
Active Clustering: Robust and Efficient Hierarchical Clustering using Adaptively Selected Similarities

Brian Eriksson
Boston University

Gautam Dasarathy
University of Wisconsin

Aarti Singh
Carnegie Mellon University

Robert Nowak
University of Wisconsin

Abstract

Hierarchical clustering based on pairwise similarities is a common tool used in a broad range of scientific applications. However, in many problems it may be expensive to obtain or compute similarities between the items to be clustered. This paper investigates the hierarchical clustering of N items based on a small subset of pairwise similarities, significantly less than the complete set of $N(N - 1)/2$ similarities. First, we show that if the intracluster similarities exceed intercluster similarities, then it is possible to correctly determine the hierarchical clustering from as few as $3N \log N$ similarities. We demonstrate this order of magnitude savings in the number of pairwise similarities necessitates sequentially selecting which similarities to obtain in an adaptive fashion, rather than picking them at random. We then propose an *active clustering* method that is robust to a limited fraction of anomalous similarities, and show how even in the presence of these noisy similarity values we can resolve the hierarchical clustering using only $O(N \log^2 N)$ pairwise similarities.

1 Introduction

Hierarchical clustering based on pairwise similarities arises routinely in a wide variety of engineering and scientific problems. These problems include inferring gene behavior from microarray data [1], Internet topology discovery [2], detecting community structure in social networks [3], advertising [4], and database management [5, 6]. It is often the case that there is a sig-

nificant cost associated with obtaining each similarity value. For example, in the case of Internet topology inference, the determination of similarity values requires many probe packets to be sent through the network, which can place a significant burden on the network resources. In other situations, the similarities may be the result of expensive experiments or require an expert human to perform the comparisons, again placing a significant cost on their collection.

The potential cost of obtaining similarities motivates a natural question: Is it possible to reliably cluster items using less than the complete, exhaustive set of all pairwise similarities? We will show that the answer is yes, particularly under the condition that intracluster similarity values are greater than intercluster similarity values, which we will define as the *Tight Clustering* (TC) condition. We also consider extensions of the proposed approach to more challenging situations in which a significant fraction of intracluster similarity values may be smaller than intercluster similarity values. This allows for robust, provably-correct clustering even when the TC condition does not hold uniformly.

The TC condition is satisfied in many situations. For example, the TC condition holds if the similarities are generated by a branching process (or tree structure) in which the similarity between items is a monotonic increasing function of the distance from the root to their nearest common branch point (ancestor). This sort of process arises naturally in clustering nodes in the Internet [7]. Also note that, for suitably chosen similarity metrics, the data can satisfy the TC condition even when the clusters have complex structures. For example, if similarities between two points are defined as the length of the longest edge on the shortest path between them on a nearest-neighbor graph, then they satisfy the TC condition given the clusters do not overlap. Additionally, density based similarity metrics [8] also allow for arbitrary cluster shapes while satisfying the TC condition.

One natural approach is to attempt clustering using a small subset of randomly chosen pairwise similarities. However, we show that this is quite ineffective in

general. We instead propose an *active* approach that sequentially selects similarities in an adaptive fashion, and thus we call the procedure *active clustering*. We show that under the TC condition, it is possible to reliably determine the unambiguous hierarchical clustering of N items using at most $3N \log N$ of the total of $N(N-1)/2$ possible pairwise similarities. Since it is clear that we must obtain at least one similarity for each of the N items, this is about as good as one could hope to do. Then, to broaden the applicability of the proposed theory and method, we propose a robust active clustering methodology for situations where a random subset of the pairwise similarities are unreliable and therefore fail to meet the TC condition. In this case, we show how using only $O(N \log^2 N)$ actively chosen pairwise similarities, we can still recover the underlying hierarchical clustering with high probability.

While there have been prior attempts at developing robust procedures for hierarchical clustering [9, 10], these works do not try to optimize the number of similarity values needed to robustly identify the true clustering, and mostly require all $O(N^2)$ similarities. Other prior work has attempted to develop efficient active clustering methods [11, 12, 13, 14], but the proposed techniques are ad-hoc and do not provide any theoretical guarantees. Balcan and Gupta developed a provably robust and efficient hierarchical clustering algorithm [15]. Their method can handle random or adversarial noise/errors and requires only a small random subset of similarities, but is only able to recover clusters of size $O(N)$. In contrast, the methods developed here are based on sequentially selecting queries, rather than picking them at random, enabling the recovery of significantly smaller clusters, as small as $O(\log N)$.

2 The Hierarchical Clustering Problem

Let $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$ be a collection of N items. Our goal will be to resolve a *hierarchical clustering* of these items.

Definition 1. A *cluster* C is defined as any subset of \mathbf{X} . A collection of clusters \mathcal{T} is called a **hierarchical clustering** if $\cup_{C_i \in \mathcal{T}} C_i = \mathbf{X}$ and for any $C_i, C_j \in \mathcal{T}$, only one of the following is true (i) $C_i \subset C_j$, (ii) $C_j \subset C_i$, (iii) $C_i \cap C_j = \emptyset$.

The hierarchical clustering \mathcal{T} has the form of a tree, where each node corresponds to a particular cluster. The tree is *binary* if for every $C_k \in \mathcal{T}$ that is not a leaf of the tree, there exists proper subsets C_i and C_j of C_k , such that $C_i \cap C_j = \emptyset$, and $C_i \cup C_j = C_k$. The binary tree is said to be *complete* if it has N leaf nodes, each corresponding to one of the individual items. Without loss

of generality, we will assume that \mathcal{T} is a complete (possibly unbalanced) binary tree, since any non-binary tree can be represented by an equivalent binary tree.

Let $\mathbf{S} = \{s_{i,j}\}$ denote the collection of all pairwise similarities between the items in \mathbf{X} , with $s_{i,j}$ denoting the similarity between x_i and x_j and assuming $s_{i,j} = s_{j,i}$. The traditional hierarchical clustering problem uses the complete set of pairwise similarities to infer \mathcal{T} . In order to guarantee that \mathcal{T} can be correctly identified from \mathbf{S} , the similarities must conform to the hierarchy of \mathcal{T} . We consider the following sufficient condition.

Definition 2. The triple $(\mathbf{X}, \mathcal{T}, \mathbf{S})$ satisfies the **Tight Clustering (TC) Condition** if for every set of three items $\{x_i, x_j, x_k\}$ such that $x_i, x_j \in C$ and $x_k \notin C$, for some $C \in \mathcal{T}$, the pairwise similarities satisfies, $s_{i,j} > \max(s_{i,k}, s_{j,k})$.

In words, the TC condition implies that the similarity between all pairs within a cluster is greater than the similarity with respect to any item outside the cluster. We can consider using off-the-shelf hierarchical clustering methodologies, such as bottom-up agglomerative clustering [16], on the set of pairwise similarities that satisfies the TC condition. Bottom-up agglomerative clustering is a recursive process that begins with singleton clusters (*i.e.*, the N individual items to be clustered). At each step of the algorithm, the pair of most similar clusters are merged. The process is repeated until all items are merged into a single cluster. It is easy to see that if the TC condition is satisfied, then the standard bottom-up agglomerative clustering algorithms such as single linkage, average linkage and complete linkage will all produce \mathcal{T} given the complete similarity matrix \mathbf{S} . Various agglomerative clustering algorithms differ in how the similarity between two clusters is defined, but every technique requires all $N(N-1)/2$ pairwise similarity values since all similarities must be compared at the very first step.

To properly cluster the items using fewer similarities requires a more sophisticated adaptive approach where similarities are carefully selected in a sequential manner. Before contemplating such approaches, we first demonstrate that adaptivity is necessary, and that simply picking similarities at random will not suffice.

Proposition 1. Let \mathcal{T} be a hierarchical clustering of N items and consider a cluster of size m in \mathcal{T} for some $m \ll N$. If n pairwise similarities, with $n < \frac{N}{m}(N-1)$, are selected uniformly at random from the pairwise similarity matrix \mathbf{S} , then any clustering procedure will fail to recover the cluster with high probability.

Proof. In order for any procedure to identify the m -sized cluster, we need to measure at least $m-1$ of the $\binom{m}{2}$ similarities between the cluster items. Let

$p = \binom{m}{2} / \binom{N}{2}$ be the probability that a randomly chosen similarity value will be between items inside the cluster. If we uniformly sample n similarities, then the expected number of similarities between items inside the cluster is approximately $n \binom{m}{2} / \binom{N}{2}$ (for $m \ll N$). Given Hoeffding's inequality, with high probability the number of observed pairwise similarities inside the cluster will be close to the expected value. It follows that we require $n \binom{m}{2} / \binom{N}{2} = n \frac{m(m-1)}{N(N-1)} \geq m-1$, and therefore we require $n \geq \frac{N}{m} (N-1)$ to reconstruct the cluster with high probability. \square

This result shows that if we want to reliably recover clusters of size $m = N^\alpha$ (where $\alpha \in [0, 1]$), then the number of randomly selected similarities must exceed $N^{(1-\alpha)}(N-1)$. In simple terms, randomly chosen similarities will not adequately sample all clusters. As the cluster size decreases (i.e., as $\alpha \rightarrow 0$) this means that almost all pairwise similarities are needed if chosen at random. This is more than are needed if the similarities are selected in a sequential and adaptive manner. In Section 3, we propose a sequential method that requires at most $3N \log N$ pairwise similarities to determine the correct hierarchical clustering.

3 Active Hierarchical Clustering under the TC Condition

From Proposition 1, it is clear that unless we acquire almost all of the pairwise similarities, reconstruction of the clustering hierarchy when sampling at random will fail with high probability. In this section, we demonstrate that under the assumption that the TC condition holds, an *active clustering* method based on *adaptively selected* similarities enables one to perform hierarchical clustering efficiently. Towards this end, we consider the work in [17] where the authors are concerned with a very different problem, namely, the identification of causality relationships among binary random variables. We present a modified adaptation of prior work here in the context of our hierarchical clustering from pairwise similarities problem.

From our discussion in the previous section, it is easy to see that the problem of reconstructing the hierarchical clustering \mathcal{T} of a given set of items $X = \{x_1, x_2, \dots, x_N\}$ can be reinterpreted as the problem of recovering a binary tree whose leaves are $\{x_1, x_2, \dots, x_N\}$. In [17], the authors define a special type of test on triples of leaves called the *leadership test* which identifies the “leader” of the triple in terms of the underlying tree structure. A leaf x_k is said to be the *leader* of the triple (x_i, x_j, x_k) if the path from the root of the tree to x_k does not contain the nearest common ancestor of x_i and x_j . This prior work shows

that one can efficiently reconstruct the entire tree \mathcal{T} using only these leadership tests.

The following lemma demonstrates that given observed pairwise similarities satisfying the TC condition, an **outlier** test using pairwise similarities will correctly resolve the leader of a triple of items.

Lemma 1. *Let \mathbf{X} be a collection of items equipped with pairwise similarities \mathbf{S} and hierarchical clustering \mathcal{T} . For any three items $\{x_i, x_j, x_k\}$ from \mathbf{X} , define*

$$\text{outlier}(x_i, x_j, x_k) = \begin{cases} x_i : \max(s_{i,j}, s_{i,k}) < s_{j,k} \\ x_j : \max(s_{i,j}, s_{j,k}) < s_{i,k} \\ x_k : \max(s_{i,k}, s_{j,k}) < s_{i,j} \end{cases} \quad (1)$$

If $(\mathbf{X}, \mathcal{T}, \mathbf{S})$ satisfies the TC condition, then $\text{outlier}(x_i, x_j, x_k)$ coincides with the leader of the same triple with respect to the tree structure conveyed by \mathcal{T} .

Proof. Suppose that x_k is the leader of the triple with respect to \mathcal{T} . This occurs if and only if there is a cluster $\mathcal{C} \in \mathcal{T}$ such that $x_i, x_j \in \mathcal{C}$ and $x_k \in \mathcal{T} \setminus \mathcal{C}$. By the TC condition, this implies that $s_{i,j} > \max(s_{i,k}, s_{j,k})$. Therefore x_k is the **outlier** of the same triple. \square

In Theorem 3.1, we find that by combining our **outlier** test approach with the tree reconstruction algorithm of [17], we discover an adaptive methodology (which we will refer to as **OUTLIERcluster**) that only requires on the order of $N \log N$ pairwise similarities to exactly reconstruct the hierarchical clustering \mathcal{T} .

Theorem 3.1. *Assume that the triple $(\mathbf{X}, \mathcal{T}, \mathbf{S})$ satisfies the Tight Clustering (TC) condition where \mathcal{T} is a complete (possibly unbalanced) binary tree that is unknown. Then, **OUTLIERcluster** recovers \mathcal{T} exactly using at most $3N \log_{3/2} N$ adaptively selected pairwise similarity values.*

Proof. From Appendix II of [17], we find a methodology that requires at most $N \log_{3/2} N$ leadership tests to exactly reconstruct the unique binary tree structure of N items. Lemma 1 shows that under the TC condition, each leadership test can be performed using only 3 adaptively selected pairwise similarities. Therefore, we can reconstruct the hierarchical clustering \mathcal{T} from a set of items \mathbf{X} using at most $3N \log_{3/2} N$ adaptively selected pairwise similarity values. \square

3.1 Tight Clustering Experiments

In Table 1 we see the results of both clustering techniques (**OUTLIERcluster** and bottom-up agglomerative clustering) on various synthetic tree topologies

given the Tight Clustering (TC) condition. The performance is in terms of the number of pairwise similarities required by the agglomerative clustering methodology, denoted by n_{agg} , and the number of similarities required by our **OUTLIERcluster** method, $n_{outlier}$. The methodologies are performed on both a balanced binary tree of varying size ($N = 128, 256, 512$) and a synthetic Internet tree topology generated using the technique from [18]. As seen in the table, our technique resolves the underlying tree structure using at most 11% of the pairwise similarities required by the bottom-up agglomerative clustering approach. As the number of items in the topology increases, further improvements are seen using **OUTLIERcluster**. Due to the pairwise similarities satisfying the TC condition, both methodologies resolve a binary representation of the underlying tree structure exactly.

Table 1: Comparison of **OUTLIERcluster** and Agglomerative Clustering on various topologies satisfying the Tight Clustering condition.

Topology	Size	n_{agg}	$n_{outlier}$	$\frac{n_{outlier}}{n_{agg}}$
Balanced	$N = 128$	8,128	876	10.78%
Binary	$N = 256$	32,640	2,206	6.21%
	$N = 512$	130,816	4,561	3.49%
Internet	$N = 768$	294,528	8,490	2.88%

While **OUTLIERcluster** determines the correct clustering hierarchy when all the pairwise similarities are consistent with the hierarchy \mathcal{T} , it can fail if one or more of the pairwise similarities are inconsistent. We find that with only two **outlier** tests erroneous at random, this corrupts the clustering reconstruction using **OUTLIERcluster** significantly. This can be attributed to the greedy construction of the clustering hierarchy using this methodology, where if one of the initial items is incorrectly placed in the hierarchy, this will result in a cascading effect that will drastically reduce the accuracy of the clustering.

4 Robust Active Clustering

Suppose that most, *but not all*, of the outlier tests agree with \mathcal{T} . This may occur if a subset of the similarities are in some sense inconsistent, erroneous or anomalous. We will assume that a certain subset of the similarities produce correct **outlier** tests and the rest may not. These similarities that produce correct tests are said to be *consistent* with the hierarchy \mathcal{T} . Our goal is to recover the clusters of \mathcal{T} despite the fact that the similarities are not always consistent with it.

Definition 3. The subset of consistent similarities is denoted $\mathbf{S}_C \subset \mathbf{S}$. These similarities satisfy the following property: if $s_{i,j}, s_{j,k}, s_{i,k} \in \mathbf{S}_C$ then

outlier(x_i, x_j, x_k) returns the leader of the triple (x_i, x_j, x_k) in \mathcal{T} (i.e., the outlier test is consistent with respect to \mathcal{T}).

We adopt the following probabilistic model for \mathbf{S}_C . Each similarity in \mathbf{S} fails to be consistent independently with probability at most $q < 1/2$ (i.e., membership in \mathbf{S}_C is termed by repeatedly tossing a biased coin). The expected cardinality of \mathbf{S}_C is $\mathbb{E}[|\mathbf{S}_C|] \geq (1 - q)|\mathbf{S}|$. Under this model, there is a large probability that one or more of outlier tests will yield an incorrect leader with respect to \mathcal{T} . Thus, our tree reconstruction algorithm in Section 3 will fail to recover the tree with large probability. We therefore pursue a different approach based on a top-down recursive clustering procedure that uses voting to overcome the effects of incorrect tests.

The key element of the top-down procedure is a robust algorithm for correctly splitting a given cluster in \mathcal{T} into its two subclusters, presented in Algorithm 1. Roughly speaking, the procedure quantifies how frequently two items tend to agree on **outlier** tests drawn from a small random subsample of other items. If they tend to agree frequently, then they are clustered together; otherwise they are not. We show that this algorithm can determine the correct split of the input cluster \mathcal{C} with high probability. The degree to which the split is “balanced” affects performance, and we need the following definition.

Definition 4. Let \mathcal{C} be any non-leaf cluster in \mathcal{T} and denote its subclusters by \mathcal{C}_L and \mathcal{C}_R ; i.e., $\mathcal{C}_L \cap \mathcal{C}_R = \emptyset$ and $\mathcal{C}_L \cup \mathcal{C}_R = \mathcal{C}$. The balance factor of \mathcal{C} is $\eta_{\mathcal{C}} := \min\{|\mathcal{C}_L|, |\mathcal{C}_R|\} / |\mathcal{C}|$.

Theorem 4.1. Let $0 < \delta' < 1$ and threshold $\gamma \in (0, 1/2)$. Consider a cluster $\mathcal{C} \in \mathcal{T}$ containing n items ($|\mathcal{C}| = n$) with balance factor $\eta_{\mathcal{C}} \geq \eta$ and disjoint subclusters \mathcal{C}_R and \mathcal{C}_L , and assume the following conditions hold:

- **A1** - The pairwise similarities are consistent with probability at least $1 - q$, for some $q \leq 1 - \frac{1}{\sqrt{2(1 - \delta')}}$.
- **A2** - q, η satisfy $(1 - (1 - q)^2) < \gamma < (1 - q)^2 \eta$.

If $m \geq c_0 \log(4n/\delta')$, for a constant $c_0 > 0$, and $n > 2m$, then with probability at least $1 - \delta'$ the output of **split**(\mathcal{C}, m, γ) is the correct subclusters, \mathcal{C}_R and \mathcal{C}_L .

The proof of the theorem is given in the Appendix. The theorem above shows that the algorithm is guaranteed (with high probability) to correctly split clusters that are sufficiently large for a certain range of q and η , as specified by **A2**. A bound on the constant c_0 is given in Equation 3 in the proof, but the important fact is that it does not depend on n , the number of items in \mathcal{C} . Thus all but the very smallest clusters can

Algorithm 1 : `split`(\mathcal{C}, m, γ)

Input :

1. A single cluster \mathcal{C} consisting of n items.
2. Parameters $m < n/2$ and $\gamma \in (0, 1/2)$

Initialize :

1. Select two subsets $\mathcal{S}_V, \mathcal{S}_A \subset \mathcal{C}$ uniformly at random (with replacement) containing m items each.
2. Select a “seed” item $x_j \in \mathcal{C}$ uniformly at random and let $\mathcal{C}_j \in \{\mathcal{C}_R, \mathcal{C}_L\}$ denote the subcluster it belongs to.

Split :

- For each $x_i \in \mathcal{C}$ and $x_k \in \mathcal{S}_A \setminus x_i$, compute the *outlier fraction* of \mathcal{S}_V :

$$c_{i,k} := \frac{1}{|\mathcal{S}_V \setminus \{x_i, x_k\}|} \sum_{x_\ell \in \mathcal{S}_V \setminus \{x_i, x_k\}} \mathbf{1}_{\{\text{outlier}(x_i, x_k, x_\ell) = x_\ell\}}$$

 where $\mathbf{1}$ denotes the indicator function.

- Compute the *outlier agreement* on \mathcal{S}_A :

$$a_{i,j} := \sum_{x_k \in \mathcal{S}_A \setminus \{x_i, x_j\}} (\mathbf{1}_{\{c_{i,k} > \gamma \text{ and } c_{j,k} > \gamma\}} + \mathbf{1}_{\{c_{i,k} < \gamma \text{ and } c_{j,k} < \gamma\}}) / |\mathcal{S}_A \setminus \{x_i, x_j\}|$$

- Assign item x_i to a subcluster according to

$$x_i \in \begin{cases} \mathcal{C}_j & : \text{ if } a_{i,j} \geq 1/2 \\ \mathcal{C}/\mathcal{C}_j & : \text{ if } a_{i,j} < 1/2 \end{cases}$$

Output : subclusters $\mathcal{C}_j, \mathcal{C}/\mathcal{C}_j$.

be reliably split. Note that total number of similarities required by `split` is at most $3mn$. So if we take $m = c_0 \log(4n/\delta')$, the total is at most $3c_0 n \log(4n/\delta')$. The key point of the lemma is this: *instead of using all $O(n^2)$ similarities, `split` only requires $O(n \log n)$.*

The allowable range in **A2** is non-degenerate and covers an interesting regime of problems in which q is not too large and η is not too small. The allowable range of γ cannot be determined without knowledge of η and q , so in practice $\gamma \in (0, 1/2)$ is a user-selected parameter (we use $\gamma = 0.30$ in all our experiments in the following section), and the Theorem holds for the corresponding set of (q, η) in **A2**.

We now give our robust active hierarchical clustering algorithm, `RAcluster`. Given an initial single cluster of N items, the `split` methodology of Algorithm 1 is recursively performed until all subclusters are of size

less than or equal to $2m$, the minimum resolvable cluster size where we can overcome inconsistent similarities through voting. The output is a hierarchical clustering \mathcal{T}' . The algorithm is summarized in Algorithm 2.

Algorithm 2 : `RAcluster`(\mathcal{C}, m, γ)

Given :

1. \mathcal{C} , n items to be hierarchically clustered.
2. parameters $m < n/2$ and $\gamma \in (0, 1/2)$

Partitioning :

1. Find $\{\mathcal{C}_L, \mathcal{C}_R\} = \text{split}(\mathcal{C}, m, \gamma)$.
2. Evaluate hierarchical subtrees, $\mathcal{T}_L, \mathcal{T}_R$, of cluster \mathcal{C} using:

$$\mathcal{T}_L = \begin{cases} \text{RAcluster}(\mathcal{C}_L, m, \gamma) & : \text{ if } |\mathcal{C}_L| > 2m \\ \mathcal{C}_L & : \text{ otherwise} \end{cases}$$

$$\mathcal{T}_R = \begin{cases} \text{RAcluster}(\mathcal{C}_R, m, \gamma) & : \text{ if } |\mathcal{C}_R| > 2m \\ \mathcal{C}_R & : \text{ otherwise} \end{cases}$$

Output : Hierarchical clustering $\mathcal{T}' = \{\mathcal{T}_L, \mathcal{T}_R\}$ containing subclusters of size $> 2m$.

Theorem 4.1 shows that it suffices to use $O(n \log n)$ similarities for each call of `split`, where n is the size of the cluster in each call. Now if the splits are balanced, the depth of the complete cluster tree will be $O(\log N)$, with $O(2^\ell)$ calls to `split` at level ℓ involving clusters of size $n = O(N/2^\ell)$. An easy calculation then shows that the total number of similarities required by `RAcluster` is then $O(N \log^2 N)$, compared to the total number which is $O(N^2)$. The performance guarantee for the robust active clustering algorithm are summarized in the following main theorem.

Theorem 4.2. *Let \mathbf{X} be a collection of N items with underlying hierarchical clustering structure \mathcal{T} . Let $0 < \delta < 1$ and $\delta' = \frac{\delta}{2N^{1/\log(\frac{1}{1-\eta})}}$. If $m \geq k_0 \log(\frac{8}{\delta}N)$, for a constant $k_0 > 0$ and conditions **A1** and **A2** of Theorem 4.1 hold, then with probability at least $1 - \delta$ `RAcluster`(\mathbf{X}, m, γ) uses $O(N \log^2 N)$ similarities and recovers all clusters $\mathcal{C} \in \mathcal{T}$ that satisfy the following conditions:*

- The cluster size, $|\mathcal{C}| > 2m$
- All clusters in \mathcal{T} that contain \mathcal{C} have a balance factor $\geq \eta$

The proof of the theorem is given in the Appendix. The constant k_0 is specified in Equation 4. Roughly speaking, the theorem implies that under the conditions of the Theorem 4.1 we can robustly recover all clusters of size $O(\log N)$ or larger using only

$O(N \log^2 N)$ similarities. Comparing this result to Theorem 3.1, we note three costs associated with being robust to inconsistent similarities: 1) we require $O(N \log^2 N)$ rather than $O(N \log N)$ similarity values; 2) the degree to which the clusters are balanced now plays a role (in the constant η); 3) we cannot guarantee the recovery of clusters smaller than $O(\log N)$ due to voting.

5 Robust Clustering Experiments

To test our robust clustering methodology we focus on experimental results from a balanced binary tree using synthesized similarities and real-world data sets using genetic microarray data ([19], with 7 expressions per gene and using Pearson correlation), breast tumor characteristics (via [20], with 10 features per tumor and using Pearson correlation), and phylogenetic sequences ([21], using the Needleman-Wunsch algorithm [22] between amino acid sequences). The synthetic binary tree experiments allows us to observe the characteristics of our algorithm while controlling the amount of inconsistency with respect to the Tight Clustering (TC) condition, while the real world data gives us perspective on problems where the tree structure and TC condition is assumed, but not known.

In order to quantify the performance of the tree reconstruction algorithms, consider the non-unique partial ordering, $\pi : \{1, 2, \dots, N\} \rightarrow \{1, 2, \dots, N\}$, resulting from the ordering of items in the reconstructed tree. For a set of observed similarities, given the original ordering of the items from the true tree structure we would expect to find the largest similarity values clustered around the diagonal of the similarity matrix. Meanwhile, a random ordering of the items would have the large similarity values potentially scattered away from the diagonal. To assess performance of our reconstructed tree structures, we will consider the rate of decay for similarity values off the diagonal of the reordered items, $\hat{s}_d = \frac{1}{N-d} \sum_{i=1}^{N-d} s_{\pi(i), \pi(i+d)}$. Using \hat{s}_d , we define a distribution over the the average off-diagonal similarity values, and compute the entropy of this distribution as follows (where, $\hat{p}_{\pi_i} = \left(\sum_{d=1}^{N-1} \hat{s}_d\right)^{-1} \hat{s}_i$):

$$\hat{E}(\pi) = - \sum_{i=1}^{N-1} \hat{p}_{\pi_i} \log \hat{p}_{\pi_i} \quad (2)$$

This entropy value provides a measure of the quality of a partial ordering induced by the tree reconstruction algorithm. For a balanced binary tree with $N=512$, we find that for the original ordering, $\hat{E}(\pi_{original}) = 2.2323$, and for the random ordering, $\hat{E}(\pi_{random}) = 2.702$. This motivates examining the estimated Δ -

entropy of our clustering reconstruction-based orderings as $\hat{E}_{\Delta}(\pi) = \hat{E}(\pi_{random}) - \hat{E}(\pi)$, where we normalize the reconstructed clustering entropy value with respect to a random permutation of the items. The quality of our clustering methodologies will be examined, where the larger the estimated Δ -entropy, the higher the quality of our estimated clustering.

For the synthetic binary tree experiments, we created a balanced binary tree with 512 items. We generated similarity between each pair of items such that $100 \cdot (1 - q)\%$ of the pairwise similarities chosen at random are consistent with the TC condition ($\in \mathbf{S}_C$). The remaining $100 \cdot q\%$ of the pairwise similarities were inconsistent with the TC condition. We examined the performance of both standard bottom-up agglomerative clustering and our Robust Clustering algorithm, **RAcluster**, for pairwise similarities with $q = 0.05, 0.15, 0.25$. The results presented here are averaged over 10 random realization of noisy synthetic data and setting the threshold $\gamma = 0.30$. We used the similarity voting budget $m = 80$, which requires 65% of the complete set of similarities. Performance gains are shown using our robust clustering approach in Table 2 in terms of both the estimated Δ -entropy and r_{min} , the size of the smallest correctly resolved cluster (where all clusters of size r_{min} or larger are reconstructed correctly for the clustering). This shows a clear correlation between high Δ -entropy and high clustering reconstruction resolution.

Table 2: Clustering Δ -entropy results for synthetic binary tree with $N = 512$ for Agglomerative Clustering and **RAcluster**.

q	Agglo. Clustering		RAcluster ($m=80$)	
	Δ -Entropy	r_{min}	Δ -Entropy	r_{min}
0.05	0.37	460.8	1.02	7.2
0.15	0.09	512	1.02	15.2
0.25	0.01	512	1.01	57.6

Our robust clustering methodology was then performed on the real world data sets using the threshold $\gamma = 0.30$ and similarity voting budgets $m = 20$ and $m = 40$. In addition to the quality of our cluster reconstruction (in terms of estimated Δ -entropy), the performance is also stated in terms of the number of pairwise similarities required by the agglomerative clustering methodology, denoted by n_{agg} , against the number of similarities required by our **RAcluster** method, n_{robust} . The results in Table 3 (averaged over 20 random permutations of the datasets) again show significant performance gains in terms of both the estimated Δ -entropy and the number of pairwise similarities required. Finally, in Figure 1 we see the reordered similarity matrices given both agglomera-

tive clustering and our robust clustering methodology, **RAcluster**.

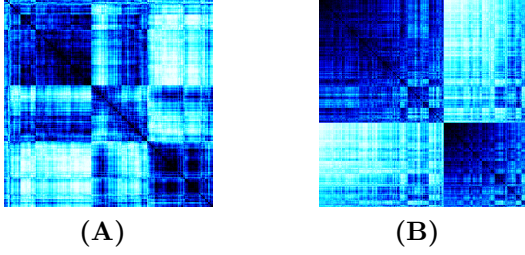


Figure 1: Reordered pairwise similarity matrices, Gene microarray data with $N = 1000$ using (A) - Agglomerative Clustering and (B) - Robust Clustering using $m = 20$ (requiring only 18.1% of the similarities).

6 Appendix

6.1 Proof of Theorem 4.1

Since the **outlier** tests can be erroneous, we instead use a two-round voting procedure to correctly determine whether two items x_i and x_j are in the same sub-cluster or not. Please refer to Algorithm 1 for definitions of the relevant quantities. The following lemma establishes that the outlier fraction values $c_{i,k}$ can reveal whether two items x_i, x_k are in the same subcluster or not, provided that the number of voting items $m = |\mathcal{S}_V|$ is large enough and the similarity $s_{i,k}$ is consistent.

Lemma 2. *Consider two items x_i and x_k . Under assumptions **A1** and **A2** and assuming $s_{i,k} \in \mathbf{S}_C$, comparing the outlier count values $c_{i,k}$ to a threshold γ will correctly indicate whether x_i, x_k are in the same subcluster with probability at least $1 - \frac{\delta_C}{2}$ for*

$$m \geq \frac{\log(4/\delta_C)}{2 \min \left((\gamma - 1 + (1 - q)^2)^2, ((1 - q)^2 \eta - \gamma)^2 \right)}.$$

Proof. Let $\Omega_{i,k} := \mathbf{1}_{\{s_{i,k} \in \mathbf{S}_C\}}$ be the event that the similarity between items x_i, x_k is in the consistent subset (see Definition 3). Under **A1**, the expected outlier fraction ($c_{i,k}$) conditioned on x_i, x_k and $\Omega_{i,k}$ can be bounded in two cases; when they belong to the same subcluster and when they do not:

$$\mathbb{E}[c_{i,k} \mid x_i, x_k \in \mathcal{C}_L \text{ or } x_i, x_k \in \mathcal{C}_R, \Omega_{i,k}] \geq (1 - q)^2 \eta$$

$$\begin{aligned} \mathbb{E}[c_{i,k} \mid x_i \in \mathcal{C}_R, x_k \in \mathcal{C}_L \text{ or } x_i \in \mathcal{C}_L, x_k \in \mathcal{C}_R, \Omega_{i,k}] \\ \leq (1 - (1 - q)^2) \end{aligned}$$

A2 stipulates a gap between the two bounds. Hoeffding's Inequality ensures that, with high probability, $c_{i,k}$ will not significantly deviate below/above the lower/upper bound. Thresholding $c_{i,k}$ at a level γ between the bounds will probably correctly determine

whether x_i and x_k are in the same subcluster or not. More precisely, if

$$m \geq \frac{\log(4/\delta_C)}{2 \min \left((\gamma - 1 + (1 - q)^2)^2, ((1 - q)^2 \eta - \gamma)^2 \right)},$$

then with probability at least $(1 - \delta_C/2)$ the threshold test correctly determines if the items are in the same subcluster. \square

Next, note that we cannot use the cluster count $c_{i,j}$ directly to decide the placement of x_i since the condition $s_{i,j} \in \mathbf{S}_C$ may not hold. In order to be robust to errors in $s_{i,j}$, we employ a second round of voting based on an independent set of m randomly selected *agreement* items, \mathcal{S}_A . The *agreement fraction*, $a_{i,j}$, is the average of the number of times the item x_i agrees with the clustering decision of x_j on \mathcal{S}_A .

Lemma 3. *Consider the following procedure:*

$$x_i \in \begin{cases} \mathcal{C}_j & : \text{if } a_{i,j} \geq \frac{1}{2} \\ \mathcal{C}_j^c & : \text{if } a_{i,j} < \frac{1}{2} \end{cases}$$

*Under assumptions **A1** and **A2**, with probability at least $1 - \frac{\delta_C}{2}$, the above procedure based on $m = |\mathcal{S}_A|$ agreement items will correctly determine if the items x_i, x_j are in the same subcluster, provided*

$$m \geq \frac{\log(4/\delta_C)}{2 \left((1 - \delta_C) (1 - q)^2 - \frac{1}{2} \right)^2}.$$

Proof. Define $\Phi_{i,j}$ as the event that similarities $s_{i,k}, s_{j,k}$ are both consistent (i.e., $s_{i,k}, s_{j,k} \in \mathbf{S}_C$) and thresholding the cluster counts $c_{i,k}, c_{j,k}$ at level γ correctly indicates if the underlying items belong to the same subcluster or not. Using Lemma 2 and the union bound, the conditional expectations of the agreement counts, $a_{i,j}$, can be bounded as,

$$\begin{aligned} \mathbb{E}[a_{i,j} \mid x_i \notin \mathcal{C}_j] &\leq P(\Phi_{i,j}^C) \leq 1 - (1 - q)^2 (1 - \delta_C) \\ \mathbb{E}[a_{i,j} \mid x_i \in \mathcal{C}_j] &\geq P(\Phi_{i,j}) \geq (1 - q)^2 (1 - \delta_C) \end{aligned}$$

Since $q \leq 1 - 1/\sqrt{2(1 - \delta')}$ and $\delta_C = \delta'/n$ (as defined below), there is a gap between these two bounds that includes the value $1/2$. Hoeffding's Inequality ensures that with high probability $a_{i,j}$ will not significantly deviate above/below the upper/lower bound. Thus, thresholding $a_{i,j}$ at $1/2$ will resolve whether the two items x_i, x_j are in the same or different subclusters with probability at least $(1 - \delta_C/2)$ provided $m \geq \log(4/\delta_C) / \left(2 \left((1 - \delta_C) (1 - q)^2 - \frac{1}{2} \right)^2 \right)$. \square

By combining Lemmas 2 and 3, we can state the following. The **split** methodology of Algorithm 1 will successfully determine if two items x_i, x_j are in the

Table 3: Δ -entropy results for real world datasets (gene microarray, breast tumor comparison, and Phylogenetics) using both Agglomerative Clustering and **RAcluster** algorithm.

	Agglo.	RAcluster ($m = 20$)		RAcluster ($m = 40$)	
Dataset	Δ -Entropy	Δ -Entropy	$\frac{n_{robust}}{n_{agg}}$	Δ -Entropy	$\frac{n_{robust}}{n_{agg}}$
Gene (N=500)	0.1561	0.1768	27%	0.1796	51%
Gene (N=1000)	0.1484	0.1674	18%	0.1788	37%
Tumor (N=400)	0.0574	0.0611	30%	0.0618	57%
Tumor (N=600)	0.0578	0.0587	24%	0.0594	47%
Phylo. (N=750)	0.0126	0.0141	21%	0.0143	41%
Phylo. (N=1000)	0.0149	0.0103	16%	0.0151	35%

same subcluster with probability at least $1 - \delta_C$ under assumption **A1** and **A2**, provided

$$m \geq \max \left(\frac{\log(4/\delta_C)}{2 \left((1 - \delta_C) (1 - q)^2 - \frac{1}{2} \right)^2}, \frac{\log(4/\delta_C)}{2 \min \left((\gamma - 1 + (1 - q)^2)^2, ((1 - q)^2 \eta - \gamma)^2 \right)} \right)$$

and the cluster under consideration has at least $2m$ items.

In order to successfully determine the subcluster assignments for all n items of the cluster \mathcal{C} with probability at least $1 - \delta'$, requires setting $\delta_C = \frac{\delta'}{n}$ (i.e., taking the union bound over all n items). Thus we have the requirement $m \geq c_0(\delta', \eta, q, \gamma) \log(4n/\delta')$ where the constant obeys

$$c_0(\delta', \eta, q, \gamma) \geq \max \left(\frac{1}{2 \left((1 - \delta') (1 - q)^2 - \frac{1}{2} \right)^2}, \frac{1}{2 \min \left((\gamma - 1 + (1 - q)^2)^2, ((1 - q)^2 \eta - \gamma)^2 \right)} \right) \quad (3)$$

Finally, this result and assumptions **A1-A2** imply that the algorithm **split**(\mathcal{C}, m, γ) correctly determine the two subclusters of \mathcal{C} with probability at least $1 - \delta'$.

6.2 Proof of Theorem 4.2

Lemma 4. *A binary tree with N leaves and balance factor $\eta_C \geq \eta$ has depth of at most $L \leq \log N / \log(\frac{1}{1-\eta})$.*

Proof. Consider a binary tree structure with N leaves (items) with balance factor $\eta \leq 1/2$. After depth of ℓ , the number of items in the largest cluster are bounded by $(1 - \eta)^\ell N$. If L denotes the maximum depth level, then there can only be 1 item in the largest cluster after depth of L , we have $1 \leq (1 - \eta)^L N$. \square

The entire hierarchical clustering can be resolved if all the clusters are resolved correctly. With a maximum depth of L , the total number of clusters M in the hierarchy is bounded by $\sum_{\ell=0}^L 2^\ell \leq 2^{(L+1)} \leq 2N^{1/\log(\frac{1}{1-\eta})}$, using the result of Lemma 4. Therefore, the probability that some cluster in the hierarchy is not resolved $\leq M\delta' \leq 2N^{1/\log(\frac{1}{1-\eta})}\delta'$ (where **split** succeeds with probability $> 1 - \delta'$). Therefore, for all clusters (which satisfy the conditions **A1** and **A2** of Theorem 4.1 and have size $> 2m$) can be resolved with probability $1 - \delta$ (by setting $\delta' = \frac{\delta}{2N^{1/\log(\frac{1}{1-\eta})}}$), from the proof of Theorem 4.1 we define $m = k_0(\delta, \eta, q, \gamma) \log(\frac{8}{\delta}N)$ where,

$$k_0(\delta, \eta, q, \gamma) \geq c_0(\delta, \eta, q, \gamma) / \left(1 + \left(1 / \log \left(\frac{1}{1-\eta} \right) \right) \right) \quad (4)$$

Given this choice of m , we find that the **RAcluster** methodology in Algorithm 2 for a set of N items will resolve all clusters that satisfy **A1** and **A2** of Theorem 4.1 and have size $> 2m$, with probability at least $1 - \delta$.

Furthermore, the algorithm only requires $O(N \log^2 N)$ total pairwise similarities. By running the **RAcluster** methodology, each item will have the **split** methodology performed at most $\log N / \log(\frac{1}{1-\eta})$ times (i.e., once for each depth level of the hierarchy). If $m = k_0(\delta, \eta, q, \gamma) \log(\frac{8}{\delta}N)$ for **RAcluster**, each call to **split** will require only $3k_0(\delta, \eta, q, \gamma) \log(\frac{8}{\delta}N)$ pairwise similarities per item. Given N total items, we find that the **RAcluster** methodology requires at most $3k_0(\delta, \eta, q, \gamma) N \frac{\log N}{\log(\frac{1}{1-\eta})} \log(\frac{8}{\delta}N)$ pairwise similarities.

Acknowledgements

This work was supported in part by AFOSR grant number FA9550-10-1-0382 and FA9550-09-1-0140. Any opinions, findings, conclusions or other recommendations expressed in this material are those of the authors and do not necessarily reflect the AFOSR.

References

- [1] H. Yu and M. Gerstein, "Genomic Analysis of the Hierarchical Structure of Regulatory Networks," in *Proceedings of the National Academy of Sciences*, vol. 103, 2006, pp. 14,724–14,731.
- [2] J. Ni, H. Xie, S. Tatikonda, and Y. R. Yang, "Efficient and Dynamic Routing Topology Inference from End-to-End Measurements," in *IEEE/ACM Transactions on Networking*, vol. 18, February 2010, pp. 123–135.
- [3] M. Girvan and M. Newman, "Community Structure in Social and Biological Networks," in *Proceedings of the National Academy of Sciences*, vol. 99, pp. 7821–7826.
- [4] R. K. Srivastava, R. P. Leone, and A. D. Shocker, "Market Structure Analysis: Hierarchical Clustering of Products Based on Substitution-in-Use," in *The Journal of Marketing*, vol. 45, pp. 38–48.
- [5] S. Chaudhuri, A. Sarma, V. Ganti, and R. Kaushik, "Leveraging Aggregate Constraints for Deduplication," in *Proceedings of SIGMOD Conference 2007*, pp. 437–448.
- [6] A. Arasu, C. Ré, and D. Suciu, "Large-Scale Deduplication with Constraints Using Dedupalog," in *Proceedings of ICDE 2009*, pp. 952–963.
- [7] R. Ramasubramanian, D. Malkhi, F. Kuhn, M. Balakrishnan, and A. Akella, "On The Treeness of Internet Latency and Bandwidth," in *Proceedings of ACM SIGMETRICS Conference*, Seattle, WA, 2009.
- [8] Sajama and A. Orlitsky, "Estimating and Computing Density-Based Distance Metrics," in *Proceedings of the 22nd International Conference on Machine Learning (ICML)*, 2005, pp. 760–767.
- [9] G. Karypis, E. Han, and V. Kumar, "CHAMELEON: A Hierarchical Clustering Algorithm Using Dynamic Modeling," in *IEEE Computer*, vol. 32, 1999, pp. 68–75.
- [10] S. Guha, R. Rastogi, and K. Shim, "ROCK: A Robust Clustering Algorithm for Categorical Attributes," in *Information Systems*, vol. 25, July 2000, pp. 345–366.
- [11] T. Hofmann and J. M. Buhmann, "Active Data Clustering," in *Advances in Neural Information Processing Systems (NIPS)*, 1998, pp. 528–534.
- [12] T. Zoller and J. Buhmann, "Active Learning for Hierarchical Pairwise Data Clustering," in *Proceedings of the 15th International Conference on Pattern Recognition*, vol. 2, 2000, pp. 186–189.
- [13] N. Grira, M. Crucianu, and N. Boujemaa, "Active Semi-Supervised Fuzzy Clustering," in *Pattern Recognition*, vol. 41, May 2008, pp. 1851–1861.
- [14] G. Kollios, D. Gunopulos, N. Koudas, and S. Berchtold, "Efficient biased sampling for approximate clustering and outlier detection in large data sets," *IEEE Transactions on Knowledge and Data Engineering*, vol. 15, pp. 1170–1187, September 2003.
- [15] M. Balcan and P. Gupta, "Robust Hierarchical Clustering," in *Proceedings of the Conference on Learning Theory (COLT)*, July 2010.
- [16] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*. Springer, 2001.
- [17] J. Pearl and M. Tarsi, "Structuring Causal Trees," in *Journal of Complexity*, vol. 2, 1986, pp. 60–77.
- [18] L. Li, D. Alderson, W. Willinger, and J. Doyle, "A First-Principles Approach to Understanding the Internet's Router-Level Topology," in *Proceedings of ACM SIGCOMM Conference*, 2004, pp. 3–14.
- [19] J. DeRisi, V. Iyer, and P. Brown, "Exploring the metabolic and genetic control of gene expression on a genomic scale," in *Science*, vol. 278, October 1997, pp. 680–686.
- [20] A. Frank and A. Asuncion, "UCI Machine Learning Repository," 2010. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [21] R. Finn, J. Mistry, and et. al., "The Pfam Protein Families Database," in *Nucleic Acids Research*, vol. 38, 2010, pp. 211–222.
- [22] S. Needleman and C. Wunsch, "A General Method Applicable to the Search for Similarities in the Amino Acid Sequence of Two Proteins," in *Journal of Molecular Biology*, vol. 48, 1970, p. 443453.