
Maximizing Agreements for Ranking, Clustering and Hierarchical Clustering via MAX-CUT

Vaggos Chatziafratis
Google Research NY

Mohammad Mahdian
Google Research NY

Sara Ahmadian
Google Research NY

Abstract

In this paper, we study a number of well-known combinatorial optimization problems that fit in the following paradigm: the input is a collection of (potentially inconsistent) local relationships between the elements of a ground set (e.g., pairwise comparisons, similar/dissimilar pairs, or ancestry structure of triples of points), and the goal is to aggregate this information into a global structure (e.g., a ranking, a clustering, or a hierarchical clustering) in a way that maximizes agreement with the input. Well-studied problems such as rank aggregation, correlation clustering, and hierarchical clustering with triplet constraints fall in this class of problems. We study these problems on stochastic instances with a hidden embedded ground truth solution. Our main algorithmic contribution is a unified technique that uses the maximum cut problem in graphs to approximately solve these problems. Using this technique, we can often get approximation guarantees in the stochastic setting that are better than the known worst case inapproximability bounds for the corresponding problem. On the negative side, we improve the worst case inapproximability bound on several hierarchical clustering formulations through a reduction to related ranking problems.

1 Introduction

In many learning/optimization problems, the input data is in the form of a number of *ordinal* judgements about the local relationships among a set of n items. A prominent example is the problem of ranking n al-

ternatives, where the input is often pairwise comparisons between these items. For example, sports teams are often ranked by aggregating the results of matches played between pairs of teams, and election outcomes are decided by aggregating individual votes.

Learning from comparisons has been prevalent across different domains, as humans are typically good at quickly answering *ordinal* questions (“which movie/restaurant/candidate do you prefer”), but often respond slowly and inaccurately to *cardinal* questions (“how much do you like this option”). In the psychology literature, the method of *paired comparisons* that has been in use since the 1920’s is based on this principle (see (Thurstone, 1959, Chapter 7)). Moreover, modern online platforms can organically extract such ordinal preferences by observing the users (e.g., “which movie did they first watch”, or “did they skip a search result and click on the next one”) and later use them for improving search or recommendation rankings (see, for example, Joachims (2002)). The same principle applies to settings other than ranking. For example, when trying to learn a clustering of n items, it is easier for a human judge to answer questions of the form “should x and y be in the same cluster” than to measure the similarity of x and y . Or, to reconstruct the evolutionary tree (also known as the phylogenetic tree) between n species, biologists often start by answering questions of the form “between three species x, y , and z , which two are evolutionarily closer”.

At the heart of each of these examples is the non-trivial algorithmic task of reconciling potentially inconsistent judgements into a global solution. This defines a number of algorithmic problems that we study in this paper. Though seemingly unrelated, all of these problems seek to find a global structure that has the maximum number of agreements with the given collection of local ordinal relationships. As we shall see later in the paper, the problems are also linked in that we can apply a common technique (based on graph max cuts) to them all. The problems, shown in Figure 1, fall under the three categories of ranking, clustering, and hierarchical clustering:

- **Ranking:** The goal is to find an ordering of n items. In the *Maximum Acyclic Subgraph* (MAS), the input is a number of pairwise comparisons of the form $a < b$. In *Betweenness*, the input is a number of triples $a|b|c$ meaning that b is between a and c in the ordering. In *Non-Betweenness*, the input is a number of triples $b|ac$ meaning that b is not between a and c .
- **Clustering:** In the *Correlation Clustering* problem, the goal is to find a partitioning of n items, and the input is a number of pairs of the form ab , meaning that a and b should be in the same cluster, and a number of pairs of the form $a|b$, meaning that a and b should be in different clusters.
- **Hierarchical clustering:** The goal is to find a (rooted or unrooted) tree with the set of n items as its leaves. In the *Desired Triplets* problem, the input is a number of triplets $ab|c$, meaning that the least common ancestor of a and b is a descendant of the least common ancestor of a, b , and c . In the *Desired Quartets* problem, the input is a number of quartets $ab|cd$, meaning that the unique path connecting a and b in the tree does not intersect with the unique path connecting c and d . The *Forbidden Triplets* and *Forbidden Quartets* problems are defined similarly with the opposite requirements.

These problems come from a variety of applications: MAS is a formulation of the rank aggregation problem and has many applications, e.g., in search ranking. Correlation Clustering is a central problem in unsupervised learning and data analysis (Bansal et al., 2004). Hierarchical clustering problems are motivated by applications in reconstructing phylogenetic trees (Felsenstein, 2004), and are also related to the objective-driven formulations of Dasgupta (2016), Moseley and Wang (2017) and Cohen-Addad et al. (2019) for hierarchical clustering. In fact, the Desired Triplets formulation described above is tightly connected with objective-based approaches for Hierarchical Clustering as can be seen in Charikar et al. (2019a,b). Betweenness and Non-Betweenness are motivated by applications in genome sequencing in bioinformatics (Slonim et al., 1997). We are interested in algorithms that can provide an approximation guarantee, i.e., a provable bound on the multiplicative factor between the solution found by the algorithm and the optimal solution. We will consider this problem both in the worst case and under a stochastic model with an embedded ground-truth solution. Our contribution is two-fold (see Table 1 for a summary):

On the positive side, in Section 3, under a simple stochastic model akin to the well-known *stochastic*

block model, we are able to improve upon worst-case approximations for all problems and in some cases (e.g., for problems on rankings and hierarchies) even overcome impossibility results. Interestingly, our algorithms are all based on variants of MAXCUT on graphs that can have both positive and negative weights and may also be directed. Some approaches for tree reconstruction based on MAXCUT had been used in previous experimental works (Snir and Rao, 2006, 2008, 2012), and in this way our work provides concrete proof for why these heuristics are reported to perform well on “real-world” instances. Our natural stochastic model captures “real-world” instances via an embedded ground-truth from which we generate “noisy” constraints, similar to the Stochastic Block Model (Mossel et al., 2012) in community detection.

On the negative side, we obtain new hardness of approximation results for four problems on hierarchical clustering: Forbidden Triplets, Desired Triplets, Forbidden Quartets, Desired Quartets. Briefly, we may refer to them as triplets/quartets consistency problems. These are instances of Constraint Satisfaction Problems (CSP) on trees (Bodirsky and Mueller, 2010; Bodirsky et al., 2016), analogous to SAT formulas in complexity. Even though such problems on hierarchies have been studied for decades, the current best approximations are achieved by trivial baseline algorithms. Our hardness results give some explanation why previous approaches were not able to obtain anything better. Our result on the Forbidden Triplets problem is tight and is the first tight hardness for CSPs on trees, extending analogous hardness results by Guruswami et al. (2011) from linear orderings (i.e., rankings) to trees. This is carried out in Section 4.

Our stochastic model for collecting information is the simplest form of embedded model on n items, and is motivated by crowdsourcing and biological applications (Vaughan, 2017; Kleindessner and von Luxburg, 2017; Ghoshdastidar et al., 2019; Snir and Yuster, 2012). We simply choose items at random and include a pairwise/triplet/quartet constraint depending on the task. For example, to generate constraints for the MAS problem on rankings, let π^* denote a ground-truth ranking (e.g., of chess players or ads to show a user). We select uniformly at random m pairs of items a_i, b_i and then we generate m pairs $a_i < b_i$; if a_i precedes b_i in π^* the constraint is included with probability $(1 - \epsilon)$, otherwise the opposite constraint is generated. Thus, some fraction of the constraints can be erroneous. After generating m (noisy) constraints in this way, our goal is to find a global solution (ranking, partition, or tree) that satisfies as many as possible.

Techniques: Our hardness reductions for Maximum Forbidden Triplets consistency are based on mapping


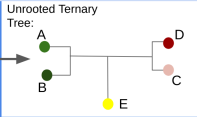
Problem Name	Types of Constraints	Types of Solutions
Max Acyclic Subgraph (MAS) Betweenness (BTW): Non-Betweenness (Non-BTW):	A<B, B<C, B<D, D<C A B C, B C D, D C A, B D C A BC, B DC, D CA, C DB	Ranking: A, B, C, D Ranking: A, B, C, D Ranking: A, B, C, D
Correlation Clustering:	Must-Link & Cannot-Link (AB) (A D) (BC) (B D) (CD) (D E)	Partition: A B C D E
Desired/Forbidden Triplets:	Desired & Forbidden AB C AC B CD B AD C AC D AB C	Rooted Binary Tree: 
Desired/Forbidden Quartets:	Desired & Forbidden AB CD AC BD AE CD AD CB ED AB AE CB AE BC AB ED AC BE AB CD	Unrooted Ternary Tree: 

Figure 1: A schematic representation of all problems considered in the paper. The left column has the problem names, the middle the types of constraints and the right column has a candidate solution. With green are constraints that are correctly resolved in the given candidate solution, whereas with red are those that are incorrect. For more examples, see Section 2.

trees to permutations on their leaves and back, and showing that any constant factor improvement over trivial baselines would refute the Unique Games Conjecture¹ (UGC) (Khot, 2002). Regarding our MAXCUT algorithm (see Algorithm 1), it is based on MAXCUT variations on directed and undirected graphs with negative weights and is conceptually simple. Briefly, given an instance for any of the problems we consider, we map it to a graph where edges encode the underlying constraints; perhaps the most intuitive such construction is for Correlation Clustering where a “must-link” or “cannot-link” constraint between items i, j is captured by a negative or positive edge (i, j) respectively. Then, we show how large (positive) cuts in this graph yield partitions that satisfy many of the constraints. The existence of a large cut can be guaranteed by analyzing our stochastic model and so an approximate MAXCUT algorithm can yield improvements over previous results. An interesting ingredient that we need for the case of MAS, is how to approximate the MAXCUT problem on *directed* graphs with both positive and negative weights which, to the best of our knowledge, hadn’t been analyzed before.

More broadly, we justify theoretically why prior experimental heuristics work and we extend them to work for new problems with provable approximation guarantees. Our work also presents the first case of a CSP on trees that is approximation resistant; recall that many important CSPs, including Max3SAT, are approximation resistant, i.e., it is NP-hard to approximate them better than a random assignment. This echoes the striking result by Håstad (2001) on ap-

¹Khot’s UGC is a major open question in complexity. We will not define it here as we only use some of its consequences on ordering problems (Guruswami et al., 2011).

proximation resistance of boolean CSPs to CSPs on trees and shows why no algorithmic improvement had been made in the worst-case, despite significant efforts (Byrka et al., 2010; Jiang et al., 2001; Bryant, 1997; He et al., 2006; Steel, 1992).

Table 1: Shown in bold are our improved hardness (column “Hardness”) and approximations under our stochastic model (column “Stochastic”). Column “Approx.” has prior approx. ratios. Also see Section 3 and Appendix A for the dependence on error parameter ϵ .

	Approx.	Hardness	Stochastic
MAS	1/2	1/2	0.642
BTW	1/3	1/3	0.402
NON-BTW	2/3	2/3	0.84
Correl. Cl.	0.76	APX-hard	0.82(*)
Forb. Triplet	2/3	2/3 (tight)	0.78(*)
Des. Triplet	1/3	2/3	0.64(*)
Forb. Quartet	2/3	8/9	0.672
Des. Quartet	1/3	2/3	0.425

Remark 1. We want to point out that all our approximation results here hold with high probability as a standard concentration argument about the stochastic process guarantees that the weight of the cuts is well-concentrated around its mean (as long as the number of generated constraints $m \geq \Omega(\log n)$).

Remark 2. Our results for ranking and quartets hold with no assumption on the optimal solution. For the positive results (denoted with $(*)$ in Table 1) via MAXCUT for correlation clustering and triplets however, we need a mild balancedness assumption, roughly stating that the optimal solution contains a relatively balanced $(\frac{1}{3} : \frac{2}{3})$ partition, to ensure the existence of a good cut in the ground-truth (see Appendix, Assumption 1). Usually, such assumptions are common in generative graph models for clustering, e.g., the Stochastic Block Model (Mossel et al., 2012; Abbe et al., 2015) and for hierarchical clustering, e.g., the Hierarchical Stochastic Block Model (Lyzinski et al., 2016; Cohen-Addad et al., 2019; Ghoshdastidar et al., 2019), where we expect to see at least two large communities emerge.

2 Background and Related Work

As the paper discusses multiple problems on rankings, partitions and hierarchies, we devote this section in describing the multitude of problems. A familiar reader can skip this section and proceed to Section 3.

There are 3 categories of problems we study here, depending on the type of the output: ranking (also called a *permutation* or a *leaf ordering* in biology (Bar-Joseph et al., 2001)), clustering (partitioning of the data points) and hierarchical clustering (also called

phylogenetic tree). There has been significant amounts of work on each of these tasks, that we only partially cover here as we go over our problems and results.

2.1 Optimization Problems and Types of Constraints

In all problems, we are given m constraints and we want to maximize the number of constraints satisfied by our output, whether it be a ranking, a partition or a hierarchy. We describe below the types of different constraints (see also Figure 1):

Ranking (i.e., a permutation or leaf ordering): Given n labels $\{1, 2, \dots, n\}$, we want to find a permutation that maximizes the number of satisfied constraints of the following form:

- **Pairwise comparisons:** A constraint here is of the form “ $a < b$ ”, indicating that in the output permutation, item a should precede b . If this information is encoded as a directed graph G with arcs $a \rightarrow b$, this gives rise to the Maximum Acyclic Subgraph (MAS) or Feedback Arc Set (FAS), two fundamental problems in computer science (Karp, 1972).
- **Betweenness (BTW) and Non-Betweenness (Non-BTW) constraints:** In the BTW problem (Opatrny, 1979; Chor and Sudan, 1998; Makarychev, 2012), we are given relative ordering constraints of the form $a|b|c$ indicating “ b should be between a and c ”. This allows for abc or cba out of the 6 possible orderings for the 3 labels. As the name suggests, NON-BTW is the complement of BTW, where a constraint $bc|a$ (equivalently $a|bc$) indicates that in the output permutation “ a should *not* lie between b and c ”. This allows for 4 valid relative orderings abc, acb, bca, cba . Generally, these are the two most common examples of ordering Constraint Satisfaction Problems (ordering CSPs) of arity 3 and are mainly motivated by applications in bioinformatics (Slonim et al., 1997). They have also played a major role in complexity (Guruswami et al., 2011; Austrin et al., 2013).

Just to give a sense of the approximability of these problems in the worst-case, the current best constant factor is a $\frac{1}{2}$ -approximation for MAS, a $\frac{1}{3}$ -approximation for BTW, and a $\frac{2}{3}$ -approximation for NON-BTW, all achieved by a *random* permutation. We also know that under the Unique Games Conjecture (UGC) of Khot (2002), the first two results are tight, whereas the third is tight under $P \neq NP$. Such problems, where a random output is provably the best, are called *approximation resistant* and have been studied extensively by theoreticians (Charikar et al., 2009; Guruswami et al., 2008; Hästad, 2001; Austrin

and Mossel, 2009). Our work gives strong evidence pointing to the fact that important CSPs on trees (triplets/quartets) may be approximation resistant.

Clustering: Here we want to maximize agreements with **Must-Link/Cannot-Link** constraints: The input is a graph with “+” or “-” edges indicating if the two endpoints should belong to the same cluster or not. Such constraints give rise to Correlation Clustering, an important paradigm for data analysis both in practice (Davidson and Basu, 2007; Wagstaff and Cardie, 2000; Wagstaff et al., 2001) and theory (Bansal et al., 2004; Ailon et al., 2008; Charikar et al., 2005; Swamy, 2004). The current best for maximizing agreements is a 0.7666 multiplicative approximation via semidefinite programs (Swamy, 2004) and an APX-hardness is known (Charikar et al., 2005). Here we will improve upon 0.7666, under our stochastic model for generating constraints.

Hierarchical Clustering (i.e., phylogenetic trees): There are two common types of trees: rooted and unrooted. Given n data points, a rooted binary tree on n leaves, where each leaf corresponds to a data point, is usually called a *hierarchical clustering* and is a standard tool for data analysis across different disciplines (Steinbach et al., 2000; Leskovec et al., 2014; Tumminello et al., 2010; Sørlie et al., 2001). Unrooted ternary trees (all nodes have degree 3, except the leaves that have degree 1) are usually called *phylogenetic trees* and are prevalent in computational biology as they describe speciation events throughout the evolution of species (Bryant, 1997; Felsenstein, 2004). Here we will use the two terms interchangeably to describe hierarchies on n leaves. Since in a hierarchy all data are eventually separated at the leaves, pairwise constraints no longer make sense and the analogue of “must-link/cannot-link” are so-called “must-link-before/cannot-link-before” constraints:

- **Desired/Forbidden Triplets:** The output here is a rooted binary tree T on n leaves. We say a triplet relation “ $t = ab|c$ ” is *obeyed* by T (or T *obeys* t), if the lowest common ancestor (LCA) of a, b is a descendant of the LCA of a, c in T . Otherwise T *disobeys* $ab|c$. A triplet can be *desired* (we write $t \in \mathcal{T}_D$) and we want the output T to obey it² or *forbidden* (we write $t \in \mathcal{T}_F$) and we want T to disobey/avoid it, giving rise to important optimization problems studied in computational biology and graph theory under the name of rooted triplets consistency (Steel, 1992; Bryant, 1997; Byrka et al., 2010; He et al., 2006). Notice that a forbidden triplet $ab|c$ is less restrictive, since it only specifies that T should either obey

²For example, “penguin, dolphin| tiger” could be a desired triplet as the tiger is the least relevant item.

$ac|b$ or $bc|a$, but not $ab|c$. This is reflected in the complexity of the problems: given a set of forbidden triplets, it is NP-complete to check consistency (i.e., if there is a tree avoiding all of them), whereas checking consistency of desired triplets in polynomial time was established long ago by Aho et al. (1981).

- **Desired/Forbidden Quartets:** The desired output here is a ternary unrooted tree T . We say a quartet $q = ab|cd$ is *obeyed* by T (or T *obeys* q) if the (unique) path from a to b in T does not share any vertices with the (unique) path from c to d in T . Otherwise T *disobeys* q . Similarly to triplets, a quartet can be *desired* ($q \in \mathcal{Q}_D$) or *forbidden* ($q \in \mathcal{Q}_F$), giving rise to important quartets consistency problems in biology and graph theory (Felsenstein, 2004; Bryant, 1997; Jiang et al., 2001; Snir and Rao, 2006). For both problems, even if the input is consistent, checking consistency is NP-complete.

Once again, just to give a sense of the approximability, for desired triplets or quartets, the current best is a $\frac{1}{3}$ -approximation and for forbidden triplets or quartets, the current best is a $\frac{2}{3}$ -approximation. Embarrassingly, in all four cases these are achieved by a random (rooted or unrooted) tree or a simple greedy construction (He et al., 2006).

2.2 Further Motivation and Related Work

Here, we further make a comparison to other relevant works. For ranking, many different types of probabilistic models have been considered (Braverman and Mossel, 2009; Shah et al., 2016; Shah and Wainwright, 2017; Negahban et al., 2012; Falahatgar et al., 2017) giving statistical guarantees for reconstructing the desired permutation. Instead of pairwise comparisons, the problem has also been studied in the case where partial rankings or complete information (“tournaments”) is provided (Fagin et al., 2006; Ailon, 2010; Kenyon-Mathieu and Schudy, 2007). Clustering with constraints and qualitative information (both max and min versions) were studied in Bansal et al. (2004); Charikar et al. (2005) where approximations via linear programs were derived or practical improvements were made possible (Wagstaff et al., 2001; Wagstaff and Cardie, 2000). In crowdsourcing and biological applications, both triplet and quartets queries have been deployed (Vinayak and Hassibi, 2016; Vaughan, 2017; Kleindessner and von Luxburg, 2017; Ghoshdastidar et al., 2019; Snir and Rao, 2006; Bryant, 1997) as they can be more intuitive for non-expert users compared to pairwise comparisons. Semi-supervised models, where triplet queries depend on answers to previous queries have been studied in Emamjomeh-Zadeh and Kempe (2018); Vikram and Dasgupta (2016).

To further motivate our stochastic model and results, we include a slightly more detailed comparison with 3 important prior works Braverman and Mossel (2009); Emamjomeh-Zadeh and Kempe (2018); Snir and Yuster (2012) that study “ground-truth” stochastic models similar to ours. The authors in Braverman and Mossel (2009) study the ranking problem and assume that there exists a ground-truth ranking π^* , as we do. However, their stochastic model assumes either that we have access to *all* $\binom{n}{2}$ pairwise comparisons, or that we have access to *complete* rankings σ on the n items, where each complete ranking σ is generated with probability inverse exponential in the Kemeny distance between π^* and σ (Kemeny distance is the number of inversions, i.e., the number of pairs ordered in π^* differently from σ).

As it will become obvious, their assumptions are much stricter than our simple stochastic model that generates m pairwise comparisons uniformly at random. Moreover, notice that our approximation guarantees hold for *any* number m of given constraints without requiring it to be $\Omega(n^2)$. Given their more refined model, they are of course in a position to analyze the maximum likelihood estimator and prove approximate recovery results, e.g., that no element is misplaced by more than $\log n$ positions with high probability; however no guarantees are given for the number of violated pairwise constraints, which is the focus of our paper.

For triplets hierarchical clustering, the authors in Emamjomeh-Zadeh and Kempe (2018) assume there exists a ground-truth binary tree T , as we do. However, they are allowed *adaptive* triplet queries and show that $\approx n \log n$ such queries suffice to recover T using a clever partition algorithm similar to Quicksselect and Quicksort. Once again, our model is not adaptive, and we do not pose any constraints on the number m of given constraints. For quartets hierarchical clustering, our model is similar to Snir and Yuster (2012), but we generalize their results to hold both for forbidden and desired quartets.

Finally, our constrained version of Hierarchical Clustering based on triplet constraints was studied in Chatziafratis et al. (2018) under the assumption that the input contains pairwise similarities as well as triplet constraints.

3 Using MaxCut on instances with embedded ground-truth

We present our main strategy MAXCUT behind our positive results. As we will see, by modifying the graphs, our method is flexible to allow for combinations of constraints, e.g., both BTW and NON-BTW constraints for rankings, or both desired and forbidden triplets (or quartets) for trees.

Stochastic Model for Generating Constraints:

Since our goal is to beat the worst-case approximation and hardness results, we use a simple stochastic model with an embedded ground-truth solution on n items. The form of the ground-truth changes depending on which problem we consider; it can be a ranking (for MAS, BTW, NON-BTW), a partition (for Correlation Clustering) or a hierarchical tree (rooted for Triplets and unrooted for Quartets). For generating the m input constraints, we simply choose items at random and with probability $(1 - \epsilon)$ we add a pairwise/triplet/quartet constraint that is consistent with the ground-truth, otherwise with probability ϵ we add an erroneous constraint on the selected items. For example, in the introduction, we saw the MAS constraints. Similarly, for BTW, we would uniformly at random pick m triples of items a, b, c and then add w.p. $(1 - \epsilon)$ the constraint $a|b|c$ if b appears in between a and c in the ground-truth ordering. Also, for the Triplets Consistency problem, we would again uniformly at random pick m triples of items a, b, c and then add w.p. $(1 - \epsilon)$ the constraint $ab|c$ if c is separated first from a, b in the ground-truth (rooted binary) tree. For all problems, after getting m (noisy) constraints in the analogous manner, our goal is to find a global solution that satisfies as many constraints as possible.

Positive Results: Using our stochastic model we can escape worst-case impossibility results and for all 3 categories of problems, we present improved approximation algorithms. At a high-level, we first construct a graph by encoding each of the local constraints on the items as a set of positive or negative edges between them. The graph captures the desired relationships and then, we find a good first split maximizing the ratio of satisfied over violated constraints by the cut. Naturally, our algorithm MAXCUT (see Algorithm 1) is based on variants of MAXCUT on graphs with negative weights. An interesting building block in our analysis when solving for better Maximum Acyclic Subgraphs, is the directed MAXCUT problem on graphs with negative weights which, to the best of our knowledge, hadn't been analyzed before. We note that for the triplets problem on trees, analogous MAXCUT heuristics had been successfully used before in experimental work for computational biology, however with no theoretical guarantees (Snir and Rao, 2006, 2012, 2008). An exception is the work of Snir and Yuster (2012), where they focus only on the desired quartets problem, however their analysis is a special case of ours for when $\mathcal{Q}_{\mathcal{F}} = \emptyset$ (i.e., the input contains no forbidden quartets). Our final approximations circumvent known hardness results for the case of rankings (Guruswami et al., 2011) and our new hardness results for trees described in detail later in Section 4.

3.1 Better Approximations for MAS

We start with MAS as it is perhaps the easiest to describe (see also Algorithm 1):

Theorem 1. *Given m constraints generated according to our stochastic model on n items, MAXCUT satisfies at least $(0.642 - 0.4285\epsilon)m$ on average, where ϵ is the fraction of erroneous comparisons. If moreover $m \geq \Omega(\log n)$, the result holds w.h.p.*

Remark 3. *For example, if the error parameter $\epsilon = 0.1$, hence 10% of the m generated constraints are erroneous, we still satisfy $\approx 60\%$ of them, and we still beat the previous best $\frac{1}{2}$ -approximation together with the known hardness (Guruswami et al., 2008).*

Our general proof template has 5 steps:

- Building a graph: For a sampled constraint $a < b$ indicating that a should precede b in the ranking, we add two directed edges:

$$+1 \text{ directed from } a \rightarrow b, -1 \text{ directed from } b \rightarrow a$$

Since the problem has orientation, we define the weight of a directed cut (S, \bar{S}) as the sum of all (positively or negatively) weighted arcs going from S to \bar{S} (and we ignore the arcs going from \bar{S} to S).

- Cuts and constraints: The goal of constructing the graph is to use information about its cuts and relate them to the pairwise constraints. Notice that a cut (S, \bar{S}) can either obey, disobey or leave unaffected the status of a $a < b$ constraint, depending on if a or b belongs to S or \bar{S} . Let m_s, m_v denote the satisfied, violated constraints by the cut, respectively. The weight of any directed (S, \bar{S}) cut is thus:

$$w(S, \bar{S}) = m_s(S, \bar{S}) - m_v(S, \bar{S}) \quad (1)$$

as satisfied pairs m_s (with $a \in S, b \in \bar{S}$) contribute $+1$ and violated pairs m_v (with $a \in \bar{S}, b \in S$) contribute -1 .

- Lower Bounding MAXCUT: The constructed graph from the first step, is directed and has both positive and negative weights. Based on eq. (1), we should find a large cut in this graph as this translates to many satisfied constraints. In order to find the cut, we use a MAXCUT variant that finds a cut comparable to the optimal max cut in graphs that are directed and contain both positive and negative weights. However, we cannot use the standard Goemans-Williamson algorithm and guarantees Goemans and Williamson (1995), as the graph is directed with positive and negative weights. A new ingredient in our proof is a semidefinite programming relaxation and analysis for this variant that achieves:

$$\mathbb{E}(w(S, \bar{S})) \geq 0.857w(\text{OPT}) - 0.143 \cdot W^- \quad (2)$$

where $w(\text{OPT})$ is the weight of the optimum cut and W^- is the total negative weight in the graph in absolute value. Based on the graph construction in this case, $W^- = m$ as every constraint contributed a -1 edge. We just note that the numerical values 0.143 and 0.857 sum to 1, and they just arise from the rounding scheme used to obtain an integral solution from the relaxation.

- Now that we have a lower bound for $w(S, \bar{S})$ based on the optimum cut, in order to conclude the algorithm’s cut is large (and hence satisfies many constraints), we need to lower bound the optimum’s cut weight $w(\text{OPT})$. To do this we consider the weight of a *median* directed cut: the median cut is defined to be the one that assigns the first $n/2$ labels in the optimum ordering for MAS, on one side of the cut, and the rest $n/2$ labels to the other side of the cut. Since the labels for the constraints according to our stochastic model were chosen at random, a counting argument implies that with high probability $\approx \frac{1}{2}m$ of the generated constraints are satisfied by the median cut and hence also by OPT . To see this, observe that for nearly half of the $a < b$ constraints, a belongs to the first $n/2$ labels of the median cut, whereas b belongs to the remaining $n/2$ labels. Since OPT is by definition even better than the median cut, we get that it has a large cut value. If we wanted to be slightly more precise, we should say that due to errors in an ϵ fraction of the generated constraints, we actually lose a small ϵ fraction of the constraints (we defer details to Appendix A) but this discounts the optimum cut only by a small amount.
- Output of MAXCUT : Finally, we need to find a good permutation overall, not just a good top split. Our algorithm starts by finding an approximate MAXCUT (S, \bar{S}) in G and then proceeds by outputting a *random* permutation on the items in S and in \bar{S} and concatenating them. Finally, we can compute the overall value of ALG (dropping the notation with (S, \bar{S})):

$$\begin{aligned} \text{ALG} &= m_s + \frac{1}{2}m_u = \\ &= m_s + \frac{1}{2}(m - m_s - m_v) = \frac{1}{2}m + \frac{1}{2}(w(S, \bar{S})) \end{aligned} \quad (3)$$

where m_u are the constraints that were unaffected by the (S, \bar{S}) cut. By eq. (3), we already see that we get some advantage over the $\frac{1}{2}m$ baseline which is optimal in the worst-case (and is achieved by a random permutation on all n items).

Remark 4. A natural question is to attempt to use MaxCut repeatedly on each of the two generated parts

of the first split. However analyzing the repeated MaxCut approach is not that simple, as once the first approximate MaxCut is performed, there is no randomness in the two generated subgraphs that we can exploit. Analogous difficulties arise in dissimilarity-based and quartets-based hierarchical clustering Charikar et al. (2019a); Snir and Yuster (2012); Ahmadian et al. (2020). Finally, we want to point out that such analyses are also known to be challenging from the literature on Random Forests for decision trees (e.g., Scornet et al. (2015)) where a similar (data-dependent) two-step analysis has been elusive.

3.2 Extensions to Other Problems

The same proof template as presented here can be modified to deal with the remaining problems: BTW, NON-BTW, forbidden and desired triplets, forbidden and desired quartets. As each of these constraints, involve 3 or 4 points, the construction and analyses become more involved. We present briefly the main modifications for the graph construction (see Appendix A for details).

For a BTW constraint $\{a|b|c\}$, we add undirected edges: $+2$ for (a, c) and -1 for $(b, a), (b, c)$. The edges capture that a cut violates the constraint if it separates b from a, c . For a NON-BTW constraint $\{ab|c\}$ indicating that c should not be between a, b in the final ordering, we add the following 3 undirected edges: $+1$ for pairs $(c, a), (c, b)$ and -2 for the pair (a, b) . Recall, that for BTW and NON-BTW, the ultimate goal is to beat the factors $\frac{1}{3}$ and $\frac{2}{3}$ which are currently optimal in the worst-case:

Theorem 2. Given $m = \Omega(\log n)$ noisy constraints on n items, variations of MAXCUT satisfy at least $(0.402 - 0.329\epsilon)m$ and $(0.845 - 0.329\epsilon)m$ constraints w.h.p. for BTW and NON-BTW, respectively, where ϵ is the fraction of erroneous constraints.

For Correlation Clustering, for each CANNOT-LINK constraint ab , we add a $+1$ for (a, b) , and for each MUST-LINK constraint ab , we add -3.2735 for edge (a, b) . The chosen numerical value -3.2735 depends on the current best 0.766-approximation for Correlation Clustering (Swamy, 2004) (see Appendix A).

Theorem 3. Given $m = \Omega(\log n)$ noisy “must-link/cannot-link” constraints on n items, MAXCUT (modified appropriately) satisfies at least $(0.8226 - 0.775\epsilon)m$ constraints w.h.p., where ϵ is the fraction of erroneous constraints.

Analogous theorems hold for the Triplets/Quartets consistency problems. Due to space constraints, we omit the statements but we refer the reader to Table 1 for the final ratios and to Appendix A for the proofs.

Algorithm 1 Our MAXCUT template as instantiated for MAS.

Input: m pairwise constraints for MAS.

1. For each $a < b$ constraint, insert a $+1$ arc directed from $a \rightarrow b$ and another arc with negative weight -1 directed from $b \rightarrow a$. Call the resulting graph G .
2. Run our approximate MAXCUT algorithm suitable for directed graphs with negative weights to get a first split (S, \bar{S}) , satisfying eq. (2).
3. Construct a random permutation π_1 on the nodes in S and a random permutation π_2 on the nodes in \bar{S} . Let π be the ranking obtained by concatenating π_1 and then π_2 .
4. Return π .

4 Hardness for CSPs on Trees

Negative Results: As mentioned, previous work (Byrka et al., 2010; Jiang et al., 2001; Bryant, 1997; He et al., 2006; Steel, 1992) tried to get better approximations for triplets/quartets consistency compared to trivial baselines. Recall, that the trivial baseline is to simply output a random tree (either rooted or unrooted depending on the problem). In our paper, near optimal hardness of approximation results for the maximum desired/forbidden triplets/quartets consistency problems (4 problems in total) are presented shedding light to why, despite significant efforts from different communities, no improvement had been made for nearly thirty years. As a consequence, we get the first tight hardness for an ordering problem on trees, thus extending the work of Guruswami et al. (2011) from orderings on the line to hierarchical clustering.

Specifically, for maximizing forbidden triplets, we show that no polynomial time algorithm can achieve a constant better than $\frac{2}{3}$ -approximation. Similar to Guruswami et al. (2008, 2011) this is assuming the Unique Games Conjecture, however for maximizing desired triplets, we show a threshold of $\frac{2}{3}$, assuming $P \neq NP$. The above also implies that forbidden triplets is approximation resistant as a random tree also achieves a $\frac{2}{3}$ factor. In fact our hardness results for all 4 problems are stronger, as we show it's not possible to distinguish almost perfectly consistent inputs from inputs where the optimum solution achieves almost the same as a random solution.

Technically, in order to get the hardness results, we give algorithms to obtain permutations on the leaves of a tree, such that if the tree obeyed many triplet/quartet constraints, then the permutation would also obey a large fraction of them when viewed as appropriate ordering constraints. Specifically, we prove that under the UGC, it is hard to approximate the Forbidden Triplets Consistency problem

better than a factor of $\frac{2}{3}$, even in the unweighted case.

Fact 1. *Let K be the total number of triplets constraints in an instance of BTW. For any $\epsilon > 0$, it is UGC-hard to distinguish between BTW instances of the following two cases:*

YES: $val(\pi^*) \geq (1 - \epsilon)K$, i.e. the optimal permutation satisfies almost all constraints.

NO: $val(\pi^*) \leq (\frac{1}{3} + \epsilon)K$, i.e. the optimal permutation does not satisfy more than $1/3$ fraction.

Given the above fact from Guruswami et al. (2011), we prove our $\frac{2}{3}$ -inapproximability result for Forbidden Triplets:

Theorem 4. *Let K be the total number of the triplet constraints in an instance of Forbidden Triplets Consistency. For any $\delta > 0$, it is UGC-hard to distinguish between the following two cases:*

YES: $val(T^*) \geq (1 - \delta)K$, i.e. the optimal tree satisfies almost all the triplet constraints.

NO: $val(T^*) \leq (\frac{2}{3} + \delta)K$, i.e. the optimal tree does not satisfy more than $\frac{2}{3}$ fraction of triplets.

Proof. Start with a YES instance of the BTW problem with optimal permutation π^* and $val(\pi^*) \geq (1 - \epsilon)K$. Viewing each BTW constraint $a|b|c$ as a forbidden triplet $ac|b$, we show how to construct a tree T such that $val(T) \geq (1 - \delta(\epsilon))K$. In fact, the construction is straightforward: simply assign the n labels, in the order they appear in π^* , as the leaves of a caterpillar tree (every internal node has its left child being a leaf). Observe that this caterpillar tree satisfies: $val(T) \geq (1 - \epsilon)K$. This is because if a BTW constraint $a|b|c$ was obeyed by π^* , it will also be avoided (viewed as a forbidden triplet $ac|b$) by the caterpillar tree above: if a appears first in the permutation then the caterpillar will avoid $ac|b$ as a gets separated first, otherwise if c appears first, then again the caterpillar tree will avoid $ac|b$ as c gets separated first.

The NO instance is more challenging. Start with a NO instance of the BTW problem with optimal π^* of value $val(\pi^*) \leq (\frac{1}{3} + \epsilon)K$. Viewing the BTW constraints as forbidden triplets, we show that the optimum tree T^* cannot achieve better than $> (2/3 + 2\epsilon)K$, because this would imply that $val(\pi^*) > (\frac{1}{3} + \epsilon)K$, which is a contradiction. For this, assume that some tree T scored a value $val(T) > (2/3 + 2\epsilon)K$. We will construct a permutation π from the tree T with value $val(\pi) > (1/3 + \epsilon)K$, a contradiction. Notice that there are forbidden triplets that may be avoided by the tree, yet obeyed by the permutation: for example for a forbidden triplet $t = ac|b$, the tree R that first removes a and then splits b, c will successfully avoid t , however the permutation acb can come from R by projection, however acb does not obey the BTW constraint $a|b|c$. Hence directly projecting the leaves of T onto a line may not satisfy $> (1/3 + 2\epsilon)K$, since every

forbidden triplet $ac|b$ avoided by T , can be ordered by this projected permutation in a way that would not obey the corresponding BTW constraint $a|b|c$. However, just by randomly swapping each left and right child for every internal node in the tree before we do the projection to the permutation, would satisfy $1/2 \cdot (2/3 + 2\epsilon)K = (1/3 + \epsilon)K$ number of constraints. To see this, note that with probability $\frac{1}{2}$ a forbidden $ac|b$ avoided by T will be mapped to the desired abc (and not acb) or cba (and not cab) ordering.

Finally, we get $val(\pi^*) \geq val(\pi) > (1/3 + \epsilon)K$, a contradiction that we were given a NO instance. To conclude, $\frac{2}{3}$ -inapproximability follows from the gap of these two instances. \square

For the Desired Triplets problem, the proof proceeds in a similar fashion. One main difference is that we prove hardness of $\frac{2}{3}$ under $P \neq NP$, without assuming UGC. The reason is that we reduce from the NON-BTW problem that is known to be approximation resistant, subject only to $P \neq NP$. Of course, one open question is to close the gap between this $\frac{2}{3}$ factor and the current best approximation of $\frac{1}{3}$.

Theorem 5. *Let K be the total number of the triplet constraints in an instance of Desired Triplets Consistency. For any $\delta > 0$, it is NP-hard to distinguish:*

YES: $val(T^*) \geq (\frac{1}{2} - \delta)K$

NO: $val(T^*) \leq (\frac{1}{3} + \delta)K$

Switching to quartet problems, our reductions are more challenging. The first challenge is that constraints are on 4 items so we need to resort to an ordering CSP of arity 4, that we term 4-SEPARATEDNESS. Next, trees are unrooted and we want to generate an ordering on their leaves. To do this we first root the tree at some internal node and then follow a similar strategy for randomly reordering their children. For desired quartets we show hardness of $\frac{2}{3}$ and for forbidden quartets a hardness of $\frac{8}{9}$ (see App. A for statements). Recall that the best approximations are $\frac{1}{3}$ and $\frac{2}{3}$ respectively, achieved by a random (unrooted) tree.

Remark 5. *Note that our hardness results give optimal results when restricted to (rooted or unrooted) caterpillar trees, an important tree family, where each internal node has at least one leaf as a child.*

5 Conclusion

We studied ranking, correlation clustering and hierarchical clustering under qualitative constraints and we presented a simple algorithm based on MAXCUT that is able to overcome known hardness results under a random model. We also provided the first tight hardness of approximation for CSPs on trees shedding light to basic problems in computational biology and extending previous results by Guruswami et al. (2011)

from ordering CSPs to trees. We believe that a nice open question is to prove that the two most important families of CSPs on trees (triplets and quartets consistency) are approximation resistant. Here we showed this for the case of forbidden triplets. More generally, it is conceivable that all non-trivial CSPs on trees are in fact approximation resistant, implying that the inapproximability results of Guruswami et al. (2011) can be extended from linear orderings to trees.

References

- E. Abbe, A. S. Bandeira, and G. Hall. Exact recovery in the stochastic block model. *IEEE Transactions on Information Theory*, 62(1):471–487, 2015.
- S. Ahmadian, V. Chatziafratis, A. Epasto, E. Lee, M. Mahdian, K. Makarychev, and G. Yaroslavtsev. Bisect and conquer: Hierarchical clustering via max-uncut bisection. *The 23rd International Conference on Artificial Intelligence and Statistics*, 2020.
- A. V. Aho, Y. Sagiv, T. G. Szymanski, and J. D. Ullman. Inferring a tree from lowest common ancestors with an application to the optimization of relational expressions. *SIAM Journal on Computing*, 10(3):405–421, 1981.
- N. Ailon. Aggregation of partial rankings, p-ratings and top-m lists. *Algorithmica*, 57(2):284–300, 2010.
- N. Ailon, M. Charikar, and A. Newman. Aggregating inconsistent information: ranking and clustering. *Journal of the ACM (JACM)*, 55(5):1–27, 2008.
- P. Austrin and E. Mossel. Approximation resistant predicates from pairwise independence. *Computational Complexity*, 18(2):249–271, 2009.
- P. Austrin, R. Manokaran, and C. Wenner. On the NP-hardness of approximating ordering constraint satisfaction problems. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 26–41. Springer, 2013.
- N. Bansal, A. Blum, and S. Chawla. Correlation clustering. *Machine Learning*, 56(1-3):89–113, 2004.
- Z. Bar-Joseph, D. K. Gifford, and T. S. Jaakkola. Fast optimal leaf ordering for hierarchical clustering. *Bioinformatics*, 17(suppl_1):S22–S29, 2001.
- M. Bodirsky and J. K. Mueller. The complexity of rooted phylogeny problems. In *Proceedings of the 13th International Conference on Database Theory*, pages 165–173, 2010.
- M. Bodirsky, P. Jonsson, and T. Van Pham. The complexity of phylogeny constraint satisfaction. In *33rd Symposium on Theoretical Aspects of Computer Science*, 2016.
- M. Braverman and E. Mossel. Sorting from noisy information. *arXiv preprint arXiv:0910.1191*, 2009.

- D. Bryant. Building trees, hunting for trees, and comparing trees: theory and methods in phylogenetic analysis. *PhD Thesis*, 1997.
- J. Byrka, S. Guillelot, and J. Jansson. New results on optimizing rooted triplets consistency. *Discrete Applied Mathematics*, 158(11):1136–1147, 2010.
- M. Charikar and V. Chatziafratis. Approximate hierarchical clustering via sparsest cut and spreading metrics. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 841–854. SIAM, 2017.
- M. Charikar, V. Guruswami, and A. Wirth. Clustering with qualitative information. *Journal of Computer and System Sciences*, 71(3):360–383, 2005.
- M. Charikar, V. Guruswami, and R. Manokaran. Every permutation csp of arity 3 is approximation resistant. In *2009 24th Annual IEEE Conference on Computational Complexity*, pages 62–73. IEEE, 2009.
- M. Charikar, V. Chatziafratis, and R. Niazadeh. Hierarchical clustering better than average-linkage. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2291–2304. SIAM, 2019a.
- M. Charikar, V. Chatziafratis, R. Niazadeh, and G. Yaroslavtsev. Hierarchical clustering for euclidean data. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2721–2730, 2019b.
- V. Chatziafratis, R. Niazadeh, and M. Charikar. Hierarchical clustering with structural constraints. In *International Conference on Machine Learning*, pages 774–783, 2018.
- B. Chor and M. Sudan. A geometric approach to betweenness. *SIAM Journal on Discrete Mathematics*, 11(4):511–523, 1998.
- V. Cohen-Addad, V. Kanade, F. Mallmann-Trenn, and C. Mathieu. Hierarchical clustering: Objective functions and algorithms. *Journal of the ACM (JACM)*, 66(4):1–42, 2019.
- S. Dasgupta. *A Cost Function for Similarity-Based Hierarchical Clustering*, page 118–127. Association for Computing Machinery, New York, NY, USA, 2016.
- I. Davidson and S. Basu. A survey of clustering with instance level constraints. *ACM Transactions on Knowledge Discovery from data*, 1(1-41):2–42, 2007.
- E. Emamjomeh-Zadeh and D. Kempe. Adaptive hierarchical clustering using ordinal queries. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 415–429. SIAM, 2018.
- R. Fagin, R. Kumar, M. Mahdian, D. Sivakumar, and E. Vee. Comparing partial rankings. *SIAM Journal on Discrete Mathematics*, 20(3):628–648, 2006.
- M. Falahatgar, A. Orlitsky, V. Pichapati, and A. T. Suresh. Maximum selection and ranking under noisy comparisons. In *International Conference on Machine Learning*, pages 1088–1096. PMLR, 2017.
- U. Feige and M. Goemans. Approximating the value of two power proof systems, with applications to max 2sat and max dicut. In *Proceedings Third Israel Symposium on the Theory of Computing and Systems*, pages 182–189. IEEE, 1995.
- J. Felsenstein. *Inferring phylogenies*, volume 2. Sinauer associates Sunderland, MA, 2004.
- D. Ghoshdastidar, M. Perrot, and U. von Luxburg. Foundations of comparison-based hierarchical clustering. In *Advances in Neural Information Processing Systems*, pages 7454–7464, 2019.
- M. X. Goemans and D. P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *Journal of the ACM (JACM)*, 42(6):1115–1145, 1995.
- V. Guruswami, R. Manokaran, and P. Raghavendra. Beating the random ordering is hard: Inapproximability of maximum acyclic subgraph. In *2008 49th Annual IEEE Symposium on Foundations of Computer Science*, pages 573–582. IEEE, 2008.
- V. Guruswami, J. Håstad, R. Manokaran, P. Raghavendra, and M. Charikar. Beating the random ordering is hard: Every ordering csp is approximation resistant. *SIAM Journal on Computing*, 40(3):878–914, 2011.
- J. Håstad. Some optimal inapproximability results. *Journal of the ACM (JACM)*, 48(4):798–859, 2001.
- Y.-J. He, T. N. Huynh, J. Jansson, and W.-K. Sung. Inferring phylogenetic relationships avoiding forbidden rooted triplets. *Journal of Bioinformatics and Computational Biology*, 4(01):59–74, 2006.
- T. Jiang, P. Kearney, and M. Li. A polynomial time approximation scheme for inferring evolutionary trees from quartet topologies and its application. *SIAM Journal on Computing*, 30(6):1942–1961, 2001.
- T. Joachims. Optimizing search engines using clickthrough data. In *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '02, page 133–142, New York, NY, USA, 2002. Association for Computing Machinery.
- R. M. Karp. Reducibility among combinatorial problems. In *Complexity of computer computations*, pages 85–103. Springer, 1972.

- C. Kenyon-Mathieu and W. Schudy. How to rank with few errors. In *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pages 95–103, 2007.
- S. Khot. On the power of unique 2-prover 1-round games. In *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*, pages 767–775. ACM, 2002.
- M. Kleindessner and U. von Luxburg. Kernel functions based on triplet comparisons. In *Advances in Neural Information Processing Systems*, pages 6807–6817, 2017.
- J. Leskovec, A. Rajaraman, and J. D. Ullman. *Mining of massive datasets*. Cambridge university press, 2014.
- V. Lyzinski, M. Tang, A. Athreya, Y. Park, and C. E. Priebe. Community detection and classification in hierarchical stochastic blockmodels. *IEEE Transactions on Network Science and Engineering*, 4(1):13–26, 2016.
- Y. Makarychev. Simple linear time approximation algorithm for betweenness. *Operations research letters*, 40(6):450–452, 2012.
- B. Moseley and J. Wang. Approximation bounds for hierarchical clustering: Average linkage, bisecting k-means, and local search. In *Advances in Neural Information Processing Systems*, pages 3094–3103, 2017.
- E. Mossel, J. Neeman, and A. Sly. Stochastic block models and reconstruction. *arXiv preprint arXiv:1202.1499*, 2012.
- S. Negahban, S. Oh, and D. Shah. Iterative ranking from pair-wise comparisons. In *Advances in neural information processing systems*, pages 2474–2482, 2012.
- J. Opatrny. Total ordering problem. *SIAM Journal on Computing*, 8(1):111–114, 1979.
- E. Scornet, G. Biau, J.-P. Vert, et al. Consistency of random forests. *The Annals of Statistics*, 43(4):1716–1741, 2015.
- N. Shah, S. Balakrishnan, A. Guntuboyina, and M. Wainwright. Stochastically transitive models for pairwise comparisons: Statistical and computational issues. In *International Conference on Machine Learning*, pages 11–20, 2016.
- N. B. Shah and M. J. Wainwright. Simple, robust and optimal ranking from pairwise comparisons. *The Journal of Machine Learning Research*, 18(1):7246–7283, 2017.
- D. Slonim, L. Kruglyak, L. Stein, and E. Lander. Building human genome maps with radiation hybrids. *Journal of Computational Biology*, 4(4):487–504, 1997.
- S. Snir and S. Rao. Using max cut to enhance rooted trees consistency. *IEEE/ACM transactions on computational biology and bioinformatics*, 3(4):323–333, 2006.
- S. Snir and S. Rao. Quartets maxcut: a divide and conquer quartets algorithm. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 7(4):704–718, 2008.
- S. Snir and S. Rao. Quartet maxcut: a fast algorithm for amalgamating quartet trees. *Molecular phylogenetics and evolution*, 62(1):1–8, 2012.
- S. Snir and R. Yuster. Reconstructing approximate phylogenetic trees from quartet samples. *SIAM Journal on Computing*, 41(6):1466–1480, 2012.
- T. Sørlie, C. M. Perou, R. Tibshirani, T. Aas, S. Geisler, H. Johnsen, T. Hastie, M. B. Eisen, M. Van De Rijn, S. S. Jeffrey, et al. Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications. *Proceedings of the National Academy of Sciences*, 98(19):10869–10874, 2001.
- M. Steel. The complexity of reconstructing trees from qualitative characters and subtrees. *Journal of classification*, 9(1):91–116, 1992.
- M. Steinbach, G. Karypis, and V. Kumar. A comparison of document clustering techniques. In *KDD workshop on text mining*, volume 400, pages 525–526. Boston, 2000.
- C. Swamy. Correlation clustering: maximizing agreements via semidefinite programming. In *Proceedings of the fifteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 526–527. Society for Industrial and Applied Mathematics, 2004.
- L. L. Thurstone. *The Measurement of Values*. The University of Chicago Press, 1959.
- M. Tumminello, F. Lillo, and R. N. Mantegna. Correlation, hierarchies, and networks in financial markets. *Journal of economic behavior & organization*, 75(1):40–58, 2010.
- J. W. Vaughan. Making better use of the crowd: How crowdsourcing can advance machine learning research. *The Journal of Machine Learning Research*, 18(1):7026–7071, 2017.
- S. Vikram and S. Dasgupta. Interactive bayesian hierarchical clustering. In *International Conference on Machine Learning*, pages 2081–2090, 2016.
- R. K. Vinayak and B. Hassibi. Crowdsourced clustering: Querying edges vs triangles. In *Advances in Neural Information Processing Systems*, pages 1316–1324, 2016.
- K. Wagstaff and C. Cardie. Clustering with instance-level constraints. *AAAI/IAAI*, 1097:577–584, 2000.

K. Wagstaff, C. Cardie, S. Rogers, and S. Schrödl.
Constrained k-means clustering with background
knowledge. In *ICML*, volume 1, pages 577–584,
2001.