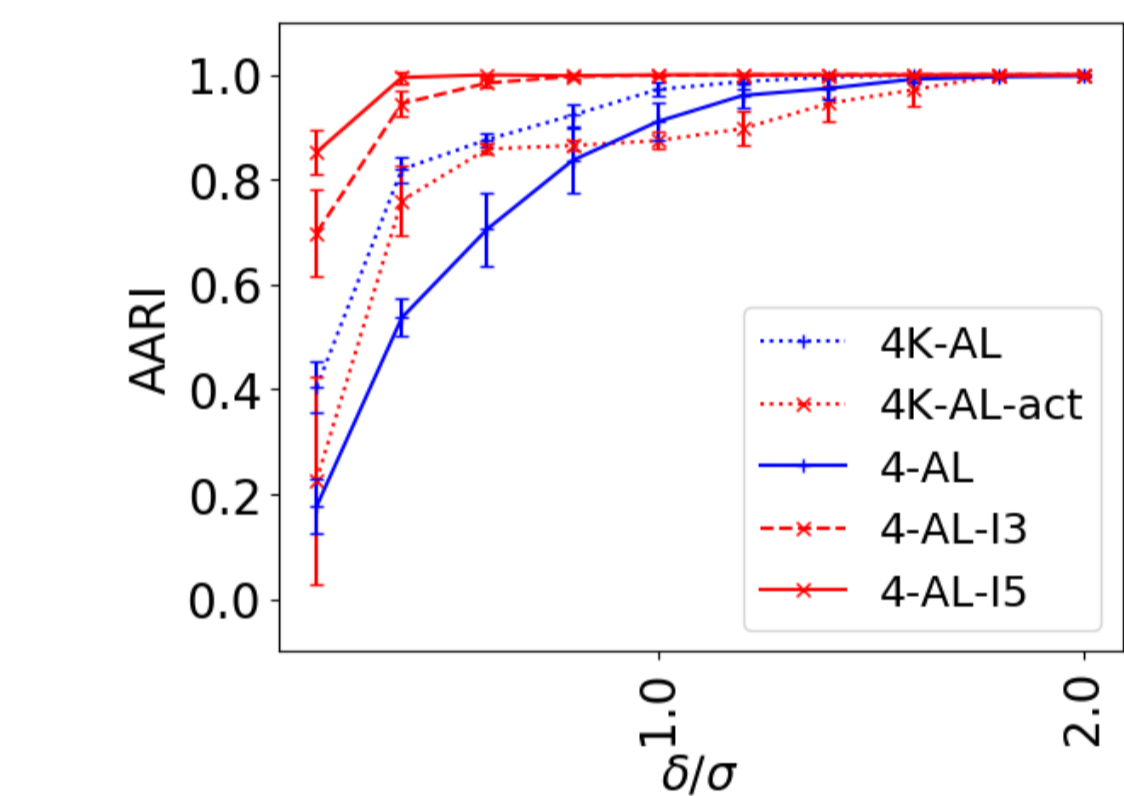
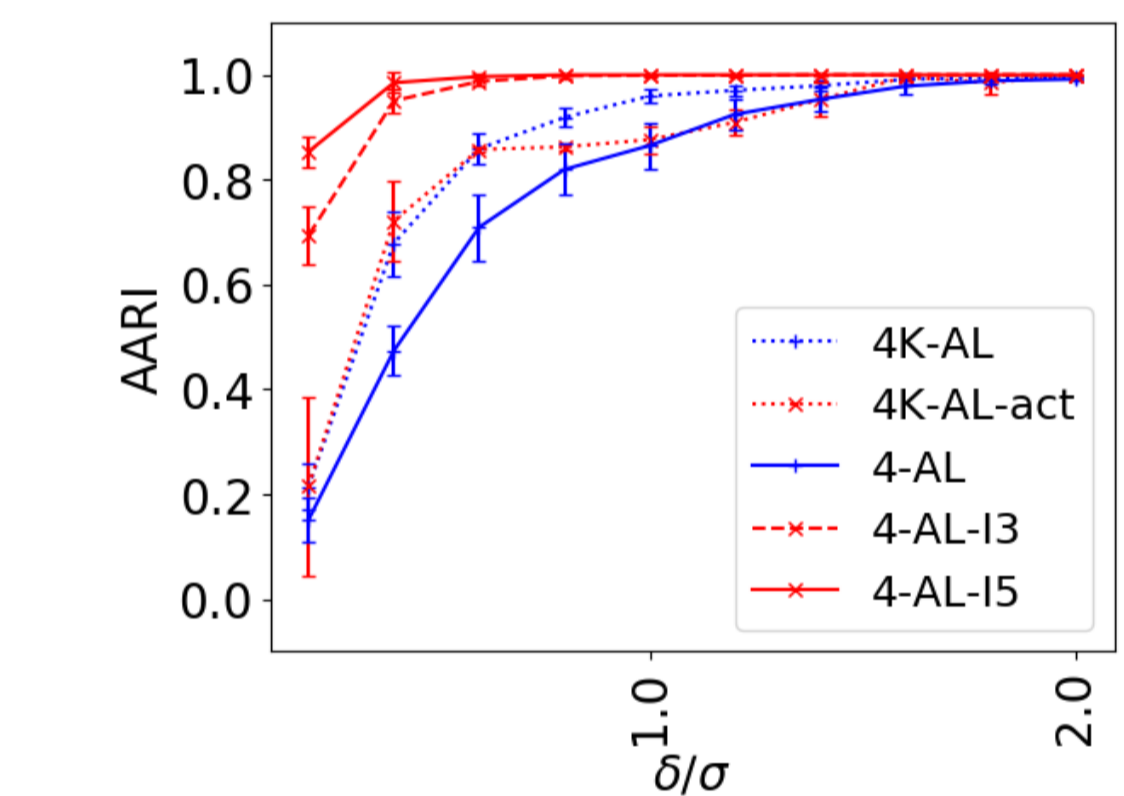
Recovery Guarantees ($L = \log N$)

| Method | Queries | Number of queries | Sufficient conditions | Remarks |
|--------|---------|--------------------------|-----------------------|---|
| SL | Active | N^2 | $\frac{1}{\ln N}$ | Tight! |
| CL | Active | N^2 | $\frac{1}{\ln N}$ | |
| 4K-AL | Active | $\mathcal{O}(N \ln N)$ | $\mathcal{O}(1)$ | Near-optimal number of queries. |
| 4K-AL | Passive | $\mathcal{O}(N^2 \ln N)$ | $\mathcal{O}(1)$ | |
| 4-AL | Passive | $N^3 \ln N$ | $\mathcal{O}(1)$ | Needs initial clusters of size $\geq \frac{1}{p} N_0$. |

Experiments

Planted Model: SL and CL only recover the hierarchy for large signal to noise ratios while 4K-AL and 4-AL exactly recover the hierarchy for smaller signal to noise ratios.

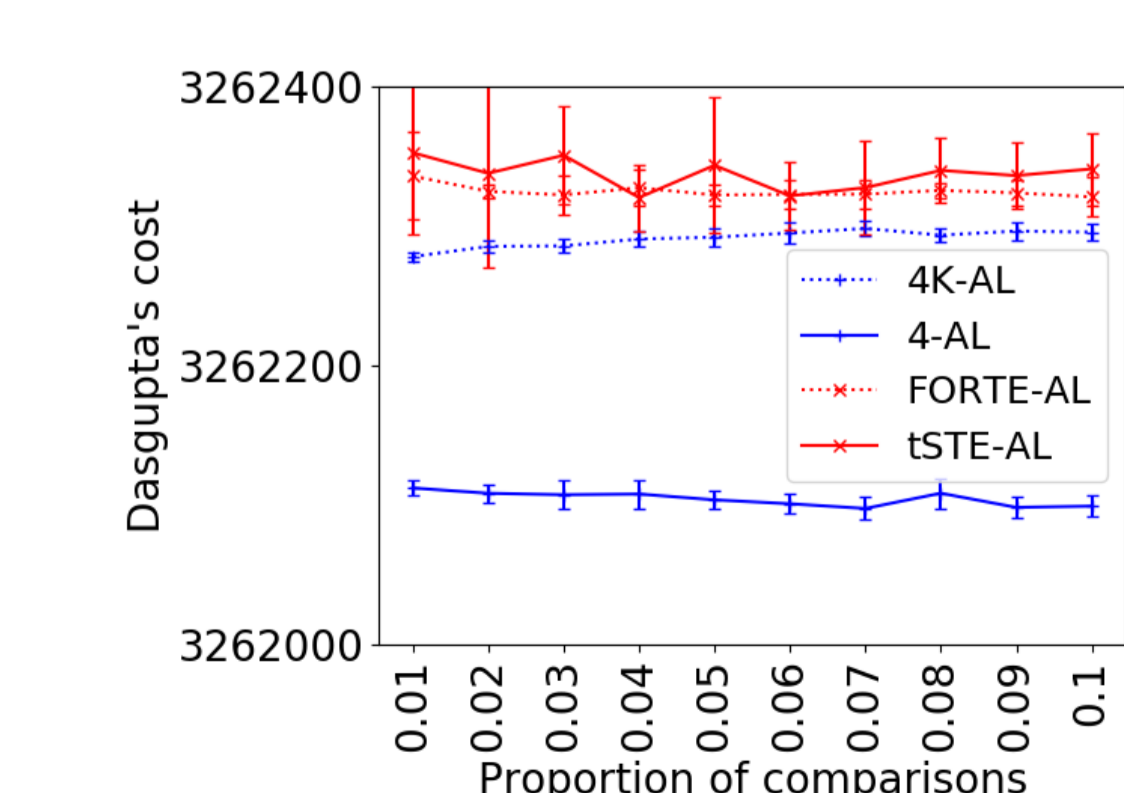
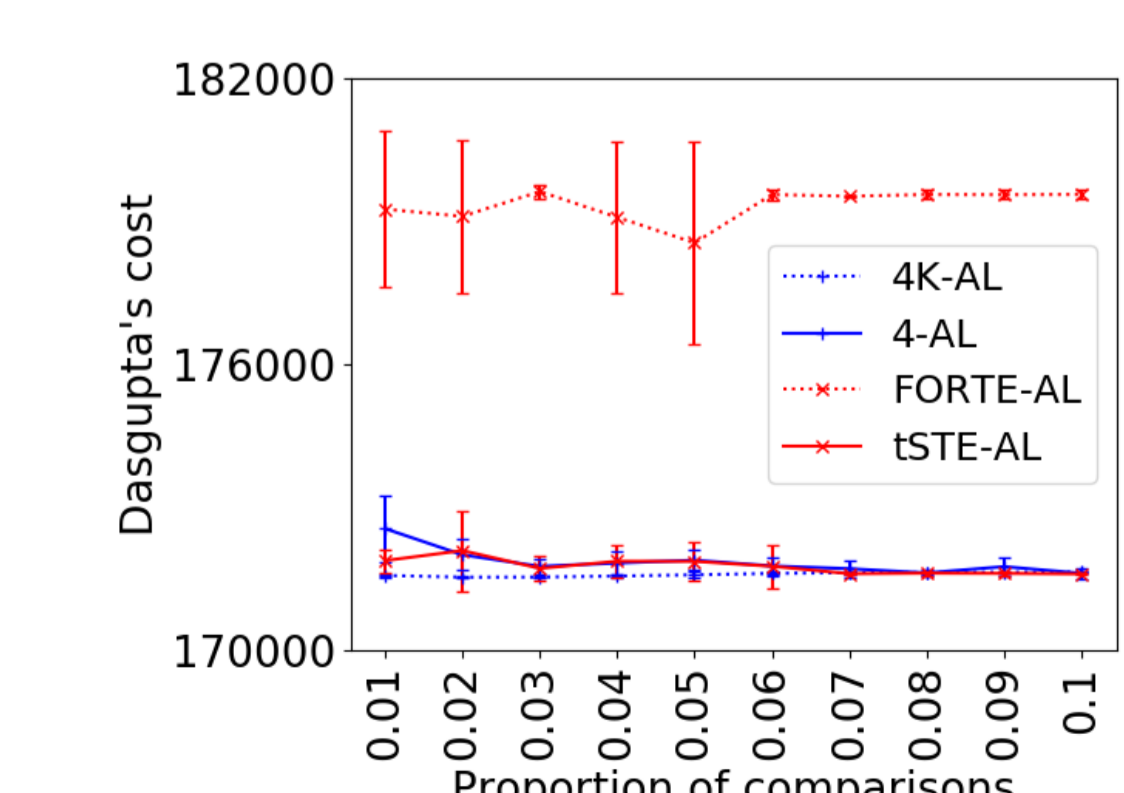
Evaluation: Average Adjusted Rand Index (AARI, higher is better).



Code available online!

Standard Datasets: 4K-AL and 4-AL are on average better than embedding based methods.

Evaluation: Dasgupta's cost (lower is better).



Zoo (100 examples, 16 features).

Glass (214 examples, 9 features).