EL-Clustering: Combining Upper- and Lower-Bounded Clusterings for Equitable Load Constraints

Rajni Dabas* rajni.dabas@pgdav.du.ac.in

Department of Computer Science, Northwestern University P.G.D.A.V. College, University of Delhi

Neelima Gupta[†] ngupta@cs.du.ac.in

Department of Computer Science, University of Delhi

Rudra Bhardwaj rudramsc24@cs.du.ac.in

Department of Computer Science, University of Delhi

Sapna Grover

Separtment of Computer Science, University of Delhi

Separtment of Computer Science, University of Delhi

Department of Computer Science, University of Deini Dyal Singh College, University of Delhi

Reviewed on OpenReview: https://openreview.net/forum?id=EkjDfnJ1gU

Abstract

The application of an ordinary clustering algorithm may yield a clustering output where the number of points per cluster (cluster size) varies significantly. In settings where the centers correspond to facilities that provide a service, this can be highly undesirable as the cluster size is essentially the service load for a facility. While prior work has considered imposing either a lower bound on the cluster sizes or an upper bound, imposing both bounds simultaneously has seen limited work, especially for the k-median objective, despite its strong practical motivation. In this paper, we solve the equitable load (EL) clustering problem where we minimize the k-median objective subject to the cluster sizes not exceeding an upper bound or falling below a lower bound. We solve this problem using a modular approach. Specifically, given a clustering solution that satisfies the lower bound constraints and another that satisfies the upper bound constraints, we introduce a combination algorithm which essentially combines both solutions to produce one that satisfies both constraints simultaneously at the expense of a bounded degradation in the k-median objective and a slight violation of the upper bound. Our combination algorithm runs in $O(k^3 + n)$ time, where n is the number of points and is faster than standard k-median algorithms that satisfy either the lower or upper bound constraints. Interestingly, our results can be generalized to various other clustering objectives, including the k-means objective. We also do empirical evaluation for k-Median objective on benchmark datasets to show that both, the cost as well as the violation factor are significantly smaller in practice than the theoretical worst-case guarantees.

1 Introduction

Decision-making using algorithms powered by machine learning has become ubiquitous. Routinely, algorithms are used in consequential applications such as loan approval (Sheikh et al. (2020); Kadam et al. (2021)), recidivism prediction (Travaini et al. (2022); Kovalchuk et al. (2023)), and kidney exchange (Ashlagi & Roth (2021); McElfresh et al. (2020)). This has naturally brought greater attention to the broader

^{*}This work was done when the auhtor was a research scholar at Department of Computer Science, University of Delhi

 $^{^{\}dagger}$ Corresponding author.

impact of these deployed algorithms and their vulnerability to noise and adversarial attacks, as well as their societal consequences in terms of fairness and privacy. These additional considerations imply that there is a real need to solve non-standard variants of many problems to overcome these possible harmful consequences.

Variants of the standard clustering problem that take such considerations into account have received significant attention from the research community; this is unsurprising since clustering is a fundamental problem in unsupervised learning and a classical problem in operations research. Examples of such works include Jones et al. (2021); Gupta et al. (2010); Kaplan & Stemmer (2018), who show algorithms for solving the k-median and k-means problems that preserve the privacy of individuals using differential privacy. Further, Chhabra et al. (2020) and Cinà et al. (2022) study the performance of clustering algorithms when the dataset is affected by adversarial corruptions. Moreover, fairness considerations in clustering have received even greater attention comparatively (Chierichetti et al. (2017); Bercea et al. (2019); Bera et al. (2019); Kleindessner et al. (2019); Ahmadi et al. (2022); Chen et al. (2019); Li et al. (2021); Chakrabarti et al. (2022); Awasthi et al. (2022)¹). The fair clustering literature has introduced a number of well-motivated fairness notions. At least seven different fairness notions have been introduced so far in clustering.

Despite the significant attention that has resulted in many variants of the classical clustering problem, we identify a simple notion that has not received much attention from the community, even though it is wellmotivated and has clear societal consequences. Specifically, from the operations research point of view, the selected centers in a clustering could represent facilities such as schools, with the cluster associated with each center (school) being the students assigned to that school. Naturally, a school requires a minimum number of students to maintain a good teaching quality² at the same time the number of students should not exceed a certain threshold as the school's resources might be over-consumed leading to a degradation in the teaching quality. One can also find a similar motivation if the schools were instead service centers providing services to clients instead of students. Each service center would want a minimum number of clients to bring in revenue; at the same time, the number of clients should not exceed a threshold, as that would lead to issues such as higher waiting time and lower service quality. At a more precise level, this notion, which we call equitable service load (EL), simply states that the size of each cluster (number of points in the cluster) should be both lower and upper bounded by some pre-set values simultaneously. This EL notion can also be motivated in machine learning applications such as market segmentation. Specifically, since points in the same cluster would receive the same ads, we might want to have a level of equity between the different ads (centers) so that none receive too little revenue or dominate the market.

Though clustering under **EL** constraints has not received much attention, this notion is not entirely new. Specifically, the literature has considered variants of the standard clustering problem where lower and upper bounds on the cluster sizes have to be satisfied simultaneously (Friggstad et al. (2016); Gupta et al. (2021); Ding et al. (2017); Rösner & Schmidt (2018)). However, the k-median variant of this problem (with lower and upper bounds) remains unsolved. Only heuristics have been introduced for the stringent case where the lower and upper bounds coincide (i.e., set to the same value); see de Maeyer et al. (2023) and citations within. Forcing the upper and lower bounds to be exactly equal is not practical since in most settings a small difference would be tolerated even if an exact equality was desired.

In this paper, we solve the k-median problem under \mathbf{EL} constraints. Unlike the prior work, we follow a modular approach. Specifically, using a solution where the cluster sizes are all lower bounded and another where the cluster sizes are upper bounded, we introduce a post-processing algorithm that combines the two to give a new solution that satisfies the \mathbf{EL} constraints. Our combination algorithm runs in $O(k^3+n)$ time, where n is the number of points, and as such does not present a heavy computational burden. In fact, in comparison to existing algorithms for the k-median problem subject to either a lower or upper bound constraint on the cluster size, our combination algorithm does not present the computational bottleneck in the algorithmic pipeline. Interestingly, although our main target is the k-median problem, we show how we can use our combination algorithm to solve other clustering variants under the \mathbf{EL} constraints, including k-means clustering.

¹See the references therein for more.

²It is well-known and documented that interaction between students can improve the educational and social outcomes (Soller (2001); Hurst et al. (2013)) but this would not be possible with a very small student body.

Organization of the Paper. In Section 2, we give our notation along with the formal statement of the EL problem and some background on relevant prior work. In Section 3, we state our main theoretical results. In Section 4, we give an overview of additional related work to our problem. In Section 5, we give a high-level discussion of our main algorithmic techniques. In Section 6, we present our algorithm for EL Clustering for the k-median objective along with detailed technical proofs that establish the guarantees of the algorithm. This is followed by a conclusion and a discussion of future work in Section 7. In addition, Appendix A discusses improvements to some of the theoretical guarantees when the gap between the lower and upper bounds is sufficiently large. Appendix B presents the algorithmic modifications required for other variants of the problem, while Appendix C contains the experimental setup and empirical evaluation.

2 Notation, Problem Statement, and Background

In our problem, we are given a set of locations P in a metric space with metric $c: P \times P \to \mathbb{R}_{\geq 0}$, a subset $\mathcal{C} \subseteq P$ of n many points to be clustered. Further, we are given a set of (potential) centers³ $\mathcal{F} \subseteq P$ and a positive integer k. Following the standard terminology in clustering and facility location, we will also refer to the given n points \mathcal{C} as clients and to the centers \mathcal{F} as facilities. As in standard k-median clustering, our objective is to find a set of facilities $\mathcal{F}' \subseteq \mathcal{F}$ of at most k facilities (i.e., $|\mathcal{F}'| \leq k$, note that this is called the cardinality constraint) and an assignment function $\sigma: \mathcal{C} \to \mathcal{F}'$ which assigns clients to the selected facilities so as to minimize the sum of distances between the clients and their assigned facilities. More formally, we want to obtain a solution $S = (\mathcal{F}', \sigma)$ that minimizes the objective function $Cost(S) = \sum_{j \in \mathcal{C}} c(j, \sigma(j))$. Furthermore, in **EL** Clustering we are additionally given two parameters L and U where L is a lower bound on the cluster size and U is the upper bound. It follows that in a valid **EL** Clustering, the size of any cluster is constrained to lie in the range [L, U]. More precisely, denoting the set of points assigned to a center $i \in \mathcal{F}'$ by $\sigma^{-1}(i)$, it follows that $L \leq |\sigma^{-1}(i)| \leq U$. From the above description, the formal and concise definition of the **EL** Clustering problem is as follows:

EL Clustering

Input: Instance $I = (P, c, C, \mathcal{F}, k, U, L)$ Optimization: $\min_{\mathcal{F}', \sigma} \sum_{j \in C} c(j, \sigma(j))$ subject to

- subset $\mathcal{F}' \subseteq \mathcal{F}$ of size at most k,
- assignment $\sigma: \mathcal{C} \to \mathcal{F}'$ such that for each facility $i \in \mathcal{F}'$, $L \leq |\sigma^{-1}(i)| \leq U$

We use three standard notions from the literature on approximation algorithms:

- 1. **Approximation Factor.** An algorithm is said to have an α -approximation factor if, for every instance, the cost of the solution it produces is at most α times the cost of an optimal solution.
- 2. Violation Factor. A solution is said to *violate* the upper bound constraint by a factor $\beta \geq 1$ if the solution exceeds the constraint by at most a multiplicative factor of β . For example, in **EL** clustering, a violation factor β means that no cluster has more than βU clients.
- 3. **Bi-criteria Approximation.** A solution is an (α, β) -bi-criteria approximation if it returns a solution of cost at most α times optimal and violates upper bound constraint by at most a factor of β .

We next define two k-MEDIAN problems satisfying partial **EL** constraints (dropping one of the two, lower or upper bound constraints) whose solutions will be combined to obtain a solution for **EL** Clustering. In UPPER BOUNDED k-MEDIAN (UkM), we drop the lower bounds (which is equivalent to setting L=0), whereas in LOWER BOUNDED k-MEDIAN (LkM), we drop the upper bounds (by setting U=n). Both of these problems have been extensively studied in the theoretical computer science literature. Though interesting, upper bound constraints are notoriously hard to handle in problems like k-median. For example, finding a constant factor approximation for UkM is one of the famous and long-standing open questions in the literature of

³Note that it is common in the classical k-median clustering to select the centers from the same set of points to be clustered, this can be captured in our formulation by simply setting setting $\mathcal{F} = \mathcal{C}$.

approximation algorithms. On the other hand, some heuristics do not provide any approximation ratio guarantees (i.e., a bound on the cost of the solution as compared to the cost of the optimal solution). On the positive side, there are several papers that solve the problem by giving bi-criteria approximations that have an approximation ratio for the clustering objective, but also violate either the upper bounds or the cardinality constraint by a small multiplicative factor (Byrka et al. (2015; 2016); Charikar et al. (1999); Demirci & Li (2016); Korupolu et al. (2000); Li (2014; 2015; 2016)). For the $\mathbf{L}k\mathbf{M}$ problem, both constant factor approximations and heuristics have been obtained (Arutyunova & Schmidt (2021); Guo et al. (2020); Han et al. (2020b;a)).

3 Our Results

In this paper, we study **EL** Clustering to obtain Theorem 3.1.

Theorem 3.1. Given a solution S_U for UPPER BOUNDED k-MEDIAN (UkM) with an upper bound violation of factor β and a solution S_L for LOWER BOUNDED k-MEDIAN (LkM). If the clustering costs of the solutions are $Cost(S_U)$ and $Cost(S_L)$, respectively. Then a solution of cost at most $(7Cost(S_U) + 2Cost(S_L))$ can be obtained for **EL** Clustering at a violation of the upper bound by a factor of $(\beta + 1)$ and in a run-time of $O(k^3 + n)$.

Note that the $O(k^3 + n)$ runtime stated above refers only to the combination step that merges the given solutions S_U and S_L ; it does not include the time required to compute S_U and S_L , which depends on the algorithms used for UPPER BOUNDED k-MEDIAN and LOWER BOUNDED k-MEDIAN. The theorem also establishes bounds on the clustering cost and the violation in the upper bound constraints using any two given solutions, even if these solutions result from heuristics. It follows that if we obtain these solutions using approximation algorithms for UPPER BOUNDED k-MEDIAN and LOWER BOUNDED k-MEDIAN, then we can establish approximation ratio guarantees for our **EL** solution as shown in the corollary below:

Corollary 3.2. Given an α_U approximation for UPPER BOUNDED k-MEDIAN with β violation in the upper bounds and an α_L approximation for LOWER BOUNDED k-MEDIAN, a $7\alpha_U + 2\alpha_L$ approximation can be obtained for **EL** Clustering with a $(\beta + 1)$ violation in upper bounds in $O(k^3 + n)$ time.

Applying existing algorithms in Corollary 3.2 immediately yields concrete bounds: using the 16-approximation with 3-factor upper bound violation of Charikar et al. (1999) for $\mathbf{U}k\mathbf{M}$ and the 387-approximation for $\mathbf{L}k\mathbf{M}$ of Han et al. (2020a), we obtain an 886-approximation with a 4-factor violation. Alternatively, with the $O(1/\varepsilon^2)$ -approximation of Byrka et al. (2016) for $\mathbf{U}k\mathbf{M}$ and small $\varepsilon > 0$, the violation factor approaches 2. The constants here are driven by the underlying algorithms for $\mathbf{U}k\mathbf{M}$ and $\mathbf{L}k\mathbf{M}$, and any future improvements in these directly improve our guarantees.

Moreover, our combination approach works independently of the specific technique used in these algorithms; the result holds even if the input solutions are produced by heuristics without formal guarantees, in which case our algorithm preserves their practical performance. For instance, by applying the Fixed Parameter Tractable (FPT) approximation algorithms⁴ for UkM and LkM by Goyal et al. (2020), the approximation factor for EL clustering can be reduced to $(27+\epsilon)$ with a 2-factor violation in upper bounds, while maintaining an FPT runtime in k.

Furthermore, the computational overhead incurred by our algorithm in combining the solutions is only $O(k^3 + n)$ whereas all algorithms⁵ for **U**k**M** and **L**k**M** require solving a linear programming problem and hence takes at least $\omega(n^4)$ time (Vaidya (1989); Jiang et al. (2020); Cohen et al. (2021); van den Brand (2020)).

Interestingly, our technique can be extended to other clustering variants, such as k-MEANS, k-CENTER, FACILITY LOCATION and KNAPSACK MEDIAN in the presence of **EL** constraints. Although we are able to extend the result to the k-Means problem, the constants associated with the cost of generating a k-means clustering with equitable load are relatively high in our paper. Improving these constants remains an

⁴An algorithm is FPT if its runtime is upper bounded by $O(f(k) \cdot n^c)$ where c is a constant but f(k) can be exponential in k, see Cygan et al. (2015) for more details.

⁵Except Korupolu et al. (2000) that uses local search but opens $(5 + \epsilon)k$ facilities instead of k.

interesting open question for future work. We mainly focus on k-median in the paper; modifications in the algorithm for other problems can be found in the supplementary material (Appendix B).

In the supplementary material (Appendix A), we show an improvement in the upper bound violation for a particular scenario when the gap between the lower and the upper bounds is not too small, specifically when $2L \leq U$. Note that this is a reasonable scenario that is likely to occur in real applications. For this special case, we reduce the violation in the upper bounds to $(\beta + \epsilon)$ at the expense of an increase of a factor of $O(1/\epsilon)$ in the cost for a constant $\epsilon > 0$.

Finally, we complement our theoretical guarantees with an empirical evaluation (Section C) on benchmark datasets, showing that the cost as well as the violation factor are substantially better than the worst case guarantees in most of the cases. It is observed that the cost overhead is within 20% of the cost of maximum of the $\mathbf{L}k\mathbf{M}$ and $\mathbf{U}k\mathbf{M}$ solutions and it is typically within 10–11%, whereas the violation factor is within 1.54 in $\approx 80\%$ of the cases.

4 Additional Related Work

As mentioned earlier, **EL** Clustering has not received much attention from the community. Heuristics are known for the problem when the lower and upper bounds coincide (Höppner & Klawonn (2008); Dinler & Tural (2016); de Maeyer et al. (2023); Lin et al. (2019); Ganganath et al. (2014); Chakraborty & Das (2019); Tang et al. (2019)). However, forcing the upper and lower bounds to be exactly equal is highly impractical since in most settings only lower and upper bounds are desired. Lei et al. (2013) provide heuristics for k-means clustering with **EL** constraints. Approximation algorithms have been obtained for clustering objectives other than the k-median and k-means. For example, Friggstad et al. (2016) gave an approximation algorithm for Facility location⁶ with **EL** constraints violating both the bounds by a constant factor with a trade-off in them whereas Gupta et al. (2021) gave an approximation algorithm violating the upper bounds by a factor of 5/2. For k-center with **EL** constraints, Ding et al. (2017) and Rösner & Schmidt (2018) independently gave constant factor approximations. To the best of our knowledge, the k-median problem with **EL** constraints has not been studied before in the literature.

Some prior works in fair clustering bear some resemblance to our work. For example, in settings where the points in the dataset belong to different demographic groups, the works of Chierichetti et al. (2017); Bercea et al. (2019); Bera et al. (2019); Esmaeili et al. (2020); Ahmadian et al. (2019) have considered a fairness notion where each cluster is constrained to have close to population-level proportions of each group. While this notion is similar to ours, there is a considerable difference since the bounds are not imposed on the cluster sizes as we do, but rather the proportions of the groups in each cluster. Further, another notion in fair clustering imposes lower and upper bounds not on the proportions of the demographic groups in each cluster but on the number of centers selected from each demographic group (Kleindessner et al. (2019); Jones et al. (2020); Hotegni et al. (2023)), i.e., in a dataset that consists of 50% from a "blue" group and 50% from a "red" group, then if we cluster with k=10, it may be desired to have at least 3 centers selected from each group and at most 7 centers from one group, thereby ensuring both a measure of diversity and restricted dominance in the selected centers. However, this notion is also different from the EL notion. Interestingly, Dickerson et al. (2024) present a modular approach to combine both demographic fairness notions mentioned above simultaneously. Although the above-mentioned demographic notions are different from imposing lower and upper bounds on the cluster sizes, the objective of Dickerson et al. (2024) of combining two notions simultaneously is similar to our objective. Further, their modular approach of post-processing existing solutions is also similar to our approach, although at a high level, our constraints and techniques are very different.

 $^{^{6}}$ In the facility location problem, instead of a hard bound k on the number of facilities, every facility has a facility opening cost and the goal now is to minimize the total cost of opening a subset of facilities and serving the clients from these opened facilities

⁷We denote the demographic groups with colors as done in fair clustering papers. Concretely, these colors could denote attributes such as age, gender, or race.

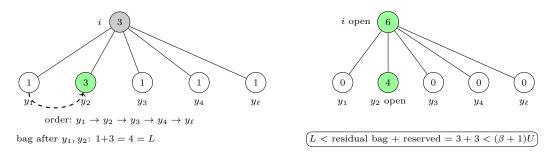


Figure 1: Toy examples for one star with L=4, U=5 and $\beta=1$; the numbers in solid (grey) circle represents the reserved clients at star center i. Numbers in circles at the spokes represents unsettled unreserved clients. A bag of unsettled clients is accumulated while scanning spokes in decreasing distance from the center i. Whenever the bag first reaches L, the current spoke opens (except for last spoke). At the end, the center i opens for all the remaining clients at last spoke along with the reserved clients.

5 High Level Idea of Our Algorithm

Our approach begins by decomposing the EL Clustering instance into two simpler subproblems: (i) A lower-bounded k-median instance I_L , obtained by dropping the upper bounds, and (ii) An upper-bounded k-median instance I_U , obtained by dropping the lower bounds. We solve (i) and (ii) using any available algorithms (approximate or heuristic), obtaining solutions S_L and S_U respectively.

Before presenting the details of the combination algorithm, we introduce the following key terms that will be used throughout the algorithmic descriptions.

Star: A subgraph formed by a facility i (the *star-center*) in S_L and a set of facilities from the S_U (its *spokes*) for which i is the closest facility among all facilities in the lower-bounded solution.

Open facility: A facility selected in the final EL Clustering solution to serve clients.

Closed facility: A facility not selected in the final EL Clustering solution.

We now describe our combination algorithm.

Grouping facilities into Stars: We group facilities into stars. For each facility i opened in the solution S_L , we form a star S_i consisting of i (the star-center) and all facilities i' opened in the solution S_U whose closest S_L facility is i.

Processing stars: Opening and Closing Facilities within a Star: Within each star S_i , we consider its spokes in decreasing order of distance from i. We maintain a "bag" collecting unsettled clients from S_U assignments; when the bag size reaches the lower bound L, we open the current spoke facility (except for the last spoke) and assign all clients in the bag to it. For all the clients accumulated at the last spoke, we open i. To achieve claimed lower and upper bounds exactly enough (max $\{0, L - |clients \ at \ last \ spoke|\}$), clients are reserved before processing the star and assigned to the facility i when it is opened. Refer Figure 1 for a simple toy example.

Determine Processing Order: The naive approach of processing stars in arbitrary order can lead to client conflicts, where clients intended for one star-center in S_L have already been assigned to spokes in another star. To avoid this, we build a *dependency graph* whose nodes are stars and whose edges represent shared clients between center of one star and the spokes of another star. We transform the built graph into an *almost-DAG* (removing cycles except self-loops) via a careful reassignment step. Processing stars in topological order ensures no premature client assignments.

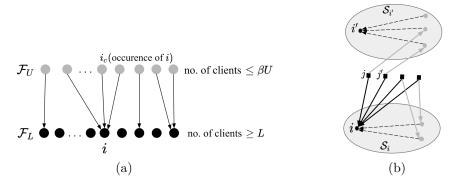


Figure 2: (a) Graph G_1 : I_L is an instance of $\mathbf{L}k\mathbf{M}$ and I_U that of $\mathbf{U}k\mathbf{M}$. (b) Let L=4. Black and grey edges show the assignment of clients in S_L and S_U , respectively. Star $S_{i'}$ is processed before star S_i . Clients j, j' assigned to i in S_L have already been assigned to facilities in $\eta^{-1}(i')$ and hence are not available while processing S_i .

6 Our Algorithm for EL Clustering

Let I = (P, c, C, F, k, U, L) be an instance of **EL** Clustering. We first create an instance I_L of **L**k**M** from I by dropping the upper bounds and then an instance I_U of **U**k**M** by dropping the lower bounds from I. Let $S_L = (\mathcal{F}_L, \sigma_L)$ and $S_U = (\mathcal{F}_U, \sigma_U)$ be solutions to I_L and I_U , respectively. Let β denote the violation in upper bounds, if any, in S_U . In the next section, we combine solutions S_L and S_U to obtain a solution $S_I = (\mathcal{F}_I, \sigma_I)$ to I with $(\beta + 1)$ factor violation in upper bounds.

6.1 Combining solutions S_L and S_U to obtain S_I

To obtain a solution $S_I = (\mathcal{F}_I, \sigma_I)$ to I, we will open some facilities in $\mathcal{F}_L \cup \mathcal{F}_U$. We construct a directed graph G_1 on the set of facilities in $\mathcal{F}_L \cup \mathcal{F}_U$. For a facility $i \in \mathcal{F}_U$, let $\eta(i)$ denote the facility in \mathcal{F}_L nearest to i (assuming that the distances are distinct). Add an edge $(i, \eta(i))$ in the graph. Note that a facility i may be open in both S_L and S_U , in that case $i \in \eta^{-1}(i)$. In order to avoid self loops, when $i = \eta(i)$, we denote the occurrence of i in \mathcal{F}_U by i_c so that $\eta(i_c) = i$. Thus, we obtain a forest of trees where-in each tree is a star. Formally, we define a star \mathcal{S}_i to be a collection of nodes in $\{i\} \cup \eta^{-1}(i)$ with $i \in \mathcal{F}_L$ as the star-center and $\eta^{-1}(i) \subseteq \mathcal{F}_U$. See Figure 2-(a).

We process the stars to decide the set of facilities to open in $\mathcal{F}_L \cup \mathcal{F}_U$. Consider a star \mathcal{S}_i centered at facility i. Clearly, the total assignments on i in S_L satisfy the lower bound but may violate the upper bound arbitrarily. On the other hand, the total assignments on a facility $i' \in \eta^{-1}(i)$ in S_U satisfy the upper bound (within β factor) but may violate the lower bound arbitrarily. We close some facilities in $\eta^{-1}(i)$ by transferring their clients (in S_U) to other facilities in $\eta^{-1}(i)$ if possible (or to i, if required) and open those at which the lower bound is satisfied. We may also have to open i in the process. We make sure that upper bound is violated within the claimed bounds and the total number of facilities opened in S_i is at most $|\eta^{-1}(i)|$. The cardinality constraint is, hence, satisfied.

Suppose we consider the facilities in $\eta^{-1}(i)$ in the order of decreasing distance from i. Let the order be $y_1, y_2, ..., y_l$. We wish to collect the clients assigned to them, by S_U , in a bag looking for a facility t at which we would have collected at least L clients so that we can open t, empty the bag by assigning all the clients in the bag to t and start the process again with the next facility in the order. The problem occurs when at the last facility (y_l) , in the order, the bag has less than L clients. In this case, we would like to assign these clients to the star-center t making use of the fact that t was assigned at least t clients in t. The problem here is that the clients assigned to t in t might have been assigned to the facilities in t

 S_i , then the clients assigned to i in S_L are not assigned to facilities in $\eta^{-1}(i')$ in S_U for a star $S_{i'}$ processed earlier. For this, we construct a weighted directed (dependency) graph G_2 (possibly with directed cycles) on stars and convert it into a directed acyclic graph (DAG) (except possibly for self-loops), before processing the stars. A topological ordering in the graph, then gives us the order in which the stars must be processed. We will denote the graph by $G_2(\sigma_L, \sigma_U)$ to show that it is a function of the assignments in S_L and S_U .

The graph $G_2(\sigma_L, \sigma_U)$ has the stars $\{S_i : |\eta^{-1}(i)| > 0\}$ as the vertices. Let $\mathcal{X}(i_1, i_2) = \{j \in \mathcal{C} : \sigma_U(j) = i' \in \eta^{-1}(i_2) \text{ and } \sigma_L(j) = i_1\}$ i.e., $\mathcal{X}(i_1, i_2)$ is the set of clients that are served by i_1 in S_L and by some facility at the spoke of the star centered at i_2 in S_U . We include the directed edge (S_{i_1}, S_{i_2}) from star S_{i_1} to S_{i_2} if $|\mathcal{X}(i_1, i_2)| > 0$. Let $w(S_{i_1}, S_{i_2}) = |\mathcal{X}(i_1, i_2)|$ denote the weight on the edge (S_{i_1}, S_{i_2}) . Refer to Algorithm 1 and Figure 3-(a) - (c) for the construction of graph G_2 . Initially, $\mathcal{X}(i_1, i_2) = \emptyset$ and $w(S_{i_1}, S_{i_2}) = 0$ for all pairs of stars S_{i_1} and S_{i_2} (i_1 may be same as i_2). If the resulting graph has no directed cycle except possibly the self-loops, we are done. The graph G_2 is an almost-DAG. A directed graph is called an almost-DAG, if the only cycles in it are self loops. However, if there are non-trivial directed cycles in the graph, we redefine the assignments in S_L to obtain another solution $\hat{S}_L = (\mathcal{F}_L, \hat{\sigma}_L)$ to break the cycles. The dependency graph for $(\hat{\sigma}_L, \sigma_U)$ will then be an almost-DAG.

```
Algorithm 1: Constructing Graph G_2(\sigma_L, \sigma_U)
```

```
Input: Stars S_i: i \in \mathcal{F}_L
Output: Weighted Directed Graph G_2(\sigma_L, \sigma_U) = (V, E)

1 V \leftarrow \{S_i: |\eta^{-1}(i)| > 0\}, E \leftarrow \emptyset
2 for j \in \mathcal{C} do
3 |i' \leftarrow \sigma_U(j), i_1 \leftarrow \sigma_L(j), i_2 \leftarrow \eta(i')
4 \mathcal{X}(i_1, i_2) \leftarrow \mathcal{X}(i_1, i_2) \cup \{j\}
5 |w(S_{i_1}, S_{i_2}) \leftarrow w(S_{i_1}, S_{i_2}) + 1
6 for each pair of stars S_{i_1}, S_{i_2} do
7 |if w(S_{i_1}, S_{i_2}) > 0 then
8 |E \leftarrow E \cup (S_{i_1}, S_{i_2})|
```

Algorithm 2: Breaking Cycles: Constructing an almost-DAG $G_2(\hat{\sigma}_L, \sigma_U)$

```
Input: Graph G_2(\sigma_L, \sigma_U)
     Output: G_2(\hat{\sigma}_L, \sigma_U)
 1 \hat{\sigma}_L(j) \leftarrow \sigma_L(j) \ \forall j \in \mathcal{C}
     while \exists a \text{ directed cycle} < S_{i_1}, S_{i_2}, \dots, S_{i_q} > (q > 1) \text{ in } G_2 \text{ do}
            \kappa \leftarrow w(\mathcal{S}_{i_1}, \mathcal{S}_{i_2}) // assume (\mathcal{S}_{i_1}, \mathcal{S}_{i_2}) as the minimum weight edge in the cycle
            for r = 1 to q do
 4
                    count \leftarrow 0, s \leftarrow (r \mod q) + 1
  5
                    for j \in \mathcal{X}(i_r, i_s) do
  6
                           if count < \kappa then
  7
                                  \hat{\sigma}_L(j) \leftarrow i_s
  8
                                  count + +
  9
                                  \mathcal{X}(i_r, i_s) \leftarrow \mathcal{X}(i_r, i_s) \setminus \{j\}
10
                                 \mathcal{X}(i_s, i_s) \leftarrow \mathcal{X}(i_s, i_s) \cup \{j\}
11
                    w(\mathcal{S}_{i_s}, \mathcal{S}_{i_s}) \leftarrow w(\mathcal{S}_{i_s}, \mathcal{S}_{i_s}) + \kappa
12
13
                    w(\mathcal{S}_{i_r}, \mathcal{S}_{i_s}) \leftarrow w(\mathcal{S}_{i_r}, \mathcal{S}_{i_s}) - \kappa
                    if w(S_{i_r}, S_{i_s}) = 0 then
14
                           E \leftarrow E \setminus (\mathcal{S}_{i_r}, \mathcal{S}_{i_s}) // Remove edge (\mathcal{S}_{i_r}, \mathcal{S}_{i_s}) from G_2
15
                    E \leftarrow E \cup (\mathcal{S}_{i_s}, \mathcal{S}_{i_s}) // Add edge (\mathcal{S}_{i_s}, \mathcal{S}_{i_s}) in G_2
```

Breaking the cycles: For graph $G_2(\sigma_L, \sigma_U)$, let $SC = \langle S_{i_1}, S_{i_2}, \dots, S_{i_q} \rangle$ be a non-trivial directed cycle with q > 1. Without loss of generality, let (S_{i_1}, S_{i_2}) be the minimum weight edge in the cycle. We reassign any $\kappa = w(S_{i_1}, S_{i_2})$ clients in $\mathcal{X}(i_r, i_{(r \mod q)+1})$ from i_r to $i_{(r \mod q)+1}$, increment the weight of the edge

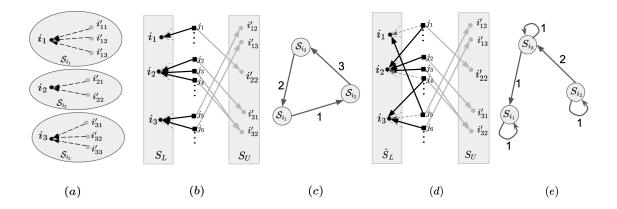


Figure 3: (a) Stars S_{i_1} , S_{i_2} and S_{i_3} ; (b) $\mathcal{X}(i_1, i_2) = \{j_1\}$, $\mathcal{X}(i_2, i_3) = \{j_2, j_3, j_4\}$, $\mathcal{X}(i_3, i_1) = \{j_5, j_6\}$; (c) Its directed cycle $G_2(\sigma_L, \sigma_U)$; (d) Breaking a cycle: assign j_1 to i_2 , j_4 to i_3 and j_5 to i_1 , that is, $\hat{\sigma}_L(j_1) = i_2$, $\hat{\sigma}_L(j_4) = i_3$ and, $\hat{\sigma}_L(j_5) = i_1$; (e) The sub-graph $G_2(\hat{\sigma}_L, \sigma_U)$ after breaking the cycle.

 $w(\mathcal{S}_{i(r \mod q)+1}, \mathcal{S}_{i(r \mod q)+1})$ by κ and, reduce the weight of the edge $w(\mathcal{S}_{i_r}, \mathcal{S}_{i(r \mod q)+1})$ by κ for $r=1\dots q$. Note that this adds new self-loops in the graph; however, no new non-trivial edge is added. Also, observe that $|\hat{\sigma}_L^{-1}(i)| = |\sigma_L^{-1}(i)|$ and hence $|\hat{\sigma}_L^{-1}(i)| \geq L$ is maintained for all $i \in \mathcal{F}_L$ after the reassignments. The weight of the edge $(\mathcal{S}_{i_1}, \mathcal{S}_{i_2})$ becomes zero and we remove it, thereby breaking the cycle. See Algorithm 2 and Figure 3-(d)-(e). Note that a client j gets reassigned at most once in all the cycles as during re-assignment, it moves its contribution from a non-trivial edge to a self-loop and not to any other non-trivial edge. Next, we bound the cost of solution \hat{S}_L in the Lemma 6.1.

Lemma 6.1. The cost, $Cost(\hat{S}_L)$, of solution \hat{S}_L is bounded by $Cost(S_L) + 2Cost(S_U)$.

Proof. Let $j \in \mathcal{C}$. The cost paid by j in solution \hat{S}_L is (see Figure 4-(a).): $c(j, \hat{\sigma}_L(j)) \leq c(j, \sigma_U(j)) + c(\sigma_U(j), \hat{\sigma}_L(j)) \leq c(j, \sigma_U(j)) + c(\sigma_U(j), \sigma_L(j)) \leq c(j, \sigma_U(j)) + (c(\sigma_U(j), j) + c(j, \sigma_L(j))) = c(j, \sigma_L(j)) + 2c(j, \sigma_U(j))$, where the second inequality holds since $\eta(\sigma_U(j)) = \hat{\sigma}_L(j)$. Summing over all $j \in \mathcal{C}$, we get the desired claim.

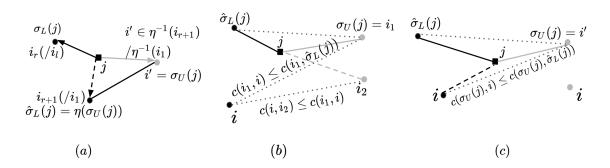


Figure 4: (a) $c(j, \hat{\sigma}_L(j)) \le c(j, \sigma_L(j)) + 2c(j, \sigma_U(j))$; (b) Cost bound of Type-I assignments; (c) Cost bound of Type-II assignments.

Graph $G_2(\hat{\sigma}_L, \sigma_U)$ has the following properties:

- 1. $G_2(\hat{\sigma}_L, \sigma_U)$ is an almost-DAG.
- 2. $|\hat{\sigma}_L^{-1}(i)| \ge L \ \forall \ i \in \mathcal{F}_L$.

Now that we have an almost-DAG on the stars, we process the stars in the sequence $\langle S_{i_1}, S_{i_2}, \dots S_{i_t} \rangle$ defined by a topological ordering of the vertices in $G_2(\hat{\sigma}_L, \sigma_U)$ (ignoring the self-loops). While processing the stars, we maintain partition of our clients into two sets, C_s and C_u of settled and unsettled clients respectively. We say that a client is settled if it has been assigned to an open facility in S_I and unsettled otherwise. Initially $C_s = \emptyset$ and $C_u = C$. As we process the stars, more and more clients get settled.

Algorithm 3: Process(S_i)

```
Input : S_i, i \in \mathcal{F}_L
 1 reserved(i) \leftarrow \emptyset, Baq \leftarrow \emptyset
 2 for i' \in \eta^{-1}(i) do
         N_{i'} \leftarrow \mathcal{C}_u \cap \sigma_U^{-1}(i')
 4 Arrange the facilities in \eta^{-1}(i) in the sequence \langle y_1, \dots y_l \rangle such that
       c(y_{l'}, i) \ge c(y_{l'+1}, i) \ \forall \ l' = 1 \dots l - 1
 5 if |N_{y_l}| < L then
          reserved(i) \leftarrow \text{set of any } L - |N_{y_l}| \text{ clients from } \hat{\sigma}_L^{-1}(i) \setminus N_{y_l}
          for i' \in \eta^{-1}(i) do
 7
               N_{i'} \leftarrow N_{i'} \setminus reserved(i)
 9 for l' = 1 \ to \ l - 1 \ do
          Bag \leftarrow Bag \cup N_{y_{i}}
10
          if |Bag| \ge L then
11
                Open facility y_{l'}
12
                for j \in Bag do
13
                     Assign j to y_{l'}, C_s \leftarrow C_s \cup \{j\}, C_u \leftarrow C_u \setminus \{j\}
14
15
17 if |Bag \cup N_{y_l} \cup reserved(i)| > (\beta + 1)U then
18
         t \leftarrow y_l
19 Open t
20 for j \in Bag \cup N_{y_l} \cup reserved(i) do
          Assign j to t, C_s \leftarrow C_s \cup \{j\}, C_u \leftarrow C_u \setminus \{j\}
```

To make sure that we do not open more than $|\eta^{-1}(i)|$ facilities in S_i , we open only one of i and y_l for the remaining $(Bag \cup N_{y_l} \cup reserved(i))$ clients. This also ensures that we do not open i more than once. We prefer to open i and give all the remaining clients to i because (as we will show later) the cost of assigning clients from $Bag \cup N_{y_l}$ to i is bounded whereas we do not know how to bound the cost of assigning clients in reserved(i) to y_l . However, in case, it leads to more than acceptable violation in the capacity at i, we open y_l and assign the remaining $(Bag \cup N_{y_l} \cup reserved(i))$ clients to it. We show that reserved(i) is empty in the latter case. Algorithm 4 summarizes our combination algorithm for constructing S_l from S_l and S_l .

6.2 Analysis

Recall that the assignments done in lines 9-15 are Type-I assignments. Let the assignment of clients to facility i when t=i in lines 20-21 be called as Type-II assignments and those to facility y_l when $t=y_l$ be

Algorithm 4: Constructing S_I

 $\overline{\textbf{Input}} : \langle S_L = \mathcal{F}_L, \sigma_L \rangle, \langle S_U = \mathcal{F}_U, \sigma_U \rangle$

Output: S_I

- 1 Construct graph $G_1 = \langle \mathcal{F}_L \cup \mathcal{F}_U, E \rangle$ where $E = \{(i', \eta(i')) : i' \in \mathcal{F}_U\}$.
- **2** Construct graph $G_2(\sigma_L, \sigma_U)$.
- **3** Construct an almost-DAG $G_2(\hat{\sigma}_L, \sigma_U)$ from $G_2(\sigma_L, \sigma_U)$ using Algorithm 2.
- 4 Obtain a topological ordering $\langle S_{i_1}, S_{i_2} \dots S_{i_t} \rangle$ of stars in the almost-DAG $G_2(\hat{\sigma}_L, \sigma_U)$.
- 5 for r = 1 to t do
- Process S_{i_r} using Algorithm 3

called as Type-III assignments. To prove our main theorem, we need to show that in the obtained solution S_I , the lower bounds are respected, the upper bounds are violated by a factor of at most $(\beta + 1)$, the cost of the solution is bounded and the running time is $O(k^3 + n)$. We first prove that the lower bounds are respected at the opened facilities in Lemma 6.2.

Lemma 6.2. Number of clients assigned to an open facility i in \mathcal{F} is at least L.

Proof. We will bound the lower bounds for all three type of assignments separately.

- 1. Observe that the facilities opened by Algorithm 3 in line 12 (Type-I assignment) satisfy the lower bounds by design of the algorithm.
- 2. In Type-II assignments, the star-center i satisfies the lower bound (if opened at line 19) as $|Bag \cup$ $N_{u_l} \cup reserved(i) | \geq L$ where the inequality follows because $|reserved(i)| = \max\{0, L - |N_{u_l}|\}$.
- 3. In Type-III assignments, facility y_l (if opened at line 19) also satisfies the lower bound as $|Baq \cup$ $N_{y_l} \cup reserved(i) | > (\beta + 1)U \ge 2L \text{ because } U \ge L \text{ and } \beta \ge 1.$

We next, show that the upper bounds are violated by a factor of at most $(\beta + 1)$ at the opened facilities in Lemma 6.3.

Lemma 6.3. Number of clients assigned to an open facility i in \mathcal{F} is no more than $(\beta + 1)U$.

Proof. We will bound the violations in the upper bounds for three type of assignments separately.

- 1. Consider the facilities in $\eta^{-1}(i)$. These facilities receive clients only in Type-I assignments (lines 13-14). Note that for l'=2,...,l-1, we have |Bag|< L just before line 10 and hence $|Bag|< L+\beta U$ (just after line 10) $\leq (1+\beta)U$ because $L \leq U$. For l'=1, |Bag|=0 just before line 10 and hence $|Baq| < \beta U$ (just after line 10).
- 2. For Type-II assignments, the bound holds trivially because the star-center i receives clients only when $|Bag| + |N_{y_l}| + |reserved(i)| \le (\beta + 1)U$.
- 3. The maximum number of clients received by facility y_l in Type-III assignments is, $|Bag| + |N_{y_l}| +$ $|reserved(i)| = |Bag| + |N_{y_l}| + max\{0, L - |N_{y_l}|\} = |Bag| + max\{L, |N_{y_l}|\} \le L + \beta U \le (\beta + 1)U.$

The next lemma (Lemma 6.4) bounds the cost of our solution (S_I) in terms of cost of solution S_U and S_L . **Lemma 6.4.** The cost of solution S_I is bounded by $7Cost(S_U) + 2Cost(S_L)$.

Proof. Consider a star S_i .

1. Type-I assignments: Let $j \in \mathcal{C}$ be assigned to a facility $i_2 \in \eta^{-1}(i)$ in our solution and to $i_1 \in \eta^{-1}(i)$ in S_U i.e., $i_1 = \sigma_U(j)$ and $i_2 = \sigma_I(j)$. The cost paid by j is (see Figure 4-(b)):

$$c(i_{2}, j) \leq c(i_{1}, j) + c(i_{1}, i) + c(i, i_{2})$$

$$\leq c(i_{1}, j) + 2c(i_{1}, i)$$

$$\leq c(i_{1}, j) + 2c(i_{1}, \hat{\sigma}_{L}(j)) \quad \text{(as } \eta(i_{1}) = i)$$

$$= 3c(i_{1}, j) + 2c(j, \hat{\sigma}_{L}(j)) \quad \text{(by triangle inequality)}$$

$$\leq 2c(j, \sigma_{L}(j)) + 7c(j, \sigma_{U}(j)) \quad \text{(by Lemma 6.1)}.$$

2. Type-II assignments: Let $j \in reserved(i)$ be assigned to i. Also, let $j \in N_{i'}$: $i' \in \eta^{-1}(i)$ be such that $i' = \sigma_U(j)$. Then, the cost (see Figure 4-(c)) is:

$$c(i, j) = c(\sigma_{U}(j), j) + c(\sigma_{U}(j), i)$$

$$= c(\sigma_{U}(j), j) + c(\sigma_{U}(j), \eta(\sigma_{U}(j)))$$

$$= c(\sigma_{U}(j), j) + c(\sigma_{U}(j), \hat{\sigma}_{L}(j))$$

$$\leq c(\sigma_{U}(j), j) + c(\sigma_{U}(j), j) + c(j, \hat{\sigma}_{L}(j))$$

$$= 2c(j, \sigma_{U}(j)) + c(j, \hat{\sigma}_{L}(j))$$

$$\leq 4c(j, \sigma_{U}(j)) + c(j, \sigma_{L}(j))$$

where the second and third equality follow because $\hat{\sigma}_L(j) = i = \eta(\sigma_U(j))$ and the last inequality follows by Lemma 6.1.

3. Type-III assignments: Note that $|Bag \cup N_{y_l} \cup reserved(i)| > (\beta + 1)U \Rightarrow |reserved(i)| = 0$, for otherwise $|N_{y_l} \cup reserved(i)| = L$ and thus $|Bag \cup N_{y_l} \cup reserved(i)| < L + L \le 2U$ because $L \le U$. Hence, the cost of assigning $|Bag \cup N_{y_l}|$ clients to y_l is bounded in the same manner as the cost of Type-I assignments.

By summing the cost over all the assignments of Type-I, Type-II and Type-III, we get, $Cost(S_I) \leq 7Cost(S_U) + 2Cost(S_L)$

We finally show bounds on the running time of our algorithm in the following lemma.

Lemma 6.5. Running time of our combination algorithm (Algorithm 4) is $O(k^3 + n)$.

Proof. Constructing G_1 takes $O(k^2)$ time and the graph G_2 can be constructed in time $O(n+k^2)$: for each client j, one can determine the edge (i_1, i_2) to which j contributes in constant time. G_2 can be converted into almost-DAG in $O(k^3+n)$ time using DFS and Algorithm 2: computing minimum weight edges takes at most $O(k^3)$ time over the entire algorithm and every client is re-assigned at most once. The time taken by Algorithm 3 when executed on all stars is no more than $O(n+k\log k)$; note that in this case also, a client is re-assigned at most once; $k\log k$ comes from sorting in step 4. Thus, having obtained solutions to LkM and UkM, combining the two solutions take $O(k^3+n)$ time.

Since any solution to I is feasible for I_L and I_U , we have $Cost(S_L) \leq Cost(S_I)$ and $Cost(S_U) \leq Cost(S_I)$. Therefore, the proof of Theorem 3.1 follows from Lemmas 6.2, 6.3, 6.4 and, 6.5. Furthermore, we have $Cost(O_L) \leq Cost(O)$ and $Cost(O_U) \leq Cost(O)$, where O, O, and O, denote optimal solution to I, I, and I, respectively. Therefore, the proof of Corollary 3.2 follows from Theorem 3.1 and by using approximation algorithms such as Byrka et al. (2016) for UkM to obtain S, and S and S and S and S and S are the proof of S are the proof of S and S are the proof of S and

7 Conclusion and Future Work

In this paper, we presented a modular approach for solving the **EL** Clustering problem by combining a solution of the k-median problem where the cluster sizes are lower bounded with another where the cluster sizes are upper bounded. Our solution introduces a bounded degradation over the costs of the given solutions.

Further, given a solution to the upper bounded instance where the upper bounds are violated by β our solution only incurs a bounded additional violation leading to at most a $\beta+1$ violation. An advantage of our method is that it gains from any improvements in the upper bounded and lower bounded solutions. Specifically, solutions for the upper and lower bounded instances with better approximation ratios enable us to obtain solutions with a better approximation ratio for **EL**. A similar note follows for solutions with smaller β violations in the upper bound. Interestingly, we note that Lemma 6.2 and Lemma 6.3 and hence our results hold for a more general scenario where the lower and upper bounds are not necessarily the same across the facilities, the only restriction is that $\max_{i \in \mathcal{F}} L_i \leq \min_{i \in \mathcal{F}} U_i$. Furthermore, we discussed how our algorithm can be applied to other clustering variants including k-means clustering. Moreover, for the special case when the gap between the lower and upper bounds is large enough (specifically, $2L_i \leq U_i, \forall i \in \mathcal{F}$) the violation in the upper bound can be reduced to $\beta + \epsilon$ for a given $\epsilon > 0$.

One direction for future work would be to get rid of the plus 1 violation in the upper bounds. Another interesting direction is to extend the results general lower and upper bounds. We acknowledge that the constants associated with the cost of generating a k-means clustering with equitable load are rather high in our paper. Improving these constants is another useful direction for future work.

References

- Ankit Aggarwal, Anand Louis, Manisha Bansal, Naveen Garg, Neelima Gupta, Shubham Gupta, and Surabhi Jain. A 3-approximation algorithm for the facility location problem with uniform capacities. *Journal of Math. Prog.*, 141(1-2):527–547, 2013.
- Gagan Aggarwal, Rina Panigrahy, Tomás Feder, Adam Meyerson, and Samir Khuller. Achieving anonymity via clustering. In *Proceedings of the 25th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS)*, pp. 153–162, 2006. doi: 10.1145/1142351.1142375.
- Saba Ahmadi, Pranjal Awasthi, Samir Khuller, Matthäus Kleindessner, Jamie Morgenstern, Pattara Sukprasert, and Ali Vakilian. Individual preference stability for clustering. arXiv preprint arXiv:2207.03600, 2022.
- Sara Ahmadian and Chaitanya Swamy. Improved approximation guarantees for lower-bounded facility location. In *Approximation and Online Algorithms*, pp. 257–271, 2013.
- Sara Ahmadian, Alessandro Epasto, Ravi Kumar, and Mohammad Mahdian. Clustering without over-representation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 267–275, 2019.
- M. Almanza, A. Epasto, A. Panconesi, and G. Re. k-Clustering with Fair Outliers. In *Proceedings of the* 15th ACM International Conference on Web Search and Data Mining (WSDM), pp. 5–15, 2022.
- Anna Arutyunova and Melanie Schmidt. Achieving Anonymity via Weak Lower Bound Constraints for k-Median and k-Means. In *STACS 2021*, pp. 7:1–7:17, 2021.
- Itai Ashlagi and Alvin E Roth. Kidney exchange: An operations perspective. *Management Science*, 67(9): 5455–5478, 2021.
- Pranjal Awasthi, Brian Brubach, Deeparnab Chakrabarty, John P Dickerson, Seyed A. Esmaeili, Matthäus Kleindessner, Marina Knittel, Jamie Morgenstern, Samira Samadi, Aravind Srinivasan, and Leonidas Tsepenekas. Fairness in clustering. 2022. URL https://www.fairclustering.com/.
- Suman Kalyan Bera, Deeparnab Chakrabarty, Nicolas Flores, and Maryam Negahbani. Fair algorithms for clustering. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 4955-4966, 2019. URL https://proceedings.neurips.cc/paper/2019/hash/fc192b0c0d270dbf41870a63a8c76c2f-Abstract.html.

- Ioana Oriana Bercea, Martin Groß, Samir Khuller, Aounon Kumar, Clemens Rösner, Daniel R. Schmidt, and Melanie Schmidt. On the cost of essentially fair clusterings. In Dimitris Achlioptas and László A. Végh (eds.), Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM 2019, September 20-22, 2019, Massachusetts Institute of Technology, Cambridge, MA, USA, volume 145 of LIPIcs, pp. 18:1–18:22. Schloss Dagstuhl Leibniz-Zentrum für Informatik, 2019. doi: 10.4230/LIPIcs.APPROX-RANDOM.2019.18. URL https://doi.org/10.4230/LIPIcs.APPROX-RANDOM.2019.18.
- Jaroslaw Byrka, Krzysztof Fleszar, Bartosz Rybicki, and Joachim Spoerhase. Bi-factor approximation algorithms for hard capacitated k-median problems. In SODA, 2015, pp. 722–736, 2015.
- Jaroslaw Byrka, Bartosz Rybicki, and Sumedha Uniyal. An approximation algorithm for uniform capacitated k-median problem with $(1 + \epsilon)$ capacity violation. In IPCO, pp. 262–274, 2016.
- Darshan Chakrabarti, John P Dickerson, Seyed A Esmaeili, Aravind Srinivasan, and Leonidas Tsepenekas. A new notion of individually fair clustering: α -equitable k-center. In *International conference on artificial intelligence and statistics*, pp. 6387–6408. PMLR, 2022.
- Debjani Chakraborty and Suman Das. Modified fuzzy c-mean for custom-sized clusters. $S\bar{a}dhan\bar{a}$, 44(8):182, 2019.
- Moses Charikar, Sudipto Guha, Éva Tardos, and David B. Shmoys. A constant-factor approximation algorithm for the k-median problem (extended abstract). In STOC, 1999, pp. 1–10, 1999.
- Xingyu Chen, Brandon Fain, Liang Lyu, and Kamesh Munagala. Proportionally fair clustering. In *International conference on machine learning*, pp. 1032–1041. PMLR, 2019.
- Anshuman Chhabra, Abhishek Roy, and Prasant Mohapatra. Suspicion-free adversarial attacks on clustering algorithms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 3625–3632, 2020.
- Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, and Sergei Vassilvitskii. Fair clustering through fairlets. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp. 5029-5037, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/978fce5bcc4eccc88ad48ce3914124a2-Abstract.html.
- Antonio Emanuele Cinà, Alessandro Torcinovich, and Marcello Pelillo. A black-box adversarial attack for poisoning clustering. *Pattern Recognition*, 122:108306, 2022.
- Michael B Cohen, Yin Tat Lee, and Zhao Song. Solving linear programs in the current matrix multiplication time. *Journal of the ACM (JACM)*, 68(1):1–39, 2021.
- Marek Cygan, Fedor V Fomin, Łukasz Kowalik, Daniel Lokshtanov, Dániel Marx, Marcin Pilipczuk, Michał Pilipczuk, and Saket Saurabh. *Parameterized algorithms*, volume 5. Springer, 2015.
- Rieke de Maeyer, Sami Sieranoja, and Pasi Fränti. Balanced k-means revisited. 2023.
- H. Gökalp Demirci and Shi Li. Constant approximation for capacitated k-median with $(1 + \epsilon)$ capacity violation. In ICALP, pp. 73:1–73:14, 2016.
- John Dickerson, Seyed Esmaeili, Jamie H Morgenstern, and Claire Jie Zhang. Doubly constrained fair clustering. Advances in Neural Information Processing Systems, 36, 2024.
- Hu Ding, Lunjia Hu, Lingxiao Huang, and Jian Li. Capacitated center problems with two-sided bounds and outliers. In WADS, pp. 325–336, 2017.
- Derya Dinler and Mustafa Kemal Tural. A survey of constrained clustering. In *Unsupervised learning* algorithms, pp. 207–235. Springer, 2016.

- Dheeru Dua and Casey Graff. Uci machine learning repository. http://archive.ics.uci.edu/ml, 2017.
- Seyed Esmaeili, Brian Brubach, Leonidas Tsepenekas, and John Dickerson. Probabilistic fair clustering. Advances in Neural Information Processing Systems, 33:12743–12755, 2020.
- Zachary Friggstad, Mohsen Rezapour, and Mohammad R. Salavatipour. Approximating Connected Facility Location with Lower and Upper Bounds via LP Rounding. In SWAT, volume 53, pp. 1:1–1:14, 2016.
- Nuwan Ganganath, Chi-Tsun Cheng, and K Tse Chi. Data clustering with cluster size constraints using a modified k-means algorithm. In 2014 International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery, pp. 158–161. IEEE, 2014.
- Dishant Goyal, Ragesh Jaiswal, and Amit Kumar. FPT approximation for constrained metric k-median/means. CoRR, abs/2007.11773, 2020. URL https://arxiv.org/abs/2007.11773.
- Sapna Grover, Neelima Gupta, Samir Khuller, and Aditya Pancholi. Constant factor approximation algorithm for uniform hard capacitated knapsack median problem. In FSTTCS, pp. 23:1–23:22, 2018.
- Yutian Guo, Junyu Huang, and Zhen Zhang. A constant factor approximation for lower-bounded k-median. In *TAMC*, pp. 119–131, 2020.
- Anupam Gupta, Katrina Ligett, Frank McSherry, Aaron Roth, and Kunal Talwar. Differentially private combinatorial optimization. In *Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms*, pp. 1106–1125. SIAM, 2010.
- Neelima Gupta, Sapna Grover, and Rajni Dabas. Respecting lower bounds in uniform lower and upper bounded facility location problem. In *COCOON*, pp. 463–475, 2021.
- Lu Han, Chunlin Hao, Chenchen Wu, and Zhenning Zhang. Approximation algorithms for the lower-bounded knapsack median problem. In AAIM, volume 12290, pp. 119–130, 2020a.
- Lu Han, Chunlin Hao, Chenchen Wu, and Zhenning Zhang. Approximation algorithms for the lower-bounded k-median and its generalizations. In *COCOON*, pp. 627–639, 2020b.
- Frank Höppner and Frank Klawonn. Clustering with size constraints. In Lakhmi C. Jain, Mika Sato-Ilic, Maria Virvou, George A. Tsihrintzis, Valentina Emilia Balas, and Canicious Abeynayake (eds.), Computational Intelligence Paradigms, Innovative Applications, volume 137 of Studies in Computational Intelligence, pp. 167–180. Springer, 2008.
- Sedjro Salomon Hotegni, Sepideh Mahabadi, and Ali Vakilian. Approximation algorithms for fair range clustering. In *International Conference on Machine Learning*, pp. 13270–13284. PMLR, 2023.
- Beth Hurst, Randall R Wallace, and Sarah B Nixon. The impact of social interaction on student learning. Reading horizons, 2013.
- Shunhua Jiang, Zhao Song, Omri Weinstein, and Hengjie Zhang. Faster dynamic matrix inverse for faster lps. arXiv preprint arXiv:2004.07470, 2020.
- Matthew Jones, Huy Nguyen, and Thy Nguyen. Fair k-centers via maximum matching. In *International conference on machine learning*, pp. 4940–4949. PMLR, 2020.
- Matthew Jones, Huy L Nguyen, and Thy D Nguyen. Differentially private clustering via maximum coverage. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11555–11563, 2021.
- Ashwini S Kadam, Shraddha R Nikam, Ankita A Aher, Gayatri V Shelke, and Amar S Chandgude. Prediction for loan approval using machine learning algorithm. *International Research Journal of Engineering and Technology (IRJET)*, 8(04), 2021.
- Haim Kaplan and Uri Stemmer. Differentially private k-means with constant multiplicative error. arXiv preprint arXiv:1804.08001, 2018.

- Samir Khuller and Yoram J. Sussmann. The capacitated k-center problem. SIAM Journal on Discrete Mathematics, pp. 403–418, 2000.
- Matthäus Kleindessner, Pranjal Awasthi, and Jamie Morgenstern. Fair k-center clustering for data summarization. In *International Conference on Machine Learning*, pp. 3448–3457. PMLR, 2019.
- Madhukar R. Korupolu, C. Greg Plaxton, and Rajmohan Rajaraman. Analysis of a local search heuristic for facility location problems. *Journal of Alg.*, 37(1):146–188, 2000. ISSN 0196-6774.
- Olha Kovalchuk, Mikolaj Karpinski, Serhiy Banakh, Mykhailo Kasianchuk, Ruslan Shevchuk, and Nataliya Zagorodna. Prediction machine learning models on propensity convicts to criminal recidivism. *Information*, 14(3):161, 2023.
- Kai Lei, Sibo Wang, Weiwei Song, and Qilin Li. Size-constrained clustering using an initial points selection method. In *Knowledge Science*, Engineering and Management: 6th International Conference, KSEM 2013, Dalian, China, August 10-12, 2013. Proceedings 6, pp. 195–205. Springer, 2013.
- Bo Li, Lijun Li, Ankang Sun, Chenhao Wang, and Yingfan Wang. Approximate group fairness for clustering. In *International conference on machine learning*, pp. 6381–6391. PMLR, 2021.
- Shanfei Li. An Improved Approximation Algorithm for the Hard Uniform Capacitated k-median Problem. In APPROX/RANDOM, pp. 325–338, 2014.
- Shi Li. On uniform capacitated k-median beyond the natural LP relaxation. In SODA, pp. 696–707, 2015.
- Shi Li. Approximating capacitated k-median with $(1 + \epsilon)k$ open facilities. In SODA, pp. 786–796, 2016.
- Weibo Lin, Zhu He, and Mingyu Xiao. Balanced clustering: A uniform model and fast algorithm. In *IJCAI*, pp. 2987–2993, 2019.
- Duncan McElfresh, Michael Curry, Tuomas Sandholm, and John Dickerson. Improving policy-constrained kidney exchange via pre-screening. *Advances in Neural Information Processing Systems*, 33:2674–2685, 2020.
- Sérgio Moro, Paulo Cortez, and Paulo Rita. A data-driven approach to predict the success of bank telemarketing. *Decision Support Systems*, 62:22–31, 2014. ISSN 0167-9236. doi: 10.1016/j.dss.2014.03.001. URL https://www.sciencedirect.com/science/article/pii/S016792361400061X.
- Clemens Rösner and Melanie Schmidt. Privacy preserving clustering with constraints. In *ICALP*, pp. 96:1–96:14, 2018.
- Mohammad Ahmad Sheikh, Amit Kumar Goel, and Tapas Kumar. An approach for prediction of loan approval using machine learning algorithm. In 2020 international conference on electronics and sustainable communication systems (ICESC), pp. 490–494. IEEE, 2020.
- Amy Soller. Supporting social interaction in an intelligent collaborative learning system. *International journal of artificial intelligence in education*, 12:40–62, 2001.
- Beata Strack, Jonathan P. Deshazo, Chris Gennings, Juan Luis Olmo Ortiz, Sebastian Ventura, Krzysztof J. Cios, and John N. Clore. Impact of hba1c measurement on hospital readmission rates: Analysis of 70,000 clinical database patient records. *BioMed Research International*, 2014:781670, 2014. doi: 10.1155/2014/781670. URL https://doi.org/10.1155/2014/781670. Article ID 781670.
- Wei Tang, Yang Yang, Lanling Zeng, and Yongzhao Zhan. Size constrained clustering with milp formulation. *IEEE Access*, 8:1587–1599, 2019.
- Guido Vittorio Travaini, Federico Pacchioni, Silvia Bellumore, Marta Bosia, and Francesco De Micco. Machine learning and criminal justice: A systematic review of advanced methodology for recidivism risk prediction. *International journal of environmental research and public health*, 19(17):10594, 2022.

Pravin M Vaidya. Speeding-up linear programming using fast matrix multiplication. In 30th annual symposium on foundations of computer science, pp. 332–337. IEEE Computer Society, 1989.

Jan van den Brand. A deterministic linear program solver in current matrix multiplication time. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 259–278. SIAM, 2020.

A Reducing violation in upper bounds when $2L \leq U$

```
Algorithm 5: Process(S_i)
    Input : S_i : i \in \mathcal{F}_L
 1 reserved(i) \leftarrow \emptyset, Bag \leftarrow \emptyset
 \mathbf{2} \ \mathbf{for} \ i' \in \eta^{-1}(i) \ \mathbf{do}
     N_{i'} \leftarrow \mathcal{C}_u \cap \sigma_U^{-1}(i')
 4 Arrange the facilities in \eta^{-1}(i) in the sequence \langle y_1, \dots y_l \rangle such that
      c(y_{l'}, i) \ge c(y_{l'+1}, i) \ \forall \ l' = 1 \dots l-1
 5 if |N_{y_l}| < L then
         reserved(i) \leftarrow \text{set of any } L - |N_{y_l}| \text{ clients from } \hat{\sigma}_L^{-1}(i) \setminus N_{y_l}
          for i' \in \eta^{-1}(i) do
 7
           N_{i'} \leftarrow N_{i'} \setminus reserved(i)
    Prev \leftarrow null, Prev_{count} = 0
10 for l' = 1 to l - 1 do
          Bag \leftarrow Bag \cup N_{y_{l'}}
11
          if |Bag| \ge L then
12
               Open facility y_{l'}
13
               Count \leftarrow 0
14
               for j \in Bag do
15
                    if Count < \beta U then
16
                         Assign j to y_{l'}, C_s \leftarrow C_s \cup \{j\}, C_u \leftarrow C_u \setminus \{j\}, Bag \leftarrow Bag \setminus \{j\}, Count + +
17
          else
18
               Prev \leftarrow y_{l'} // Prev denotes the last unopened facility in \eta^{-1}(i)
19
               Prev_{count} = |Bag|
20
21 if |Bag \cup N_{y_l} \cup reserved(i)| \leq (\beta + \epsilon)U then
          Open i
22
          for j \in Bag \cup N_{y_l} \cup reserved(i) do
23
           Assign j to i, \mathcal{C}_s \leftarrow \mathcal{C}_s \cup \{j\}, \mathcal{C}_u \leftarrow \mathcal{C}_u \setminus \{j\}
24
25
26 if |Bag \cup N_{y_l} \cup reserved(i)| \leq (\beta + \epsilon)U then
27
          for j \in Bag \cup N_{y_l} \cup reserved(i) do
28
              Assign j to y_l, C_s \leftarrow C_s \cup \{j\}, C_u \leftarrow C_u \setminus \{j\}
29
30
31 Open Prev and y_l // |reserved(i)| = 0 when |Bag \cup N_{y_l} \cup reserved(i)| > (\beta + \epsilon)U
32 Count \leftarrow 0
33 Bag \leftarrow Bag \cup N_{y_l} \cup reserved(i)
34 for j \in Bag do
          if Count \le L then
35
               Assign j to Prev, C_s \leftarrow C_s \cup \{j\}, C_u \leftarrow C_u \setminus \{j\}, Bag \leftarrow Bag \setminus \{j\}, Count + +
36
          else
37
               Assign all remaining clients in Bag to y_l and Break
38
```

In this section, assuming $2L \leq U$, we modify Algorithm 3 to obtain Algorithm 5 that reduces the violation in upper bounds from $(\beta + 1)$ to $(\beta + \epsilon)$ for a given $\epsilon > 0$. In particular, we present the following results:

Theorem A.1. For $2L \leq U$, given a solution S_U for UPPER BOUNDED k-MEDIAN (UkM) violating the upper bound by a factor of β and a solution S_L for LOWER BOUNDED k-MEDIAN (LkM). If the clustering costs of the solutions are $Cost(S_U)$ and $Cost(S_L)$, respectively. Then, a solution of cost at most $(O(\frac{1}{\epsilon})(7Cost(S_U) + 2Cost(S_L)))$ can be obtained for **EL** Clustering that violates the upper bound by a factor of $(\beta + \epsilon)$ for a fixed $\epsilon > 0$.

We do the following modifications to Algorithm 3: (i) on arriving at a facility, say t, at which $|Bag| \geq L$, we open t and instead of emptying the bag, we assign only βU clients to t. Remaining clients are carried forward to the next facility in the order; (ii) we keep account of the last facility (in Prev), if any, that is not opened, and the number of clients in the bag at that instant (in $Prev_{count}$) i.e., Prev is the facility $y_{l'}$ for which |Bag| < L immediately after line 14 (hence at line 23) and $Prev_{count} = |Bag|$ at that time. We open Prev at the end, if required. This is done as follows: if $|Bag \cup reserved(i) \cup N_{y_l}| \leq (\beta + \epsilon)U$, we are done (we open $i(/y_l)$) and assign all clients to it). Else, we open both Prev and y_l (at line 42) (note that $Prev \neq y_l$ must exist in this case) and, distribute the clients in Prev and Prev are done and Prev a

Let the assignment of clients to facility Prev in line 47 be called as Type-IV assignments. The assignments in line 20, line 31 and lines 38 & 49 are Type-II, Type-II and Type-III assignments respectively. Before we proceed to prove our claims, note that we open at most one of y_l and i: if i is opened at line 29, we return at line 33 and thus y_l is never opened in this case. As before, this ensures that i is not opened more than once.

Clearly, lower bound is satisfied by Type-I and Type-IV assignments done in line 20 and 47 for the facilities opened in lines 16 and 42 respectively. Also, since $|Bag \cup N_{y_l} \cup reserved(i)| \geq L$, lower bound is satisfied by Type-II assignments done in line 31 for the facility i opened in line 29. For Type-III assignments done at line 38, $|Bag \cup N_{y_l} \cup reserved(i)| > (\beta + \epsilon)U \geq (\beta + \epsilon)L$. Clearly, the upper bound is violated by a factor of at most $(\beta + \epsilon)$ at the facilities opened in lines 16, 29, 36 and Prev in line 42. For the assignments done in line 49, we look at the status at line 42: |reserved(i)| = 0, for otherwise $|N_{y_l} \cup reserved(i)| = L$, hence $|Bag \cup reserved(i) \cup N_{y_l}| < L + L \leq U$. Thus, $|Bag \cup N_{y_l} \cup reserved(i)| = |Bag \cup N_{y_l}| < L + \beta U$. Also, $|Bag \cup N_{y_l} \cup reserved(i)| > (\beta + \epsilon)U > 2L$. Thus, $L < |Bag \cup N_{y_l} \cup reserved(i)| - L < \beta U$ i.e., at line 49, $L < |Bag| < \beta U$.

Costs of Type-II and Type-III assignments are bounded in the same manner as in Section 6. To bound the service cost of Type-IV assignments (line 47), observe that $|Bag \cup reserved(i) \cup N_{y_l}| > (\beta + \epsilon)U \Rightarrow |Bag| > \epsilon U$ as |reserved(i)| = 0 and $|N_{y_l}| \leq \beta U$; hence, $Prev_{count} \geq |Bag|$ (at line 25) $> \epsilon U > \epsilon L$. Note that Prev and $Prev_{count}$ do not change after exiting the for-loop at line 27. Thus, $Prev_{count} > \epsilon L$ after line 42 also and the cost of assigning at most L clients from N_{y_l} to Prev is bounded by $(1/\epsilon)$ times the cost of assignments is bounded by $(1/\epsilon)$ total cost of Type-III assignments.

B Modifications for Other related Problems

B.1 *k*-Means with Equitable Load

The k-means problem with **EL** constraints is same as the k-median problem with **EL** constraints except that the goal now is to minimize the sum of the squared distances instead of minimizing the sum of distances from the assigned facilities. Further, the facilities to be selected in the k-means problem possibly belong to an infinite space. Note that, in the k-means problem,

- 1. the distances are squared which may not satisfy triangle inequality but they satisfy α -relaxed triangle inequality, that is, $c(x, y) \leq \alpha c(x, z) + \alpha c(z, y)$ for $\alpha = 2$ and,
- 2. we can assume the set of facilities to be in a finite space by losing 2α factor in the distances for $\alpha = 2$.

We create an instance I_L of LOWER BOUNDED k-MEANS and I_U of UPPER BOUNDED k-MEANS instead of $\mathbf{L}k\mathbf{M}$ and $\mathbf{U}k\mathbf{M}$. Solution S_I is obtained by using the combination algorithm on S_U and S_L . Note that, the violation in upper bounds remains the same, that is, for β violation in upper bounds in S_U , we get $(\beta + 1)$ violation in upper bounds in S_I . We next bound the cost of the obtained solution S_I . With relaxed triangle inequality, Lemma 6.1 can be modified to bound the cost, $Cost(\hat{S}_L)$, of solution \hat{S}_L by $4Cost(S_L) + 6Cost(S_U)$.

Lemma B.1. The cost of solution S_I is bounded by $352Cost(S_U) + 192Cost(S_L)$.

Proof. We will modify the proof Lemma 6.4 to accommodate relaxed triangle equality. Due to space constraints, We will give details of Type-I assignments which have the dominating cost. Cost of Type-II can be bounded by $30c(j, \sigma_U(j)) + 16c(j, \sigma_L(j))$ in a similar manner. Cost of Type-III assignments is same as Type-I assignments.

Type-I assignments: Consider a star S_i . Let $j \in \mathcal{C}$ be assigned to a facility $i_2 \in \eta^{-1}(i)$ and to $i_1 \in \eta^{-1}(i)$ in S_U i.e., $i_1 = \sigma_U(j)$ and $i_2 = \sigma_I(j)$. The cost paid by j is: $c(i_2, j) \leq \alpha \cdot c(i_2, i) + \alpha \cdot c(i, j) \leq \alpha \cdot c(i_2, i) + \alpha^2 \cdot (c(i, i_1) + c(i_i, j)) \leq \alpha^2 \cdot c(i_1, j) + (\alpha + \alpha^2)c(i_1, i) \leq \alpha^2 \cdot c(i_1, j) + (\alpha + \alpha^2)c(i_1, \hat{\sigma}_L(j)) \leq \alpha^2 \cdot c(i_1, j) + (\alpha^2 + \alpha^3) \cdot (c(i_1, j) + c(j, \hat{\sigma}_L(j))) = (2\alpha^2 + \alpha^3) \cdot c(i_1, j) + (\alpha^2 + \alpha^3) \cdot c(j, \hat{\sigma}_L(j)) = 16c(i_1, j) + 12c(j, \hat{\sigma}_L(j)) \leq 88c(j, \sigma_U(j)) + 48c(j, \sigma_L(j))$, where the first, second, fourth inequality follow by relaxed triangle inequality, third inequality follows as $\eta(i_1) = i$, the last equality follows by setting value of α to 2 and the last inequality follows by bound on $Cost(\hat{S}_L)$.

We incur an additional multiplicative factor of 2α due to the assumption that the points lie in a finite space. Multiplying by 4 for $\alpha = 2$, we get, $Cost(S_I) < 352Cost(S_{II}) + 192Cost(S_{II})$.

B.2 k-Center with Equitable Load

The k-Center problem with \mathbf{EL} constraints is the same as the k-median with \mathbf{EL} constraints except that the goal now is to minimize the maximum distance of a client from the assigned facility instead of minimizing the total distance. We create instance I_L and I_U of Lower Bounded k-Center and Upper Bounded k-Center respectively instead of $\mathbf{L}k\mathbf{M}$ and $\mathbf{U}k\mathbf{M}$. Same bounds are obtained on the cost by taking the maximum of the cost of all the types of assignments. Bounds on violation in upper bounds remains the same.

B.3 k-Facility Location with Equitable Load

The k-Facility location with **EL** constraints is a generalization of the k-median with **EL** constraints where for every facility $i \in \mathcal{F}$, we also have a facility opening cost f_i . The objective now is to identify $\mathcal{F}' \subseteq \mathcal{F}$ of size at most k and an assignment σ of clients to \mathcal{F}' so as to minimize the sum of the distances of the clients from their assigned facilities plus the facility opening costs of the selected facilities. We create instance I_L of Lower Bounded k-Facility Location by dropping the upper bounds and cardinality constraint. An instance I_U of UPPER Bounded k-Facility Location is created by dropping the lower bounds. We then follow the same procedure as described for k-median in Section 6.1 to combine the solutions of the two instances. Cost of assignment is bounded in the same manner. There is no loss in factor due to facility opening costs as we only open facilities in $(\mathcal{F}_L \cup \mathcal{F}_U)$. The violation in the upper bounds remains the same.

B.4 Knapsack-Median with Equitable Load

KNAPSACK MEDIAN with **EL** constraints is another generalization of k-median with **EL** constraints where every facility i has weight f_i and instead of k, and we have a budget B on the total weight. Therefore, the objective is to identify $\mathcal{F}' \subseteq \mathcal{F}$ and an assignment σ of clients to \mathcal{F}' so as to minimize the sum of the distances of the clients from the assigned facility subject to the constraint $\sum_{i \in \mathcal{F}'} f_i \leq B$.

We first create an instance I_U of UPPER BOUNDED KNAPSACK MEDIAN by dropping the lower bounds and instance I_L of LOWER BOUNDED KNAPSACK MEDIAN from I by dropping the upper bounds, reducing the set of facilities to \mathcal{F}_U and setting budget to the budget of S_U (note that this can be different from given budget B if there is violation in budget in S_U , otherwise it is B only). It can be shown that $Cost(O_L) \leq (2 + Cost(S_U))Cost(O)$: if a client j is assigned in the optimal solution O to I, to a facility i not in \mathcal{F}_U , we assign it to a facility i', nearest to i, in \mathcal{F}_U . The cost $c(j, i') \leq c(j, i) + c(i, i') \leq c(j, i) + c(i, i'') \leq c(j, i''') \leq c(j, i'') \leq c(j, i'') \leq c(j, i''') \leq c(j, i''') \leq c(j, i$

c(j, i) + c(j, i) + c(j, i'') = 2c(j, i) + c(j, i'') where $i'' \in \mathcal{F}_U : \sigma_U(j) = i''$ and the second inequality holds because i' is nearest to i and not i''.

We next use the same procedure as in Section 6.1 to combine solutions S_L and S_U of instances I_L and I_U respectively. Note that since $\mathcal{F}_L \subseteq \mathcal{F}_U$, for $i \in \mathcal{F}_L$, $y_l = i$ in the star S_i . This is important to make sure that the total facility opening cost in our solution is no more than that of S_U in case we open i.

B.5	A summary of	approximation	results using	current best	S_L and	S_{II}
D .0	7 Committee of	approximation	results using	carrent best	\mathcal{L}_L and	^

Problem	Results for underlying problems		Our results	Previous results
1 Toblem	(α_L)	(β, α_U)	$(heta_U, heta_{approx})$	
	$_{ m LFL}$	UkM		
k-Median	(82.6)	$((1+\epsilon), O(1/\epsilon^2))$	$((2+\epsilon), O(1/\epsilon^2))$	Nil
	Ahmadian & Swamy (2013)	Byrka et al. (2016)		
	$\mathrm{L}k\mathrm{C}$	UkC		
k-Center	(2)	(1,6)	(2,46)	Nil
	Aggarwal et al. (2006)	Khuller & Sussmann (2000)		
	$_{ m LFL}$	UFL		$(O(1), O(1))^{\#}$ Friggstad et al. (2016)
Facility Location	(82.6)	(1,3)	(2, 186.2)	
	Ahmadian & Swamy (2013)	Aggarwal et al. (2013)		
	$_{ m LFL}$	UkFL		
k-Facility Location	(82.6)	$((2+\epsilon), O(1/\epsilon^2))$	$((3+\epsilon), O(1/\epsilon^2))$	Nil
	Ahmadian & Swamy (2013)	Grover et al. (2018)		
	LKnM	UKnM		
Knapsack Median*	(1608)	$((2+\epsilon), O(1/\epsilon^2))$	$((3+\epsilon), O(1/\epsilon^2))$	Nil
	Han et al. (2020a)	Grover et al. (2018)		

Table 1: Our results. α_L and α_U indicate the approximation factor for the lower and upper bounded variant respectively of the problem under consideration. Also, θ_U and θ_{approx} indicate the violation in the upper bounds and the approximation factor respectively for the corresponding lower and upper bounded problem. 'Nil' represents that there is no result known for the problem. (*): the result violates the budget by a factor of $(1 + \epsilon)$. (#): the result violates the lower bounds also. Note that, to the best of our knowledge there is no known polynomial time approximation algorithm for upper bounded k-means problem.

C Experiments

We empirically evaluate our combination algorithm for the k-Median objective on benchmark datasets from the UCI Machine Learning Repository, as preprocessed in Almanza et al. (2022): **Adult, Diabetes**, and **Bank**. The **Adult** dataset Dua & Graff (2017) contains census records, the **Bank** dataset Moro et al. (2014) contains data points from a marketing campaign, and the **Diabetes** dataset Strack et al. (2014) contains admission records of diabetic patients. The input points are embedded in Euclidean space, and distances are measured using the ℓ_2 metric. The algorithm is implemented in Python and executed on a Apple laptop with an Apple M1 processor, 8 GB RAM, running macOS (Tahoe 26 beta). All source code, logs, and charts are publicly available⁸.

Input solutions. The initial lower-bounded (S_L) and upper-bounded (S_U) solutions were obtained using greedy heuristics that respect the lower and upper bounds respectively; thus, $\beta = 1$. Since our combination algorithm is modular, it can operate on any feasible pair (S_L, S_U) and hence our focus is to evaluate the performance of the *combination step itself*, independent of the specific heuristics used to generate S_L and S_U .

Experimental design. We vary the number of clusters k for $k \in \{10, 25, 50, 75, 100\}$. The lower bound L and upper bound U are derived from the natural cluster size n/k, with buffers of $\pm \{10\%, 25\%, 50\%\}$. Facility pool \mathcal{F} of varying sizes ($\{5\%, 10\%, 20\%, 25\%\}$) of the input size) are selected randomly from the input points.

⁸https://github.com/0-rudra-0/el-clustering

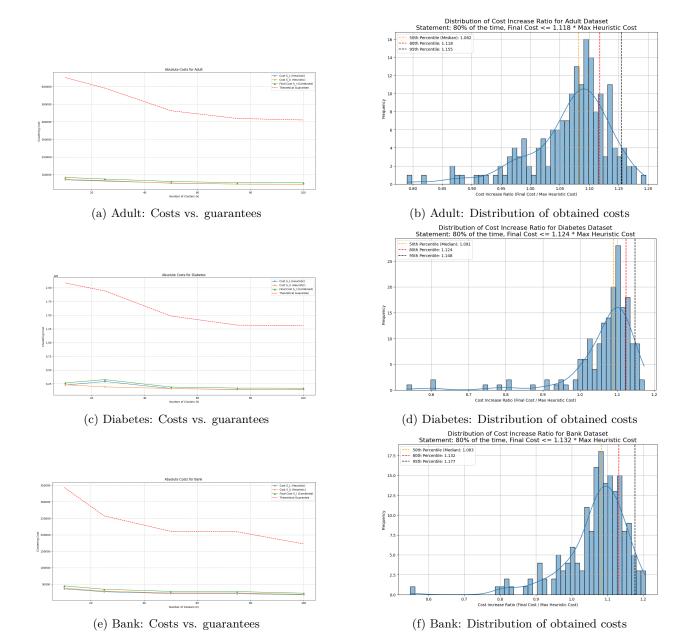


Figure 5: Left column (a, c, e) shows the plot of costs of S_L , S_U , our combined solution, and the theoretical guarantee $(7Cost(S_U) + 2Cost(S_L))$ against the values of k. For our combined solution, the y-axis reports the cost corresponding to the worst observed gap between the cost of our solution and $\max\{Cost(S_L), Cost(S_U)\}$. Right column (b, d, f): distribution of the ratio $\frac{Cost(\text{combined})}{\max\{Cost(S_L), Cost(S_U)\}}$ across all runs.

C.1 Results and Insights

Cost Analysis. Figure 5 reports cost performance across all datasets. The left column (Figure 5 (a, c, e)) shows results as k varies. For each value of k, we run the algorithm with different (L, U) settings and varying the size of the facility pool. The y-axis reports the cost corresponding to the worst observed gap between the cost of our solution and $\max\{Cost(S_L), Cost(S_U)\}$. We observe that, across all datasets, our algorithm is close to the maximum of the input costs $Cost(S_L)$ and $Cost(S_U)$.

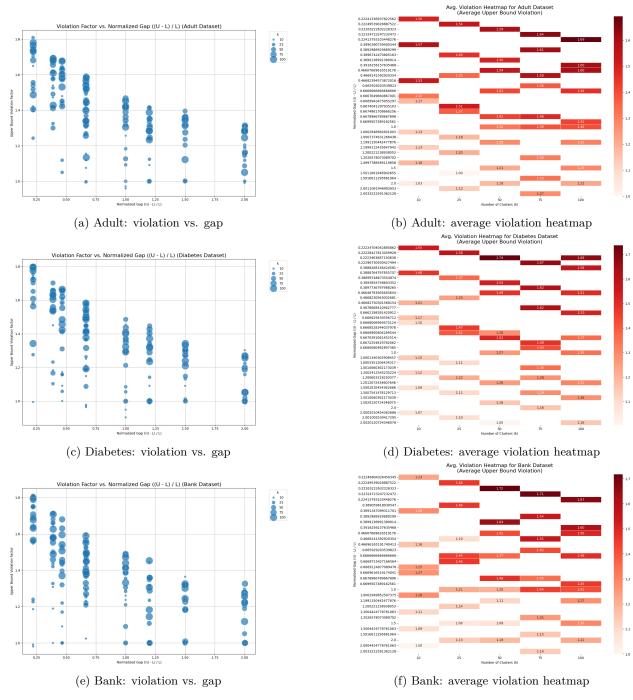


Figure 6: Left column (a, c, e): upper bound violation factor as a function of the (U/L). Right column (b, d, f): heatmap showing the combined impact of k and (L, U) setting on the violation, on the cost averaged over varying sizes of the facility pool.

The right column (Figure 5 (b, d, f)) shows the full distribution of the ratio $\frac{Cost(combined)}{\max\{Cost(S_L),Cost(S_U)\}}$ across all runs. The distribution is sharply concentrated near 1, confirming that the cost overhead is minimal in practice. In particular, the combined cost is within 1.2 of the cost of maximum of the two input solutions and it is typically 1.1.

Thus, the cost of our solution is substantially better than the worst-case theoretical bound.

Upper Bound Violation. Figure 6 analyzes the violation factor of the upper bound constraint. The left column (Figure 6 (a, c, e)) plots the violation factor as a function of the U/L and k for all sizes of the facility pool. It can be observed, the violation is within 1.54 in $\approx 80\%$ of the cases. The right column (Figure 6 (b, d, f)) presents a heatmap that shows the combined impact of k and (L, U) setting on the violation, on the cost averaged over varying sizes of the facility pool. It can be observed that the average violation over all the runs never exceeds 1.37 in all the data sets.

Further, we observed the trend that the performance degrades with higher values of k and also with tight gap between L and U. Since k is typically small, violations in practice are expected to be small. The latter observation aligns with our theoretical analysis presented in Appendix A. Thus, in practice, violations remain well below the theoretical bound of $(\beta + 1)$ in most of the cases.