

Hierarchical Clustering: a 0.585 Revenue Approximation

Noga Alon¹, Yossi Azar², and Danny Vainstein³

¹Princeton University, USA and Tel Aviv University, Israel. nalon@math.princeton.edu

²School of Computer Science, Tel-Aviv University, Israel. azar@tau.ac.il

³School of Computer Science, Tel-Aviv University, Israel. dannyvainstein@gmail.com

January 30, 2020

Abstract

Hierarchical Clustering trees have been widely accepted as a useful form of clustering data, resulting in a prevalence of adopting fields (phylogenetics, image analysis, bioinformatics and many more). Recently, Dasgupta (STOC 16[']) began the work of analyzing these types of algorithms through the lenses of approximation. Later, the dual problem was considered by Moseley and Wang (NIPS 17[']) dubbing it the Revenue goal function. In this problem, given a nonnegative weight w_{ij} for each pair $i, j \in [n] = \{1, 2, \dots, n\}$, the objective is to find a tree T whose set of leaves is $[n]$ that maximizes the function $\sum_{i,j \in [n]} w_{ij}(n - |T_{ij}|)$, where $|T_{ij}|$ is the number of leaves in the subtree rooted at the least common ancestor of i and j .

In our work we consider the revenue goal function and prove the following results. First, we prove the existence of a bisection (i.e., a tree of depth 2 in which the root has two children, each being a parent of $n/2$ leaves) which approximates the general optimal solution up to a factor of $\frac{1}{2}$ (which is tight). Second, we apply this result in order to prove a $\frac{2}{3}p$ approximation for the general revenue problem, where p is defined as the approximation ratio of the MAX-UNCUT BISECTION problem. Since p is known to be at least 0.8776 [3] [17], we get a 0.585 approximation algorithm for the revenue problem. This improves a sequence of earlier results which culminated in an 0.4246-approximation guarantee.

1 Introduction

The notion of Hierarchical Clustering (HC) trees has been introduced and subsequently studied for several decades. The notion was first accepted due to its applications to the realm of phylogenetics [15] [13]. Here, given genomic similarities between species our goal is to hierarchically cluster the species in a way that captures fine-grained relations between different species. Since then, the notion of HC trees has expanded to many more fields (for a survey on the subject see [6]).

Typically, schemes for generating HC trees fall into one of two categories: Agglomerative algorithms (i.e., bottom-up) and Divisive algorithms (i.e., top-down). Agglomerative algorithms initially start with a partition such that each data point is its own partition set. The algorithm then proceeds to recursively merge different sets, terminating once all data points are contained in the same set. Notably, the well-studied average linkage algorithm is an example of such (specifically, merging the two sets maximizing their average induced weight). On the other hand, divisive algorithms start with a single set containing all the data points, and then proceed to recursively split sets, terminating once all data points remain alone in their set.

Recently, Dasgupta [12], formally defined the notion of a "good" HC tree. Under such a definition, he elegantly bridged the gap between HC trees and the field of approximation algorithms

by defining a minimization cost function (see related work for an extended discussion on this goal function). Thereafter, Moseley and Wang [14] considered the complementary maximization variant of this problem - namely, the revenue goal function.

Both definitions considered the following model. Assume we are given n data points with some notion of similarity between them. The similarity is formally captured through similarity-edges between any two data points, defined by $G = (V, E, w)$, where $w(\cdot) : V \rightarrow \mathbb{R}^+$. Our goal is to construct an HC tree, T , such that its leaves are in a 1-1 correspondence with our data points, V . Note that given such a tree, every internal node (1) represents a set of data points (i.e., those in its subtree) and (2) partitions this set (i.e., each of the node's children corresponds to a subset of the original set).

Intuitively speaking, since higher weighted edges correspond to more similar data points, it seems that a good HC tree would split these nodes low in the tree. Therefore, Moseley and Wang defined the revenue problem as,

$$\max_T R_G(T) = \max_T \sum_{e=\{ij\} \in E} w_{ij}(|V| - |T_{ij}|),$$

where T is an HC tree, T_{ij} denotes the subtree rooted at the least-common-ancestor (LCA) of data points i and j and $|T_{ij}|$ denotes the amount of data points in T_{ij} .

In [14], Moseley and Wang considered several algorithms. Notably, they considered the random algorithm, that simply splits data points randomly at every cut (henceforth denoted by *RAND*), and the average-linkage algorithm. They showed that both yield a revenue which is a 1/3 of the optimal solution. Subsequently, Charikar et al. [8] showed that one can beat average-linkage through the use of semi-definite programming, improving the bound to 0.3364. Recently, Ahmadian et al. [1] managed to leverage the MAX-UNCUT BISECTION (MUB) problem in order to prove a 0.4246 approximation. In our paper we improve upon this result and show an improved approximation of 0.585.

Our contributions. We consider the revenue goal function and prove the following results.

- We show that for any revenue instance, there exists a bisection, X , that is, a tree of depth 2 in which the root has two children, each being the parent of $n/2$ leaves, such that $R(X) \geq \frac{1}{2}OPT$, where OPT denotes the revenue gained by the optimal tree (see Theorem 1). In order to show such existence we make use of two random processes: we randomly fix the optimal tree and then randomly generate our bisection, X . We emphasize the fact that even though OPT makes use of an arbitrarily deep tree, it is enough to consider a single cut in order to gain half the revenue.
- Using our result regarding the existence of a large revenue generating bisection, we prove a 0.585 approximation for the revenue problem. We note that in fact we show a $\frac{2}{3}p$ approximation where $p = 0.8776$ is the best known approximation for the MUB problem.

Remark 1. *The algorithm we consider is that which solves the MUB problem for the first cut and then proceeds using the random algorithm. This algorithm has already been used in [1]. They showed a bound of 0.4246 for the algorithm that outputs the revenue gained by the maximum between this algorithm and the random algorithm. Surprisingly we show that the former algorithm on its own is enough to yield an approximation of 0.585.*

Techniques. Our first result makes use of a new upper bound on the optimal solution (which may be of its own interest). Specifically, given an optimal solution, we embed its leaves on a line such that its root is above the line and we have no resulting crossing edges. This clearly yields an ordering of the leaves (see Figure 1). We then consider the distance between any two data points,

i and j within this ordering (simply the difference in rank) and observe that this is in fact a lower bound on $|T_{ij}|$. This in turn yields an upper bound on the optimal solution.

Next, we make use of this bound by (1) randomly generating a bisection that gains revenue that is "large" with respect to the bound and (2) in expectation this bound is "far" from the optimal solution. Both "large" and "far" will be formally defined later on.

Related work. Dasgupta [12] kicked off the line of work considering HC trees within the realm of approximation algorithms. In his paper, he considered similarity-edges and defined the cost of an HC tree as,

$$\min_T C_G(T) = \min_T \sum_{e=\{ij\} \in E} w_{ij} |T_{ij}|.$$

Note that the revenue goal function is in fact complementary to that of Dasguptas (that is, the optimal solution is the same for both, albeit with different goal function values).

In [12], many general properties pertaining to this goal function were discussed. More notably, they showed that this goal function is intuitive in that (1) on complete graphs (with no structure) all HC trees yield the same cost, (2) on disconnected graphs, optimal HC trees begin by splitting disconnected components and (3) they show modularity of the goal function. They further present an $O(\log^{1.5}(n))$ approximation via recursive sparsest cut. Later, both [7] and [11] showed that this algorithm is in fact an $O(\sqrt{\log n})$ approximation. In the hardness domain, [12] showed that the problem is NP-hard via a reduction to a variant of NAE-SAT. This was later improved by [7], showing that in fact no constant approximation exists (assuming the Small Set Expansion hypothesis). [10] managed to overcome the latter worst-case specific result by considering an average case defined by a stochastic block model and its hierarchical extension. Here, they managed to show an $O(1)$ approximation.

Following Dasgupta's work, Cohen-Addad et al. [11] considered the case of dissimilarity-edges. In this case, Dasgupta's cost function is now translated to a maximization problem. Given an HC tree, T , we now denote its gained value as its gained dissimilarity, $D_G(T)$. In this setting, both the random algorithm (i.e., at every step, split the data points randomly), henceforth denoted as *RAND*, and the average-linkage algorithm yield dissimilarity values of 2/3 of the optimal solution. Charikar et al. [8] improved upon this by proving an approximation of 0.6671 which makes use of a more delicate multi-phase algorithm.

Several other extensions to the formerly defined HC goal functions were also considered. One of these extensions includes that of structural constraints. Specifically, every constraint appears in the form of " $i, j|k$ " for data points i, j and k . A constraint is then considered satisfied if $k \notin T_{ij}$, for an HC tree, T . Aho et al. [2] considered this problem in the phylogenetic realm where this notion gives rise to the problem of constructing a phylogenetic tree that satisfies a set of lineage constraints on common ancestors. This notion has more recently been investigated in the domain of HC, by Charikar et al. [9]. In their paper they extended Dasguptas goal function to include structural constraints and showed an $O(k\sqrt{\log n})$ approximation where k is the number of constraints. For several more related variations of the vanilla HC see [5], [4], [16].

2 Notation

Given a revenue instance, $G = (V, E, w)$, where $V = \{1, 2, \dots, n\}$, we denote by *OPT* both the optimal revenue tree and its yielded revenue (it will be clear from context which of these definitions we will be referring to). Given an HC tree, T , we denote by $R(T)$ the revenue it yields and for any similarity edge, e , we denote by $R_T(e)$ the revenue gained by the edge e with respect to the HC tree, T (i.e., $R_T(e) = w_e(n - |T_e|)$).

3 Existence of a High Revenue Bisection

Theorem 1. *For any revenue instance and corresponding optimal solution, OPT , there exists a bisection, X , such that,*

$$R(X) \geq \frac{1}{2}OPT,$$

and this is asymptotically tight.

Before we prove the theorem, we introduce some notation. Given the optimal tree, OPT , we may fix its leaves in several different orders; each order is produced by representing OPT as a planar graph where its leaves are all embedded on a line, the root is above the line and all edges of the tree are straight lines going down from each parent to its children with no crossing edges. Denote each such ordering by the function $\pi : V \rightarrow [n]$. E.g., if we fix the tree such that leaf i is in fact first in the left-to-right resulting order, then $\pi(i) = 1$.

Now, given such an ordering, π , and an edge $e = \{i, j\}$, let $y_e^\pi = |\pi(i) - \pi(j)| + 1$, denote the distance between leaves i and j in the fixed tree, OPT . For a pictorial example, see Figure 1.

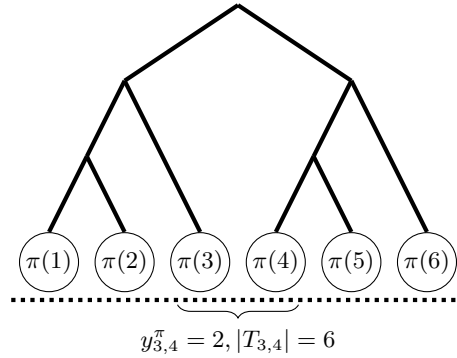


Figure 1: Embedding an HC tree with order π .

In order to prove Theorem (1) we first upper bound the revenue gained by OPT .

Observation 1. $\forall \pi, \forall e \in E : y_e^\pi \leq |T_e|$, where both π and T_e are defined with respect to OPT .

In itself, this will not result in a good enough bound on OPT . Therefore, we show that on average, y_e^π is far from T_e .

Lemma 2. $E_\pi[y_e^\pi] = \frac{|T_e|}{2} + 1$, where the expectation is taken over a uniform distribution of the orderings of OPT .

Proof. Let $e = \{i, j\}$. It can be shown through simple induction that for any tree with n leaves, and for any leaf, k , $E[\pi(k)] = \frac{n+1}{2}$. Therefore, if we denote by n_i and n_j the number of leaves in the subtrees containing i and j , each rooted at a separate child of i and j 's least-common-ancestor, then,

$$E_\pi[y_e^\pi] = (1/2)(n_i + \frac{n_j + 1}{2} - \frac{n_i + 1}{2} + 1) + (1/2)(n_j + \frac{n_i + 1}{2} - \frac{n_j + 1}{2} + 1) = \frac{|T_e|}{2} + 1.$$

□

Let $Y_\pi = \sum_{e \in E} w_e \cdot y_e^\pi$. By linearity of expectation we get the following lemma.

Lemma 3. *There exists an ordering of OPT , π^* , such that,*

$$Y_{\pi^*} \leq \sum_{e \in E} w_e \cdot \left(\frac{|T_e|}{2} + 1 \right).$$

Proof. By linearity of expectation and Lemma 2,

$$E_{\pi}[Y_{\pi}] = \sum_{e \in E} w_e \cdot E[y_e^{\pi}] = \sum_{e \in E} w_e \left(\frac{|T_e|}{2} + 1 \right).$$

Therefore, there exists an ordering as needed. \square

Next we show that there exists a distribution over the set of all bisections with high (to some degree) revenue with respect to the revenue gained by considering y_e^{π} rather than e .

Lemma 4. *Given any ordering of OPT , π , there exists a distribution, P_X , over all bisections, X , such that for any edge $e \in E$,*

$$E_{P_X}[R_X(e)] \geq \frac{1}{2} w_e (n - 2y_e^{\pi} + 2).$$

Proof. Fix π . We relabel the leaves of OPT such that $\pi(i) = i$ (this is simply to ease the notation). Next we define a distribution over all bisections, X . Consider the following random process: choose x uniformly at random from $\{1, \dots, \frac{n}{2}\}$. Then define X to be the bisection, $\{x, x+1, \dots, x+\frac{n}{2}-1\}, \{1, 2, \dots, x-1, x+\frac{n}{2}, x+\frac{n}{2}+1, \dots, n\}$.

Now consider some edge, $e \in E$. If $y_e^{\pi} \geq \frac{n}{2} + 1$, then $n - 2y_e^{\pi} + 2 \leq 0$. $R_X(e)$ is always nonnegative (since it is defined as the revenue gained by e) and thus the assertion of the lemma holds in this case.

Otherwise, $y_e^{\pi} \leq \frac{n}{2}$ for $e = \{i, j\}$. In this case, the probability that i and j are cut by the bisection is at most $\frac{y_e^{\pi}-1}{n/2} = \frac{2y_e^{\pi}-2}{n}$. Furthermore, since X is a bisection, any uncut edge yields a revenue of $n/2$. Overall,

$$E_{P_X}[R_X(e)] = w_e \left(1 - \frac{2y_e^{\pi} - 2}{n} \right) \frac{n}{2},$$

completing the proof of the lemma. \square

Recall that $R(X)$ is defined such that, $R(X) = \sum_{e \in E} R_X(e)$. Therefore, by combining lemmas (3) and (4), we may sum over all edges and get the following observation.

Observation 2. *Given the ordering, π^* , guaranteed by Lemma (3), we get,*

$$E_{P_X}[R(X)] \geq \sum_{e \in E} \frac{1}{2} w_e (n - 2y_e^{\pi^*} + 2).$$

We are now ready to prove Theorem 1.

Proof of Theorem 1. Fix an ordering of OPT to be π^* as guaranteed by Lemma (3). Therefore, by Observation (2), there exists a bisection X^* such that

$$R(X^*) \geq E_{P_X}[R(X)] \geq \sum_{e \in E} \frac{1}{2} w_e (n - 2y_e^{\pi^*} + 2).$$

Thus, by Lemma (3),

$$\begin{aligned}
R(X) &\geq \sum_{e \in E} \frac{1}{2} w_e (n - 2y_e^{\pi^*} + 2) \\
&= \sum_{e \in E} \frac{1}{2} w_e (n) - Y_{\pi^*} + \sum_{e \in E} w_e \\
&\geq \sum_{e \in E} \frac{1}{2} w_e (n) - \sum_{e \in E} w_e \cdot \left(\frac{|T_e|}{2} + 1 \right) + \sum_{e \in E} w_e \\
&= \frac{1}{2} \sum_{e \in E} w_e (n - |T_e|) \\
&= \frac{1}{2} OPT.
\end{aligned}$$

To show that the result is asymptotically tight, simply consider the instance G which is a matching with weight 1 for each of its edges. In this case any tree of depth 2, T generates a revenue of $R(T) \leq (\frac{1}{2} + o(1))OPT$. □

4 A 0.585 Approximation for the Revenue Goal Function

We define the approximation algorithm as a 2 step process: first we cut all data points using some black box algorithm that produces a bisection, denoted henceforth as ALG_{MUB} . Thereafter, we continue splitting each cluster randomly. Denote the combined algorithm by ALG .

In order to show an approximation bound for ALG , we first need to consider the MUB problem. In this problem we are given a weighted graph and our goal is to create a bisection maximizing the weights of uncut edges. We note that if we restrict ourselves to revenue trees which are bisections (i.e., the first cut splits the data points into two equal sets, then each set is cut only once using a "star" subtree - see Figure 2), the MUB optimal solution and the revenue optimal solution are in fact the same.

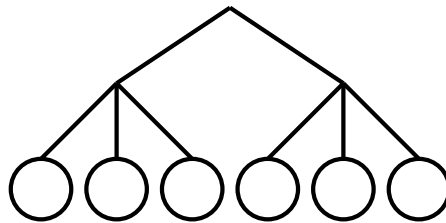


Figure 2: Example of a bisection revenue tree - revenue is only gained through edges uncut by the first cut.

Let p denote the approximation ratio algorithm ALG_{MUB} guarantees for the MUB problem. Note that p is at least 0.8776 (see [3] [17]). Thus, by defining ALG_{MUB} to be such an algorithm, the following theorem shows that algorithm ALG is a 0.585 approximation for the revenue problem.

Theorem 5. *Algorithm ALG guarantees an approximation ratio of $\frac{2}{3}p = 0.585$ for the revenue case, where p is defined as the approximation for the MUB problem.*

Proof. Let T_{ALG} denote the tree generated by algorithm ALG . It is a known simple fact that the random algorithm generates a revenue of at least $\frac{1}{3}nw_e$ for any edge e . Therefore, if we denote by

W_L and W_R the weight of the uncut edges generated by the first cut of our algorithm, then,

$$R(T_{ALG}) \geq W_L\left(\frac{n}{2} + \frac{1}{3} \cdot \frac{n}{2}\right) + W_R\left(\frac{n}{2} + \frac{1}{3} \cdot \frac{n}{2}\right) = (W_L + W_R)\frac{2n}{3}.$$

Denote by W_{L^*} and W_{R^*} the weights of the uncut edges generated by the optimal MUB solution and let p denote our first cut's approximation with respect to the MUB problem. Furthermore, let X^* denote the optimal solution to the revenue problem restricted to bisections. As noted earlier, X^* corresponds to W_{L^*} and W_{R^*} . Therefore,

$$W_L + W_R \geq p(W_{L^*} + W_{R^*}) = \frac{2p}{n}R(X^*).$$

Let OPT denote the revenue gained by the optimal solution. Thus, leveraging Theorem 1 yields,

$$\begin{aligned} R(T_{ALG}) &\geq \frac{2n}{3}(W_L + W_R) \geq \\ &\frac{4p}{3}(R(X^*)) \geq \\ &\frac{2p}{3}OPT. \end{aligned}$$

Since p is known to be at least 0.8776, we get that ALG is a $\frac{2}{3}p = 0.585$ approximation algorithm. \square

References

- [1] Sara Ahmadian, Vaggos Chatziafratis, Alessandro Epasto, Euiwoong Lee, Mohammad Mahdian, Konstantin Makarychev, and Grigory Yaroslavtsev. Bisect and conquer: Hierarchical clustering via max-uncut bisection. *CoRR*, abs/1912.06983, 2019.
- [2] Alfred V. Aho, Yehoshua Sagiv, Thomas G. Szymanski, and Jeffrey 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.
- [3] Per Austrin, Siavosh Benabbas, and Konstantinos Georgiou. Better balance by being biased: A 0.8776-approximation for max bisection. *ACM Trans. Algorithms*, 13(1):2:1–2:27, 2016.
- [4] Pranjal Awasthi, Maria Florina Balcan, and Konstantin Voevodski. Local algorithms for interactive clustering. *The Journal of Machine Learning Research*, 18(1):75–109, 2017.
- [5] Maria-Florina Balcan and Avrim Blum. Clustering with interactive feedback. In *International Conference on Algorithmic Learning Theory*, pages 316–328. Springer, 2008.
- [6] Pavel Berkhin. A survey of clustering data mining techniques. In *Grouping Multidimensional Data - Recent Advances in Clustering*, pages 25–71. 2006.
- [7] Moses Charikar and Vaggos Chatziafratis. Approximate hierarchical clustering via sparsest cut and spreading metrics. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2017, Barcelona, Spain, Hotel Porta Fira, January 16-19*, pages 841–854, 2017.

- [8] Moses Charikar, Vaggos Chatziafratis, and Rad Niazadeh. Hierarchical clustering better than average-linkage. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2019, San Diego, California, USA, January 6-9, 2019*, pages 2291–2304, 2019.
- [9] Vaggos Chatziafratis, Rad Niazadeh, and Moses Charikar. Hierarchical clustering with structural constraints. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, pages 773–782, 2018.
- [10] Vincent Cohen-Addad, Varun Kanade, and Frederik Mallmann-Trenn. Hierarchical clustering beyond the worst-case. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 6201–6209, 2017.
- [11] Vincent Cohen-Addad, Varun Kanade, Frederik Mallmann-Trenn, and Claire Mathieu. Hierarchical clustering: Objective functions and algorithms. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2018, New Orleans, LA, USA, January 7-10, 2018*, pages 378–397, 2018.
- [12] Sanjoy Dasgupta. A cost function for similarity-based hierarchical clustering. In *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2016, Cambridge, MA, USA, June 18-21, 2016*, pages 118–127, 2016.
- [13] N Jardine and R Sibson. A model for taxonomy. *Mathematical Biosciences*, 2(3-4):465–482, 1968.
- [14] Benjamin Moseley and Joshua Wang. Approximation bounds for hierarchical clustering: Average linkage, bisecting k-means, and local search. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 3094–3103, 2017.
- [15] Peter HA Sneath and Robert R Sokal. Numerical taxonomy. *Nature*, 193(4818):855–860, 1962.
- [16] Sharad Vikram and Sanjoy Dasgupta. Interactive bayesian hierarchical clustering. In *International Conference on Machine Learning*, pages 2081–2090, 2016.
- [17] Chenchen Wu, Donglei Du, and Dachuan Xu. An improved semidefinite programming hierarchies rounding approximation algorithm for maximum graph bisection problems. *J. Comb. Optim.*, 29(1):53–66, 2015.