# Constant Approximation for $k$-Median and $k$-Means with Outliers via Iterative Rounding

Ravishankar Krishnaswamy[*]     Shi Li[†]     Sai Sandeep[‡]

April 9, 2018

### Abstract

In this paper, we present a new iterative rounding framework for many clustering problems. Using this, we obtain an $(\alpha_1 + \epsilon \le 7.081 + \epsilon)$-approximation algorithm for $k$-median with outliers, greatly improving upon the large implicit constant approximation ratio of Chen [16]. For $k$-means with outliers, we give an $(\alpha_2 + \epsilon \le 53.002 + \epsilon)$-approximation, which is the first $O(1)$-approximation for this problem. The iterative algorithm framework is very versatile; we show how it can be used to give $\alpha_1$- and $(\alpha_1 + \epsilon)$-approximation algorithms for matroid and knapsack median problems respectively, improving upon the previous best approximations ratios of 8 [42] and 17.46 [9].

The natural LP relaxation for the $k$-median/$k$-means with outliers problem has an unbounded integrality gap. In spite of this negative result, our iterative rounding framework shows that we can round an LP solution to an *almost-integral* solution of small cost, in which we have at most two fractionally open facilities. Thus, the LP integrality gap arises due to the gap between almost-integral and fully-integral solutions. Then, using a pre-processing procedure, we show how to convert an almost-integral solution to a fully-integral solution losing only a constant-factor in the approximation ratio. By further using a sparsification technique, the additive factor loss incurred by the conversion can be reduced to any $\epsilon > 0$.

## 1 Introduction

Clustering is a fundamental task studied by several scientific communities (such as operations research, biology, and of course computer science and machine learning) due to the diversity of its applications. Of the many ways in which the task of clustering can be formalized, the $k$-median and $k$-means problems are arguably the most commonly studied methods by these different communities, perhaps owing to the simplicity of their problem descriptions. In the $k$-median (resp. $k$-means) problem, we are given a set $C$ of clients, a set $F$ of potential facility locations, and a metric space $d(, \cdot, ) : C \cup F \to \mathbb{R}_{\ge 0}$, and the goal is to choose a subset $S \subseteq F$ of cardinality at most $k$ so as to minimize $\sum_{j \in C} d(j, S)$ (resp. $\sum_{j \in C} d^2(j, S)$), where $d(j, S) := \min_{i \in S} d(j, i)$. Both problems are NP-hard in the worst case, and this has led to a long line of research on obtaining efficient approximation algorithms. For the $k$-median problem, the current best known factors are an upper bound of $2.671$ [10], and a lower bound of $1 + 2/e \approx 1.736$ [27]. The $k$-means problem, while originally not as well studied in the theoretical community, is now increasingly gaining attention due to its importance in machine learning and data analysis. The best known approximations are 9 [1] and 6.357 [1] for general and Euclidean metrics respectively, while the lower bounds are $1 + 8/e \approx 3.94$ and $1.0013$ [32]

---

[*]Microsoft Research India. Email: `ravishan@cs.cmu.edu`

[†]Department of Computer Science and Engineering, University at Buffalo. Email: `shil@buffalo.edu`

[‡]Microsoft Research India. Email: `saisandeep192@gmail.com`

1

for general and Euclidean metrics respectively. Both problems admit PTASs [2, 21, 19] on fixed-dimensional Euclidean metrics.

Despite their simplicity and elegance, a significant shortcoming these formulations face in real-world data sets is that they are *not robust to noisy points*, i.e., a few outliers can completely change the cost as well as structure of solutions. To overcome this shortcoming, Charikar et al. [12] introduced the robust $k$-median (RkMed) problem (also called $k$-median with outliers), which we now define.

**Definition 1.1** (The Robust $k$-Median and $k$-Means Problems)**.** *The input to the Robust $k$-Median (*RkMed*) problem is a set $C$ of clients, a set $F$ of facility locations, a metric space $(C \cup F, d)$, integers $k$ and $m$. The objective is to choose a subset $S \subseteq F$ of cardinality at most $k$, and a subset $C^* \subseteq C$ of cardinality at least $m$ such that the total cost $\sum_{j \in C^*} d(j, S)$ is minimized. In the Robust $k$-Means (*RkMeans*) problem, we have the same input and the goal is to minimize $\sum_{j \in C^*} d^2(j, S)$.*

The problem is not just interesting from the clustering point of view. In fact, such a joint view of clustering and removing outliers has been observed to be more effective [15, 39] for even the sole task of outlier detection, a very important problem in the real world. Due to these use cases, there has been much recent work [15, 23, 40] in the applied community on these problems. However, their inherent complexity from the theoretical side is much less understood. For RkMed, Charikar et al. [12] give an algorithm that violates the number of outliers by a factor of $(1 + \epsilon)$, and has cost at most $4(1 + 1/\epsilon)$ times the optimal cost. Chen [16] subsequently showed a pure $O(1)$-approximation without violating $k$ or $m$. However, Chen's algorithm is fairly complicated, and the (unspecified) approximation factor is rather large. For the RkMeans problem, very recently, Gupta et al. [23] gave a bi-criteria $O(1)$-approximation violating the number of outliers, and Friggstad et al. [20] give *bi-criteria* algorithms that open $k(1 + \epsilon)$ facilities, and have an approximation factor of $(1 + \epsilon)$ on Euclidean and doubling metrics, and $25 + \epsilon$ on general metrics.

## 1.1 Our Results

In this paper, we present a simple and generic framework for solving clustering problems with outliers, which substantially improves over the existing approximation factors for RkMed and RkMeans.

**Theorem 1.2.** *For $\alpha_1 = \inf_{\tau > 1}(3\tau - 1)/\ln \tau \leq 7.081$, we have an $\alpha_1(1 + \varepsilon)$-approximation algorithm for* RkMed*, in running time $n^{O(1/\varepsilon^2)}$.*

**Theorem 1.3.** *For $\alpha_2 = \inf_{\tau > 1} \frac{(\tau+1)(3\tau-1)^2}{2(\tau-1)\ln \tau} \leq 53.002$, we have an $\alpha_2(1 + \varepsilon)$-approximation algorithm for* RkMeans*, in running time $n^{O(1/\varepsilon^3)}$.*

In fact, our framework can also be used to improve upon the best known approximation factors for two other basic problems in clustering, namely the matroid median (MatMed) and the knapsack median (KnapMed) problems. As in $k$-median, in both problems we are given $F$, $C$ and a metric $d$ over $F \cup C$ and the goal is to select a set $S \subseteq F$ so as to minimize $\sum_{j \in C} d(j, S)$. In MatMed, we require $S$ to be an independent set of a given matroid. In KnapMed, we are given a vector $w \in \mathbb{R}^F_{\geq 0}$ and a bound $B \geq 0$ and we require $\sum_{i \in S} w_i \leq B$. The previous best algorithms had factors 8 [42] and 17.46 [9] respectively.

**Theorem 1.4.** *For $\alpha_1 = \inf_{\tau > 1}(3\tau - 1)/\ln \tau \leq 7.081$, we have an efficient $\alpha_1$-approximation algorithm for* MatMed*, and an $\alpha_1(1 + \varepsilon)$-approximation for* KnapMed *with running time $n^{O(1/\varepsilon^2)}$.*

2

## 1.2 Our techniques

For clarity in presentation, we largely focus on RkMed in this section. Like many provable algorithms for clustering problems, our starting point is the natural LP relaxation for this problem. However, unlike the vanilla $k$-median problem, this LP has an unbounded integrality gap (see Section 3.1) for RkMed. Nevertheless, we can show that all is not lost due to this gap. Indeed, one of our main technical contributions is in developing an iterative algorithm which can round the natural LP to compute an *almost-integral* solution, i.e., one with at most two fractionally open facilities. By increasing the two possible fractional $y$ values to 1, we can obtain a solution with at most $k + 1$ open facilities and cost bounded by $\alpha_1$ times the optimal LP cost. So, the natural LP is good if we are satisfied with such a *pseudo-solution*[1]. So the unbounded integrality gap essentially arises as the gap between almost-integral solutions and fully-integral ones. In what follows, we first highlight how we overcome this gap, and then give an overview of the rounding algorithm.

### 1.2.1 Overcoming Gaps Between Almost-Integral and Integral Solutions

In the following, let $y_i$ denote the extent to which facility $i$ is open in a fractional solution. Note that once the $y_i$ variables are fixed, the fractional client assignments can be done in a greedy manner until a total of $m$ clients are connected fractionally. Now, assume that we have an almost-integral solution $y$ with two strictly fractional values $y_{i_1}$ and $y_{i_2}$. To round it to a fully-integral one, we need to increase one of $\{y_{i_1}, y_{i_2}\}$ to 1, and decrease the other to 0, in a direction that maintains the number of connected clients. Each of the two operations will incur a cost that may be unbounded compared to the LP cost, and they lead to two types of gap instances described in Section 3. We handle these two types separately.

The first type of instances, correspondent to the "increasing" operation, arise because the cost of the clients connected to the facility increased to 1 can be very large compared to the LP cost. We handle this by adding the so-called "star constraints" to the natural LP relaxation, which explicitly enforces a good upper bound on the total connection to each facility. The second type of instances, correspondent to the "decreasing" operation, arise because there could be many clients very near the facility being set to 0, all of which now incur a large connection cost. We handle this by a preprocessing step, where we derive an upper bound $R_j$ on the connection distance of each client $j$, which ensures that in any small neighborhood, the number of clients which are allowed to have large connection cost is small. The main challenge is in obtaining good bounds $R_j$ such that there exists a near-optimal solution respecting these bounds.

The above techniques are sufficient to give an $O(1)$ approximation for RkMed and RkMeans. Using an additional sparsification technique, we can reduce the gap between an almost-integral solution and an integral one to an additive factor of $\varepsilon$. Thus, our final approximation ratio is essentially the factor for rounding an LP solution to an almost-integral one. The sparsification technique was used by Li and Svensson [33] and by Byrka et al. [9] for the $k$-median and knapsack median problems respectively. The detail in how we apply the technique is different from those in [33, 9], but the essential ideas are the same. We guess $O_\varepsilon(1)$ balls of clients that incur a large cost in the optimum solution ("dense balls"), remove these clients, pre-open a set of facilities and focus on the residual instance. We show that the gap between almost-integral and fully-integral solutions on the residual instance (where there are no dense balls) is at most $\epsilon$.

---

[1]Indeed, an $O(1)$ pseudo-approximation with $k + 1$ open facilities is also implicitly given in [16], albeit with much a larger constant as a bound on the approximation ratio.

### 1.2.2 Iterative Rounding to Obtain Almost-Integral Solutions

We now highlight the main ideas behind our iterative rounding framework. At each step, we maintain a partition of $C$ into $C_{\text{full}}$, the set of clients which need to be fully connected, and $C_{\text{part}}$, the set of partially connected clients. Initially, $C_{\text{full}} = \emptyset$ and $C_{\text{part}} = C$. For each client, we also have a set $F_j$ denoting the set of facilities $j$ may connect to, and a "radius" $D_j$ which is the maximum connection cost for $j$; these can be obtained from the initial LP solution. Since the client assignments are easy to do once the $y$ variables are fixed, we just focus on rounding the $y$ variables. In addition to $C_{\text{full}}$, we also maintain a subset $C^*$ of full clients such that a) their $F_j$ balls are disjoint, and b) every other full client $j \notin C^*$ is "close" (within $O(1)D_j$) to the set $C^*$. Finally, for each client $j$, we maintain an *inner-ball* $B_j$ which only includes facilities in $F_j$ at distance at most $D_j/\tau$ from $j$ for some constant $\tau > 1$.

We then define the following *core constraints*: (i) $y(F_j) = 1$ for every $j \in C^*$, where $y(S) := \sum_{i \in S} y_i$, (ii) $y(F) \leq k$, and (iii) the total number of connected clients is at least $m$. As for non-core constraints, we define $y(F_j) \leq 1$ for all partial clients, and $y(B_j) \leq 1$ for all full clients. These constraints define our polytope $\mathcal{P} \subseteq [0,1]^F$.

Then, in each iteration, we update $y$ to be a vertex point in $\mathcal{P}$ that minimizes the linear function

$$\sum_{j \in C_{\text{part}}} d(i,j)y_i + \sum_{j \in C_{\text{full}}} \left( \sum_{i \in B_j} d(i,j)y_i + (1 - y(B_j))D_j/\tau \right).$$

Now, if none of the non-core constraints are tight, then this vertex point $y$ is defined by a laminar family of equalities along with a total coverage constraint, which is almost integral and so we output this. Otherwise, some non-core constraint is tight and we make the following updates and proceed to the next iteration. Indeed, if $y(F_j) = 1$ for some $j \in C_{\text{part}}$, we make it a full client and update $C^*$ accordingly. If $y(B_j) = 1$ for some $j \in C_{\text{full}}$, we update its $D_j$ to be $D_j/\tau$ and its $F_j$ to be $B_j$. In each iteration of this rounding, the cost of the LP solution is non-increasing (since $D_j$ and $F_j$ are non-increasing), and at the end, we relate the total connection cost of the final solution in terms of the LP objective using the property that every full client is within $O(D_j)$ from some client in $C^*$.

The iterative rounding framework is versatile as we can simply customize the core constraints. For example, to handle the the matroid median and knapsack median problems, we can remove the coverage constraint and add appropriate constraints. In matroid median, we require $y$ to be in the given matroid polytope, while in knapsack median, we add the knapsack constraint. For matroid median, the polytope defined by core constraints is already integral and this leads to our $\alpha_1$-approximation. For knapsack median, the vertex solution will be almost-integral, and again by using the sparsification ideas we can get our $(\alpha_1 + \epsilon)$-approximation. The whole algorithm can be easily adapted to RkMeans to get an $(\alpha_2 + \epsilon)$-approximation.

## 1.3 Related Work

The $k$-median and the related uncapacitated facility location (UFL) problem are two of the most classic problems studied in approximation algorithms. There is a long line of research for the two problems [35, 41, 25, 17, 28, 13, 26, 27, 37, 8, 11, 5] and almost all major techniques for approximation algorithms have been applied to the two problems (see the book of Williamson and Shmoys [43]). The input of UFL is similar to that of $k$-median, except that we do not have an upper bound $k$ on the number of open facilities, instead each potential facility $i \in F$ is given an opening cost $f_i \geq 0$. The goal is to minimize the sum of connection costs and facility opening costs. For the problem, the current best approximation ratio is 1.488 due to Li [34] and there is a hardness of 1.463 [22]. For the outlier version of uncapacitated facility, there is a 3-approximation due to [12]. This suggests that the outlier version of UFL is easier than that of $k$-median,

mainly due to the fact that constraints imposed by facility costs are soft ones, while the requirement of opening $k$ facilities is a hard one.

The $k$-means problem in Euclidean space is the clustering problem that is used the most in practice, and the most popular algorithm for the problem is the Lloyd's algorithm [36] (also called "the $k$-means" algorithm). However, in general, this algorithm has no worst case guarantee and can also have super-polynomial running time. There has also been a large body of work on bridging this gap, for example, by considering variants of the Lloyd's algorithm [4], or bounding the smoothed complexity [3], or by only focusing on instances that are "stable/clusterable" [31, 6, 38, 7, 18].

The MatMed problem was first studied by Krishnaswamy et. al. [29] as a generalization of the *red-blue median* problem[24], who gave an $O(1)$ approximation to the problem. Subsequently, there has been work [14, 42] on improving the approximation factor and currently the best upper bound is an 8-approximation by Swamy [42]. As for KnapMed, the first constant factor approximation was given by Kumar [30], who showed a 2700 factor approximation. The approximation factor was subsequently improved by [14, 42], and the current best algorithm is a $17.46$-approximation [9].

## 1.4 Paper Outline

We define some useful notations in Section 2. Then in Section 3, we explain our preprocessing procedure for overcoming the LP integrality gaps for RkMed/RkMeans. Then in Section 4, we give our main iterative rounding framework which obtains good almost-integral solutions. Section 5 will show how to covert an almost-integral solution to an integral one, losing only an additive factor of $\varepsilon$. Finally, in Sections 6 and 7, we present our $\alpha_1$ and $(\alpha_1 + \varepsilon)$ approximation algorithms for MatMed and KnapMed respectively.

*Getting a pseudo-approximation:* if one is only interested in getting a pseudo-approximation for RkMed/RkMeans, i.e, an $O(1)$-approximation with $k + 1$ open facilities, then our algorithm can be greatly simplified. In particular, the preprocessing step and the conversion step in Sections 3 and 5 are not needed, and proofs in Section 4 can be greatly simplified. Such readers can directly jump to Section 4 after reading the description of the natural LP relaxation in Section 3. In Section 4 we use the term pseudo-approximation setting to denote the setting in which we only need a pseudo-approximation.

## 2 Preliminaries

To obtain a unified framework for both RkMed and RkMeans, we often consider an instance $\mathcal{I} = (F, C, d, k, m)$ which could be an instance of either problem, and let a parameter $q$ denote the particular problem: $q = 1$ for RkMed instances and $q = 2$ for RkMeans instances. Because of the preprocessing steps, our algorithms deal with *extended* RkMed/RkMeans instances, denoted as $(F, C, d, k, m, S_0)$. Here, $F, C, d, k$ and $m$ are defined as before, and $S_0 \subseteq F$ is a set of *pre-opened facilities* that feasible solutions must contain. We assume that $d(i, i') > 0$ for $i \neq i' \in F$ since otherwise we can simply remove $i'$ from $F$; however, we allow many clients to be collocated, and they can be collocated with one facility.

We use a pair $(S^*, C^*)$ to denote a solution to a (extended) RkMed or RkMeans instance, where $S^* \subseteq F$ is the set of open facilities and $C^* \subseteq C$ is the set of connected clients in the solution, and each $j \in C^*$ is connected to its nearest facility in $S^*$. Note that given $S^*$, the optimum $C^*$ can be found easily in a greedy manner, and thus sometimes we only use $S^*$ to denote a solution to a (extended) RkMed/RkMeans instance. Given a set $S^* \subseteq F$ of facilities, it will be useful in the analysis to track the nearest facility for every $p \in F \cup C$ and the corresponding distance. To this end, we define the *nearest-facility-vector-pair* for $S^*$ to be $(\kappa^* \in (S^*)^{F \cup C}, c^* \in \mathbb{R}^{F \cup C})$, where $\kappa_p^* = \arg\min_{i \in S^*} d(i, p)$ and $c_p^* = d(\kappa_p^*, p) = \min_{i \in S^*} d(i, p)$ for

every $p \in F \cup C$. Though we only connect clients to $S^*$, the distance from a facility to $S^*$ will be useful in our analysis.

We use $y$ (and its variants) to denote a vector in $[0,1]^F$. For any $S \subseteq F$, we define $y(S) = \sum_{i \in S} y_i$ (same for the variants of $y$). Finally, given a subset $P \subseteq F \cup C$, a point $p \in F \cup C$ and radius $r \in \mathbb{R}$, we let $\text{Ball}_P(p,r) := \{p' \in P : d(p,p') \leq r\}$ to denote the set of points in $P$ with distance at most $r$ from $p$.

# 3 Preprocessing the RkMed/RkMeans Instance

In this section, we motivate and describe our pre-processing step for the RkMed/RkMeans problem, that is used to reduce the integrality gap of the natural LP relaxation. First we recap the LP relaxation and explain two different gap examples which, in turn illustrate two different reasons why this integrality gap arises. Subsequently, we will describe our ideas to overcome these two sources of badness.

## 3.1 Natural LP and Gap Examples

The natural LP for the RkMed/RkMeans problem is as follows.

$$\min \quad \sum_{i \in F, j \in C} x_{i,j} d^q(i,j) \qquad \text{s.t.} \qquad \text{(LP}_\text{basic})$$

$$\sum_i y_i \leq k \qquad\qquad\qquad\qquad \sum_{i \in F} x_{i,j} \leq 1 \quad \forall j \in C$$

$$x_{i,j} \leq y_i \quad \forall i \in F, j \in C \qquad\qquad \sum_{j \in C} \sum_{i \in F} x_{i,j} \geq m$$

In the correspondent integer programming, $y_i \in \{0,1\}$ indicates whether a facility $i \in F$ is open or not, and $x_{i,j} \in \{0,1\}$ indicates whether a client $j$ is connected to a facility $i$ or not. The objective function to minimize is the total connection cost and all the above constraints are valid. In the LP relaxation, we only require all the variables to be non-negative.

Notice that once the $y$ values are fixed, the optimal choice of the $x$ variables can be determined by a simple greedy allocation, and so we simply state the values of the $y$ variables in the following gap examples. For simplicity, we focus on RkMed problem when talking about the gap instances.

**Gap Example 1.** Instance (a) in Figure 1 contains two separate clusters in the metric: in the first cluster, there is a facility $i_1$ collocated with $t^3$ clients, and in the second cluster (very far from the first cluster), there is a facility $i_2$ at unit distance from $t^3 + t^2$ clients. Suppose $k = 1$ and $m = t^3 + t$. Clearly, the integral optimal solution is forced to open a center at $i_2$, thereby incurring a cost of $m = t^3 + t$. However, note that an LP solution can open $1 - 1/t$ fraction of facility $i_1$ and $1/t$ fraction of facility $i_2$. This can give a solution with cost $1/t \times (t^3 + t^2) = t^2 + t$, and connecting $(1 - 1/t) \times t^3 + 1/t \times (t^3 + t^2) = t^3 + t = m$ clients. Hence the integrality gap is $\frac{t^3+t}{t^2+t}$, which is unbounded as $t$ increases.

**Gap Example 2.** Instance (b) contains 3 facilities $i_0, i_1$ and $i_2$, collocated with $2t, 2t$ and $t$ clients respectively. $d(i_0, i_1) = 1$ and $i_2$ is far away from $i_0$ and $i_1$. We need to open $k = 2$ facilities and connect $m = 4t + 1$ clients. An integral solution opens $i_0$ and $i_2$, connect the $3t$ clients at $i_0$ and $i_2$, and connect $t + 1$ clients at $i_1$ to the open facility at $i_0$, incurring a total cost of $t + 1$. On the other hand, in an LP solution, we can set $y_0 = 1, y_1 = 1 - 1/t$ and $y_2 = 1/t$. We fully connect the $4t$ clients at $i_0$ and $i_1$, and connect each client at $i_2$ to an extent of $1/t$. The total number of connected clients is $4t + t \times 1/t = 4t + 1 = m$. Each client
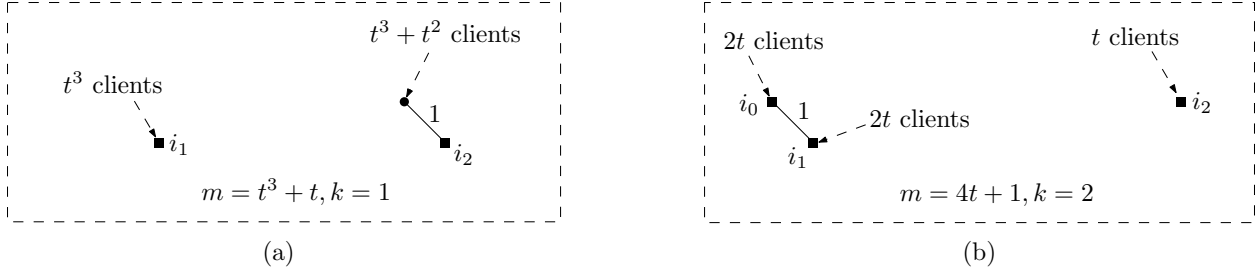
6

Figure 1: Two gap instances.

For instance (a), the optimum solution has cost $t^3 + t$. An LP solution with $y_1 = 1 - 1/t$ and $y_2 = 1/t$ has cost $t^2 + t$. For instance (b), the optimum solution has cost $t + 1$. An LP solution with $y_0 = 1, y_1 = 1 - 1/t$ and $y_2 = 1/t$ has cost 2.

at $i_1$ will be connected to $i_0$ with fraction $1 - y_1 = 1/t$, incurring a total cost of $2t \times 1/t = 2$. Thus, the integrality gap is $\frac{t+1}{2}$, which is again unbounded as $t$ increases.

## 3.2 Preprocessing

In both the examples in Section 3.1, the unbounded integrality gap essentially arises from needing to round an *almost-integral* solution $y$ into a *truly* integral one. This is no co-incidence as we show in Section 4 that we can indeed obtain an almost-integral solution with cost at most $\alpha_q$ times the optimal LP cost. Let us now examine the difficulty in rounding these last two fractional variables. As in both instances in Figure 1, suppose $i_1$ and $i_2$ are the two facilities with fractional $y$ values and let their $y$ values be $y_1 = 1 - \rho$ and $y_2 = \rho$ respectively, where $\rho > 0$ is a sub-constant. A natural idea is to increase one of these variables to 1 and decrease the other to 0. Suppose that, in order to maintain the number $m$ of connected clients, we are forced to increase $y_2$ to 1 and decrease $y_1$ to 0 (as in the two gap instances). Then two costs, corresponding to the two gap instances detailed above, will be incurred by this operation. The first cost is incurred by increasing $y_2$ to 1. Some clients were partially connected to $i_2$ in the almost-integral solution. In order to satisfy the coverage requirement, most of these clients need to be fully connected to $i_2$ in the integral solution. The instance in Figure 1(a) gives the example where this incurred cost is large compared to the optimal LP value. The second cost is incurred by decreasing $y_1$ to 0. Some clients were fully connected in the almost-integral solution, each being connected to $i_1$ with fraction $y_1 = 1 - \rho$ and to some other facility, say, $i_0 \notin \{i_1, i_2\}$ to extent $1 - y_1 = \rho$. Then, to satisfy the coverage requirement of these clients, we need to connect them to $i_0$ in the integral solution, and this cost could be large (see Figure 1(b)).

**Preprocessing Idea 1: Bounding Star Cost.** The first kind of bad instances described above are fairly straightforward to handle. Indeed, when increasing $y_2$ from $\rho$ to 1, the incurred cost is at most the cost of the "star" associated with $i_2$, i.e., the total connection cost of clients which are fractionally connected to $i_2$. We know the cost of such a star in the optimum solution is at most OPT and thus we can add the following constraints to overcome this type of bad instances: $\sum_{j \in C} x_{i,j} d^q(i,j) \leq y_i \cdot \text{OPT}$ for every $i \in F$.

In fact, we can bring down this additive error to $\rho \cdot \text{OPT}$ (where $\rho \to 0$ as $\varepsilon \to 0$) by *guessing* the set $S_1$ of centers corresponding to the expensive stars (whose connection cost exceeds $\rho \cdot \text{OPT}$) in the optimum solution and opening them. More importantly, we can strengthen (LP$_{\text{basic}}$) by enforcing the constraints bounding the connection cost to facilities $i \notin S_1$: $\sum_{j \in C} x_{i,j} d^q(i,j) \leq y_i \cdot \rho \cdot \text{OPT}$.

7

**Preprocessing Idea 2: Bounding Backup Cost.** Let us now explain how we handle the large incurred cost when we decrease $y_1$ to 0, as in Figure 1(b). Note that in the optimal solution, for any $R \geq 0$, the number of clients whose connection cost is at least $R$ is at most $\mathrm{OPT}/R$. In particular, in the example in Figure 1(b), if OPT is indeed $\Theta(1)$, then the number of clients at $i_1$ with connection cost at least 1 can be at most OPT. To this end, suppose we are able to *guess* a *specified* set of $\mathrm{OPT} = \Theta(1)$ clients located at $i_1$, and *disallow* the remaining collocated clients from connecting to facilities which are at distance 1 or greater. Then we could set $x_{ij} = 0$ for all disallowed connections which in turn will make the LP infeasible when our guess for OPT is $\Theta(1)$ for the bad instance from Figure 1(b). While it is difficult to get such good estimates on the disallowed connections of clients on a global scale, we show that we can indeed make such restrictions *locally*, which is our second pre-processing idea. We show that we can efficiently obtain a vector $(R_j)_{j \in C}$ of upper bounds on connection costs of clients which satisfies the following two properties: (i) for any $R$ and some constant $\delta > 0$, the number of clients in *any ball* of radius $\delta R$ with $R_j$ values more than some $\Theta(R)$ is at most $O(1)\mathrm{OPT}/R^q$, and (ii) there exists a feasible solution of cost at most $O(1)\mathrm{OPT}$ that respects these upper bounds $(R_j)_{j \in C}$ on client connections. This will then help us bound the re-assignment cost incurred by decreasing $y_1$ to 0 (i.e., shutting down facility $i_1$).

**Preprocessing Idea 3: Sparsification.** To reduce the additive factor we lose for handling the second type of bad instances to any small constant $\epsilon$, we use the so-called *sparsification* technique. Now, suppose we had the situation that, in the optimum solution, the number of clients in any ball of radius $\delta R$ with connection cost at least $R$ is miraculously at most $\rho \cdot \mathrm{OPT}/R^q$ for some small $0 \leq \rho < 1$. Informally, we call such instances as sparse instances. Then, we can bound total increase in cost by re-assigning these clients when shutting down one facility by $O(\rho)\mathrm{OPT}$. Indeed, our third pre-processing idea is precisely to correctly guess such "dense" regions (the clients and the facility they connect to in an optimum solution) so that the remaining clients are locally sparse. We note that this is similar to ideas employed in previous algorithms for $k$-median and KnapMed.

Motivated by these ideas, we now define the notion of *sparse extended instances* for RkMed/RkMeans.

**Definition 3.1** (Sparse Instances). *Suppose we are given an extended* RkMed/RkMeans *instance* $\mathcal{I}' = (F, C', d, k, m', S_0)$*, and parameters* $\rho, \delta \in (0, 1/2)$ *and* $U \geq 0$*. Let* $(S^*, C'^*)$ *be a solution to* $\mathcal{I}'$ *with cost at most* $U$*, and* $(\kappa^*, c^*)$ *be the nearest-facility-vector-pair for* $S^*$*. Then we say* $\mathcal{I}'$ *is* $(\rho, \delta, U)$*-sparse w.r.t the solution* $(S^*, C'^*)$*, if*

*(3.1a) for every* $i \in S^* \setminus S_0$*, we have* $\sum_{j \in C'^* : \kappa_j^* = i} (c_j^*)^q \leq \rho U$*,*

*(3.1b) for every* $p \in F \cup C'$*, we have* $\left| \mathrm{Ball}_{C'^*}(p, \delta c_p^*) \right| \cdot (c_p^*)^q \leq \rho U$*.*

Property (3.1a) requires the cost incurred by each $i \in S^* \setminus S_0$ in the solution $(S^*, C'^*)$ is at most $\rho U$. Property (3.1b) defines the sparsity requirement: roughly speaking, for every $p \in F \cup C'$, in the solution $(S^*, C'^*)$, the total connection cost of clients in $C'^*$ near $p$ should be at most $\rho U$.

We next show in the following theorem, that we can effectively reduce a general RkMed/RkMeans instance to a sparse extended instance with only a small loss in the approximation ratio.

**Theorem 3.2.** *Suppose we are given a* RkMed/RkMeans *instance* $\mathcal{I} = (F, C, d, k, m)$*, also parameters* $\rho, \delta \in (0, 1/2)$ *and an upper bound* $U$ *on the cost of the optimal solution* $(S^*, C^*)$ *to* $\mathcal{I}$ *(which is not given to us). Then there is an* $n^{O(1/\rho)}$*-time algorithm that outputs* $n^{O(1/\rho)}$ *many extended* RkMed/RkMeans *instances such that one of the instances in the set, say,* $\mathcal{I}'$*, has the form* $(F, C' \subseteq C, d, k, m' = |C^* \cap C'|, S_0 \subseteq S^*)$ *and satisfies:*

*(3.2a)* $\mathcal{I}'$ *is* $(\rho, \delta, U)$*-sparse w.r.t the solution* $(S^*, C^* \cap C')$*,*

8

(3.2b) $\frac{(1-\delta)^q}{(1+\delta)^q} \sum_{j \in C^* \setminus C'} d^q(j, S_0) + \sum_{j \in C^* \cap C'} d^q(j, S^*) \leq U.$

Before we give the proof, we remark that Property (3.2b) means that the cost of the solution $(S^*, C^* \cap C')$ for the residual sparse instance $\mathcal{I}'$ plus the approximate cost of reconnecting the $m - m'$ clients $C^* \setminus C'$ to the guessed facilities $S_0$ is upper bounded by $U$.

*Proof.* Let us first assume we know the optimum solution $(S^*, C^*)$ to the instance $\mathcal{I} = (F, C, d, k, m)$, and we will dispense with this assumption later. Let $(\kappa^*, c^*)$ be the nearest-facility-vector-pair for $S^*$. We initialize $C' = C$ and $S_0 = \emptyset$, and construct our instance iteratively until it becomes $(\rho, \delta, U)$-sparse w.r.t $(S^*, C^* \cap C')$. Our instance is always defined as $\mathcal{I}' = (F, C', d, k, m' = |C^* \cap C'|, S_0)$.

Indeed, to satisfy Property (3.1a), we add to $S_0$ the set of facilities $i \in S^*$ with $\sum_{j \in C^*: \kappa_j^* = i} (c_j^*)^q > \rho U$. There are at most $1/\rho$ many such facilities. After this, we have $\sum_{j \in C^* \cap C': \kappa_j^* = i} (c_j^*)^q \leq \sum_{j \in C^*: \kappa_j^* = i} (c_j^*)^q \leq \rho U$ for every $i \in S^* \setminus S_0$. This will always be satisfied since we shall only add facilities to $S_0$.

To guarantee Property (3.1b), we run the following iterative procedure: while there exists $p \in F \cup C$ such that $\left| \text{Ball}_{C^* \cap C'}(p, \delta c_p^*) \right| \cdot (c_p^*)^q > \rho U$, we update $S_0 \leftarrow S_0 \cup \{\kappa_p^*\}$ and $C' \leftarrow C' \setminus \text{Ball}_{C'}(p, \delta c_p^*)$. By triangle inequality, each client $j$ removed from $C'$ has $c_j^* \geq c_p^* - d(j, p) \geq (1 - \delta) c_p^*$. Also, $d(j, S_0) \leq d(j, \kappa_p^*) \leq d(p, \kappa_p^*) + d(j, p) \leq c_p^* + \delta c_p^* \leq \frac{1+\delta}{1-\delta} c_j^*$. Moreover, the total $(c_j^*)^q$ over all clients $j \in \text{Ball}_{C' \cap C^*}(p, \delta c_p^*)$ is at least $\left| \text{Ball}_{C^* \cap C'}(p, \delta c_p^*) \right| \cdot (1 - \delta)^q (c_p^*)^q > (1 - \delta)^q \rho U$. Thus, the procedure will terminate in less than $\frac{1}{\rho(1-\delta)^q} = O(1/\rho)$ iterations. After this procedure, we have $\left| \text{Ball}_{C^* \cap C'}(p, \delta c_p^*) \right| \cdot (c_p^*)^q \leq \rho U$ for every $p \in F \cup C$. That is, Property (3.1b) holds for the instance $\mathcal{I}' = (F, C', d, k, |C^* \cap C'|, S_0)$ w.r.t solution $(S^*, C^* \cap C')$. Thus Property (3.2a) holds.

Now we prove Property (3.2b).

$$\frac{(1-\delta)^q}{(1+\delta)^q} \sum_{j \in C^* \setminus C'} d^q(j, S_0) + \sum_{j \in C^* \cap C'} d^q(j, S^*) \leq \sum_{j \in C^* \setminus C'} (c_j^*)^q + \sum_{j \in C^* \cap C'} (c_j^*)^q = \sum_{j \in C^*} (c_j^*)^q \leq U.$$

The first inequality used the fact that for every $j \in C^* \setminus C'$, we have $d(j, S_0) \leq \frac{1+\delta}{1-\delta} c_j^*$. Thus Property (3.2b) holds.

Since we do not know the optimum solution $(S^*, C^*)$ for $\mathcal{I}$, we can not directly apply the above procedure. However, note that $C'$ is obtained from $C$ by removing at most $O(1/\rho)$ balls of clients, $S_0$ contains at most $O(1/\rho)$ facilities, and there are $m$ possibilities for $m' = |C^* \cap C'|$. Thus, there are at most $n^{O(1/\rho)}$ possible instances $(F, C', d, k, m', S_0)$, and we simply output all of them. □

We next show how, for $(\rho, \delta, U)$-sparse instances, we can effectively guess good upper bounds $R_j$ on the maximum connection cost of each client. The notion of sparse instances is only to get this vector of $R_j$'s and will not be needed after the proof of this theorem.

**Theorem 3.3.** *Let $\mathcal{I}' = (F, C', d, k, m', S_0)$ be a $(\rho, \delta, U)$-sparse instance w.r.t some solution $(S^* \supseteq S_0, C'^*)$ of $\mathcal{I}'$, for some $\rho, \delta \in (0, 1/2)$ and $U \geq 0$. Let $U' \leq U$ be the cost of $(S^*, C'^*)$ to $\mathcal{I}'$. Then given $\mathcal{I}', \rho, \delta, U$, we can efficiently find a vector $R = (R_j)_{j \in C'} \in \mathbb{R}_{\geq 0}^{C'}$, such that:*

*(3.3a) for every $t > 0$ and $p \in F \cup C'$, we have*

$$\left| \left\{ j \in \text{Ball}_{C'}\left( p, \frac{\delta t}{4 + 3\delta} \right) : R_j \geq t \right\} \right| \leq \frac{\rho(1 + 3\delta/4)^q}{(1 - \delta/4)^q} \cdot \frac{U}{t^q}, \tag{1}$$

*(3.3b) there exists a solution to $\mathcal{I}'$ of cost at most $(1 + \delta/2)^q U'$ where if $j$ is connected to $i$ then $d(i, j) \leq R_j$; moreover, the total cost of clients connected to any facility $i \notin S_0$ in this solution is at most $\rho(1 + \delta/2)^q U$.*

9

Property (3.3a) says that in a ball of radius $\Theta(t)$, the number of clients $j$ with $R_j \geq t$ is at most $O(\rho) \cdot \frac{U}{t^q}$. The first part of Property (3.3b) says that there is a near-optimal solution which respects these upper bounds $R_j$ on the connection distances of all clients $j$, and the second part ensures that all the star costs for facilities in $F \setminus S_0$ are still bounded in this near-optimal solution (akin to Property (3.1a)).

*Proof of Theorem 3.3.* We will now show that the following algorithm efficiently finds a vector $(\widehat{R}_j)_{j \in C'} \in \mathbb{R}_{\geq 0}^{C'}$, such that the following properties hold.

(i) for every $t > 0$ and $p \in F \cup C'$, we have

$$\left| \left\{ j \in \mathrm{Ball}_{C'} \left( p, \frac{\delta t}{4} \right) : \widehat{R}_j \geq t \right\} \right| \leq \frac{\rho}{(1 - \delta/4)^q} \cdot \frac{U}{t^q}, \tag{2}$$

(ii) there exists a solution to $\mathcal{I}'$ of cost at most $(1 + \delta/2)^q U'$ where if $j$ is connected to $i$ then $d(i, j) \leq (1 + 3\delta/4)\widehat{R}_j$; moreover, the total cost of clients connected to any facility $i \notin S_0$ in this solution is at most $\rho(1 + \delta/2)^q U$.

The proof then follows by setting $R_j = \widehat{R}_j(1 + 3\delta/4)$.

---

**Algorithm 1** Construct $\widehat{R}$

---

1: **for** every $j \in C'$ **do**: $\widehat{R}_j \leftarrow 0$
2: **for** every $t' \in \{d(i, j) : i \in F, j \in C'\} \setminus \{0\}$, in decreasing order **do**
3:     **for** every $j' \in C'$ such that $\widehat{R}_{j'} = 0$ **do**
4:         **if** letting $\widehat{R}_{j'} = t'$ will not violate (2) for $t = t'$ and any $p \in F \cup C'$ **then**
5:             $\widehat{R}_{j'} \leftarrow t'$

---

Consider Algorithm 1 that constructs the vector $\widehat{R}$. Property (i) holds at the beginning of the procedure since all clients $j$ have $\widehat{R}_j = 0$. We show that the property is always maintained as the algorithm proceeds. To this end, focus on the operation in which we let $\widehat{R}_{j'} = t'$ in Step 5.

- This will not violate (2) for $t = t'$ and any $p \in F \cup C'$ as guaranteed by the condition.

- It will not violate (2) for any $t > t'$ and $p \in F \cup C'$ since on the left side of (2) we only count clients $j$ with $\widehat{R}_j \geq t$ and setting $\widehat{R}_{j'}$ to $t' < t$ does not affect the counting.

- Focus on some $t < t'$ and $p \in F \cup C'$. Firstly, note that $\left| \mathrm{Ball}_{C'}(p, \delta t/4) : \widehat{R}_j \geq t \right| \leq \left| \mathrm{Ball}_{C'}(p, \delta t'/4) : \widehat{R}_j \geq t \right| = \left| \mathrm{Ball}_{C'}(p, \delta t'/4) : \widehat{R}_j \geq t' \right|$. The first inequality holds because we consider a bigger ball on the RHS, and the following equality holds since at this moment any $j$ has either $\widehat{R}_j \geq t'$ or $\widehat{R}_j = 0$, which implies $\widehat{R}_j \geq t \Leftrightarrow \widehat{R}_j \geq t'$. Now, since (2) holds for $t'$ and $p$, we have

$$\left| \mathrm{Ball}_{C'}(p, \delta t/4) : \widehat{R}_j \geq t \right| \leq \left| \mathrm{Ball}_{C'}(p, \delta t'/4) : \widehat{R}_j \geq t' \right| \leq \frac{\rho}{(1 - \delta/4)^q} \cdot \frac{U}{(t')^q} \leq \frac{\rho}{(1 - \delta/4)^q} \cdot \frac{U}{t^q}.$$

Thus, (2) also holds for this $t$ and $p$.

This finishes the proof of Property (i). Now consider Property (ii). Let $(\kappa^*, c^*)$ be the nearest-facility-vector-pair for $S^*$. Notice that the cost of the $(S^*, C'^*)$ is at most $U'$ as specified in the lemma statement. We build our desired solution incrementally using the following one-to-one mapping $f : C'^* \to C'$ such that the final set of clients which will satisfy Property (ii) will be $f(C'^*)$. Initially, let $f(j) = j$ for every $j \in C'^*$. To compute the final mapping $f$, we apply the following procedure. At a high-level, if a client $j \in C'^*$ has

connection cost $c_j^* > (1 + 3\delta/4)\widehat{R}_j$, then it means that there is a nearby ball with many clients with $\widehat{R}$ value at least $\widehat{R}_j$. We show that, since the instance is sparse, at least one of them is not yet mapped in $f(C'^*)$ and has a larger radius bound, and so we set $f(j)$ to $j'$.

Formally, for every client $j \in C'^*$, in non-decreasing order of $c_j^*$, if $c_j^* > (1 + 3\delta/4)\widehat{R}_j$, we update $f(j)$ to be a client in $C' \setminus f(C'^*)$ such that

(A) $d(f(j), j) \leq \delta c_j^*/2$, and

(B) $\widehat{R}_{f(j)} \geq c_j^*$.

Assuming we can find such an $f(j)$ in every iteration, then at the end of the procedure, we have that $f$ remains one-to-one. As for connections, for every $j \in C'^*$, let us connect $f(j)$ to $\kappa_j^*$. Then, we have $d(f(j), \kappa_j^*) \leq d(j, \kappa_j^*) + d(j, f(j)) \leq (1 + \delta/2)c_j^*$. Now it is easy to see that Property (ii) is satisfied. Indeed, the cost of the solution $(S^*, f(C'^*))$ is clearly at most $(1 + \delta/2)^q U'$, and the connection satisfies $d(f(j), \kappa_j^*) \leq (1 + 3\delta/4)\widehat{R}_{f_j}$. Moreover, for every facility $i \in S^* \setminus S_0$, the total cost of clients connected to $i$ is at most $(1 + \delta/2)^q \rho U$.

Thus, our goal becomes to prove that in every iteration we can always find an $f_j$ satisfying properties (A) and (B) above. To this end, suppose we have succeeded in all previous iterations and now we are considering $j \in C'^*$ with $c_j^* > (1 + 3\delta/4)\widehat{R}_j$. We thus need to update $f(j)$. Notice that at this time, we have $d(j', f(j')) \leq \delta c_{j'}^*/2$ for every $j' \in C'^*$, by the order we consider clients and the success of previous iterations. Now, focus on the iteration $t' = c_j^*$ in Algorithm 1. Then the fact that $c_j^* > (1 + 3\delta/4)\widehat{R}_j$ means that $j$ is not assigned a radius before and during iteration $t'$. Thus, there must be some $p \in F \cup C'$, such that $d(p, j) \leq \frac{\delta c_j^*}{4}$, and the set $H_j := \left\{ j' \in \mathrm{Ball}_{C'}\left( p, \delta c_j^*/4 \right) : \widehat{R}_{j'} \geq c_j^* \right\}$ such that $H_j \cup \{j\}$ has cardinality strictly more than $\frac{\rho}{(1-\delta/4)^q} \cdot \frac{U}{(c_j^*)^q}$. We claim that this set $H_j$ has a free client which $j$ can map to. Indeed, if there exists $j' \in H_j$ such that $j' \notin f(C'^*)$, then we can set $f(j) = j'$. This satisfies property (A) since $d(j', j) \leq d(j, p) + d(j', p) \leq \delta c_j^*/2$, and property (B) trivally, from the way $H_j$ is defined.

We now show there exists a $j' \in H_j \setminus f(C'^*)$, by appealing to the sparsity of the instance. Indeed, suppose $H_j \subseteq f(C'^*)$ for the sake of contradiction. Now, because $d(j', f_{j'}) \leq \delta c_j^*/2$ for every $j' \in C'^*$ at this time, we get that every client in $f^{-1}(H_j)$ has distance at most $\delta c_j^*/4 + \delta c_j^*/2 = 3\delta c_j^*/4$ from $p$. Thus,

$$\left| \mathrm{Ball}_{C'^*}(p, 3\delta c_j^*/4) \right| \geq |f^{-1}(H_j)| = |H_j| \geq \frac{\rho}{(1 - \delta/4)^q} \cdot \frac{U}{(c_j^*)^q}.$$

By triangle inequality, note that $c_p^* \geq c_j^* - d(j, p) \geq (1 - \delta/4)c_j^* \geq (3\delta/4)c_j^*$, as $\delta < 1$. Hence, we have

$$\left| \mathrm{Ball}_{C'^*}(p, \delta c_p^*) \right| \cdot (c_p^*)^q > \frac{\rho}{(1 - \delta/4)^q} \cdot \frac{U}{(c_j^*)^q} \cdot (1 - \delta/4)^q (c_j^*)^q = \rho U,$$

contradicting Property (3.1b) of the $(\rho, \delta, U)$-sparsity of $\mathcal{I}'$. $\qquad\square$

## 3.3 Putting Everything Together

Now we put everything together to prove Theorems 1.2 and 1.3. We will use, as a black-box, the following theorem, which will be the main result proved in Sections 4 and 5.

**Theorem 3.4** (Main Theorem). *Let* $\mathcal{I}' = (F, C', d, k, m', S_0)$ *be an extended* RkMed/RkMeans *instance,* $\rho, \delta \in (0, 1/2)$, $0 \leq U' \leq U$ *and* $R \in \mathbb{R}_{\geq 0}^{C'}$ *such that Properties (3.3a) and (3.3b) hold. Then there is an efficient randomized algorithm that, given* $\mathcal{I}', \rho, \delta, U$ *and* $R$, *computes a solution to* $\mathcal{I}'$ *with expected cost at most* $\alpha_q(1 + \delta/2)^q U' + O(\rho/\delta^q)U$.

11

Given a RkMed/RkMeans instance $\mathcal{I} = (F, C, d, k, m)$ and $\epsilon > 0$, we assume (by a standard binary-search idea) that we are given an upper bound $U$ on the cost of the optimum solution $(S^*, C^*)$ to $\mathcal{I}$ and our goal is to find a solution of cost at most $\alpha_q(1 + \epsilon)U$. To this end, let $\delta = \Theta(\epsilon)$ and let $\rho = \Theta(\epsilon^{q+1})$ be such that $(1 + \delta/2)^q \leq 1 + \epsilon/2$ and the $O(\rho/\delta^q)$ term before $U$ (in Theorem 3.4) is at most $\epsilon$.

Using Theorem 3.2, we obtain $n^{O(1/\rho)}$ many extended instances. We run the following procedure for each of them, and among all the computed solutions, we return the best one that is valid. Thus, for the purpose of analysis, we can assume that we are dealing with the instance $\mathcal{I}' = (F, C' \subseteq C, d, k, m' = |C^* \cap C'|, S_0 \subseteq S^*)$ satisfying properties (3.2a) and (3.2b) of Theorem 3.2. Defining $C'^* = C^* \cap C'$ and applying Theorem 3.3 and **??** in order, we obtain a solution $(\tilde{S} \supseteq S_0, \tilde{C} \subseteq C')$ to $\mathcal{I}'$ whose expected cost is $\alpha_q(1 + \delta/2)^q U' + O(\rho/\delta^q)U \leq \alpha_q(1 + \epsilon/2)U' + \epsilon U$. To extend the solution be feasible for $\mathcal{I}$, we simply greedily connect the cheapest $m - m'$ clients in $C \setminus C'$; the connection cost of these clients is at most $\sum_{j \in C^* \setminus C'} d^q(j, S_0)$. The expected cost of this solution is then at most $\sum_{j \in C^* \setminus C} d^q(j, S_0) + \alpha_q(1 + \epsilon/2)U' + \epsilon U$, which in turn is at most

$$
\leq \max\left\{\frac{(1+\delta)^q}{(1-\delta)^q}, \alpha_q(1 + \epsilon/2)\right\} \cdot \left(\frac{(1-\delta)^q}{(1+\delta)^q} \sum_{j \in C^* \setminus C'} d^q(j, S_0) + \sum_{j \in C^* \cap C'} d^q(j, S^*)\right) + \epsilon U
$$

$$
\leq \alpha_q(1 + \epsilon/2)U + \epsilon U \leq \alpha_q(1 + \epsilon)U.
$$

The second inequality uses Property (3.2b) and $\frac{(1+\delta)^q}{(1-\delta)^q} \leq \alpha_q$. The overall running time is $n^{O(1/\rho)} = n^{O(1/\epsilon^{q+1})}$, which is $n^{O(1/\epsilon^2)}$ for $q = 1$ and $n^{O(1/\epsilon^3)}$ for $q = 2$.

## 4 The Iterative Rounding Framework

This section and the next one are dedicated to the proof of the main theorem (Theorem 3.4). *For notational convenience, we replace the $\mathcal{I}', C', C'^*$ and $m'$ with $\mathcal{I}, C, C^*$ and $m$ subsequently. For clarity in presentation, we restate what we are given and what we need to prove after the notational change, without referring to properties stated in Section 3.* The input to our algorithm is an instance $\mathcal{I} = (F, C, d, k, m, S_0)$, parameters $\rho, \delta \in (0, 1/2), U \geq 0$ and a vector $R \in \mathbb{R}_{\geq 0}^C$. There is a parameter $U' \in [0, U]$ (which is not given to our algorithm) such that the following properties are satisfied:

(3.3a') for every $t > 0$ and $p \in F \cup C$, we have

$$
\left|\left\{j \in \mathrm{Ball}_C\left(p, \frac{\delta t}{4 + 3\delta}\right) : R_j \geq t\right\}\right| \leq \frac{\rho(1 + 3\delta/4)^q}{(1 - \delta/4)^q} \cdot \frac{U}{t^q}, \tag{3}
$$

(3.3b') there exists a solution to $\mathcal{I}$ of cost at most $(1+\delta/2)^q U'$ where if $j$ is connected to $i$ then $d(i, j) \leq R_j$; moreover, the total cost of clients connected to any facility $i \notin S_0$ is at most $\rho(1 + \delta/2)^q U$.

Our goal is to output a solution to $\mathcal{I}$ of cost at most $\alpha_q(1 + \delta/2)^q U' + O(\rho/\delta^q)U$. We do this in two parts: in this section, we use our iterative rounding framework and obtain an almost-integral solution for $\mathcal{I}$, with cost $\alpha_q(1 + \delta/2)^q U'$. This in fact immediately also gives us the desired pseudo-approximation solution. In the next section, we show how to convert the almost-integral solution to an integral one, with an additional cost of $O(\rho/\delta^q)U$.

We remark that in the pseudo-approximation setting, many parameters can be eliminated: we can set $S_0 = \emptyset, U = \infty, \delta = 0, R_j = \infty$ for every $j \in C$, and $\rho$ to be any positive number. Then Property (3.3a') trivially holds as $U = \infty$. Property (3.3b') simply says that there is an integral solution to $\mathcal{I}$ with cost at most $U'$. Our goal is to output an almost-integral solution of cost at most $\alpha_q U'$.

## 4.1 The Strengthened LP

For notational simplicity, we let $\tilde{U} = (1+\delta/2)^q U$ and $\tilde{U}' = (1+\delta/2)^q U'$. We now present our strengthened LP, where we add constraints (4)-(7) to the basic LP.

$$\min \quad \sum_{i \in F, j \in C} x_{i,j} d^q(i,j) \quad \text{s.t. constraints in (LP}_{\text{basic}}\text{) and} \qquad \text{(LP}_{\text{strong}}\text{)}$$

$$y_i = 1 \qquad\qquad \forall i \in S_0 \tag{4}$$

$$x_{i,j} = 0 \qquad\qquad \forall i, j \text{ s.t } d(i,j) > R_j \tag{5}$$

$$\sum_j d^q(i,j) x_{i,j} \leq \rho \tilde{U} y_i \qquad\qquad \forall i \notin S_0 \tag{6}$$

$$x_{i,j} = 0 \qquad\qquad \forall i \notin S_0, j \text{ s.t } d^q(i,j) > \rho \tilde{U} \tag{7}$$

Note that in the pseudo-approximation setting, we do not have Constraints (4) to (7), since $S_0 = \emptyset$, $R_j = \infty$ for all $j \in C$ and $\tilde{U} = \infty$. Thus (LP$_{\text{strong}}$) is the same as (LP$_{\text{basic}}$).

**Lemma 4.1.** *The cost of the optimum solution to* (LP$_{\text{strong}}$) *is at most* $\tilde{U}' = (1+\delta/2)^q U'$.

*Proof.* Consider the solution satisfying Property (3.3b'), and set $y_i = 1$ if $i$ is open and $x_{i,j} = 1$ if $j$ is connected to $i$ in the solution. Then the constraints in (LP$_{\text{basic}}$) and (4) hold as we have a valid solution. Constraints (5), (6) and (7) hold, and the cost of the solution is also at most $(1 + \delta/2)^q U' = \tilde{U}'$ from Property (3.3b').

In the pseudo-approximation setting, the lemma holds trivially since (LP$_{\text{basic}}$) is a valid LP relaxation to the instance $\mathcal{I}$. $\qquad\square$

## 4.2 Eliminating the $x_{ij}$ variables

Suppose we have an optimal solution $(x, y)$ to (LP$_{\text{strong}}$), such that for all $i, j$, either $x_{ij} = y_i$ or $x_{ij} = 0$, then we can easily eliminate the $x_{ij}$ variables and deal with an LP purely over the $y$ variables. Indeed, for the standard $k$-median problem, such a step, which is the starting point for many primal-based rounding algorithms, is easy by simply splitting each $y_i$ and making collocated copies of facilities in $F$. In our case, we need to take extra care to handle Constraint (6).

**Lemma 4.2.** *By adding collocated facilities to $F$, we can efficiently obtain a vector $y^* \in [0,1]^F$ and a set $F_j \subseteq \mathrm{Ball}_F(j, R_j)$ for every $j \in C$, with the following properties.*

*(4.2a)* *For each client $j$, we have $y^*(F_j) \leq 1$.*

*(4.2b)* *The total fractionally open facilities satisfies $y^*(F) \leq k$.*

*(4.2c)* *The total client coverage satisfies $\sum_{j \in C} y^*(F_j) \geq m$.*

*(4.2d)* *The solution cost is bounded: $\sum_{j \in C} \sum_{i \in F_j} d^q(i,j) y_i^* \leq \tilde{U}'$.*

*(4.2e)* *For each $i \in S_0$, we have $\sum_{i' \text{ collocated with } i} y_{i'}^* = 1$.*

*(4.2f)* *For every facility $i$ that is not collocated with any facility in $S_0$, the "star-cost" of $i$ satisfies $\sum_{j \in C : i \in F_j} d^q(i,j) \leq 2\rho \tilde{U}$.*

*Proof.* Let us denote the solution obtained by solving $(\text{LP}_{\text{strong}})$ as $(x, y)$. To avoid notational confusion, we first create a copy $F'$ of $F$ and transfer all the $y$ values from $F$ to $F'$, i.e., for each facility $i \in F$, there is a facility $i' \in F'$ collocated with $i$ and we set $y_{i'}^* = y_i$. Our final vector $y^*$ (and all $F_j$'s) will be entirely supported only in $F'$. We shall create a set $F_j \subseteq F'$ for each $j$, where initially we have $F_j = \emptyset$ for every $j$. The star cost of a facility $i' \in F'$ is simply $\sum_{j \in C : i' \in F_j} d^q(i', j)$. During the following procedure, we may split a facility $i' \in F'$ into two facilities $i'_1$ and $i'_2$ collocated with $i'$ with $y_{i'_1}^* + y_{i'_2}^* = y_{i'}^*$. Then, we update $F' \leftarrow F' \setminus \{i'\} \cup \{i'_1, i'_2\}$ and for each $j$ such that $i' \in F_j$, we shall update $F_j \leftarrow F_j \setminus \{i'\} \cup \{i'_1, i'_2\}$.

For every facility $i \in F$, for every client $j \in C$, with $x_{i,j} > 0$, we apply the following procedure. We order the facilities in $F'$ collocated with $i$ in non-decreasing order of their current star costs. We choose the first $o$ facilities in this sequence whose total $y^*$ value is exactly $x_{i,j}$; some facility may need to be split in order to guarantee this. We then add these facilities to $F_j$.

Properties (4.2a) to (4.2e) are easy to see since the new solution given by $y^*$ and $(F_j)_{j \in C}$ is the same as the solution $(x, y)$ up to splitting of facilities. It remains to prove Property (4.2f); this is trivial for the pseudo-approximation setting as $\tilde{U} = \infty$. We can view the above procedure for any fixed $i \in F$ as a greedy algorithm for makespan minimization on identical machines. Let $M$ be a large integer such that all considered fractional values are multiplies of $1/M$. There are $My_i$ machines, and for every $j \in C$, there are $Mx_{i,j}$ jobs of size $d^q(i, j)$ corresponding to each client $j$. There is an extra constraint that all the jobs corresponding to $j$ must be scheduled on different machines. The above algorithm is equivalent to the following: for each $j \in C$, we choose the $Mx_{i,j}$ machines with the minimum load and assign one job correspondent to $j$ to each machine. (This algorithm may run in exponential time; but it is only for analysis purpose.)

First, note that size of any job is at most $\rho\tilde{U}$ from (7). Thus, the highest load over all machines is at most the minimal one plus $\rho\tilde{U}$. The minimal load over all machines is at most the average load, which is equal to $\left(\sum_j d^q(i, j) \cdot \frac{Mx_{ij}}{My_i}\right) \leq \rho\tilde{U}$, by (6). Thus, the maximum load is at most $2\rho\tilde{U}$. Now, we redefine $F$ as $F'$ and finish the proof of the Lemma. $\qquad\square$

## 4.3 Random Discretization of Distances

Our next step in the rounding algorithm is done in order to optimize our final approximation factor. To this end, if $q = 1$, let $\tau = \arg\min \frac{3\tau-1}{\ln \tau} \approx 2.3603$, and then $\alpha_1 = \frac{3\tau-1}{\ln \tau} < 7.081$; if $q = 2$, let $\tau = \arg\min \frac{(\tau+1)(3\tau-1)^2}{2(\tau-1)\ln \tau} \approx 2.24434$ and then $\alpha_2 = \frac{(\tau+1)(3\tau-1)^2}{2(\tau-1)\ln \tau} \leq 53.002$. We choose a random offset $a \sim [1, \tau)$ such that $\ln a$ is uniformly distributed in $[0, \ln \tau)$. Then, define $D_{-2} = -1, D_{-1} = 0$, and $D_\ell = a\tau^\ell$ for every integer $\ell \geq 0$. Now, we *increase* each distance $d$ to the smallest value $d' \geq d$ belonging to the set $\{D_{-1}, D_0, D_1, D_2, \cdots, \}$. Formally, for every $i \in F, j \in C$, let $d'(i, j) = D_\ell$, where $\ell$ is the minimum integer such that $d(i, j) \leq D_\ell$. Note that the distances $d'(i, j)$ may not satisfy the triangle inequality anymore, but nevertheless this discretization will serve useful in our final analysis for optimizing the approximation factor.

**Lemma 4.3.** *For all $i, j$, we have $d'(i, j) \geq d(i, j)$ and $\mathbb{E}_a[d'^q(i, j)] = \frac{\tau^q - 1}{q \ln \tau} d^q(i, j)$.*

*Proof.* We can assume $d(i, j) \geq 1$, since the case for $d(i, j) = 0$ is trivial. Let $\beta = \log_\tau a$. $\beta$ is distibuted uniformly in $[0, 1)$. Let $d(i, j) = \tau^{\ell+p}$ where $\ell$ is an integer and $0 \leq p < 1$. When $\beta$ is less than $p$, $d(i, j)$ is rounded to $d'(i, j) = \tau^{\ell+1+\beta}$. If $\beta$ is at least $p$, then $d'(i, j) = \tau^{\ell+\beta}$.

$$
\mathbb{E}_a\left[d'^q(i, j)\right] = \int_{\beta=0}^p \tau^{q(\ell+1+\beta)} + \int_{\beta=p}^1 \tau^{q(\ell+\beta)} = \left.\frac{\tau^{q(\ell+1+\beta)}}{q \ln \tau}\right|_0^p + \left.\frac{\tau^{q(\ell+\beta)}}{q \ln \tau}\right|_p^1
$$
$$
= \frac{\tau^{q(\ell+1+p)} - \tau^{q(\ell+1)}}{q \ln \tau} + \frac{\tau^{q(\ell+1)} - \tau^{q(\ell+p)}}{q \ln \tau} = \frac{\tau^{q(\ell+1+p)} - \tau^{q(\ell+p)}}{q \ln \tau} = \tau^{q(\ell+p)} \frac{\tau^q - 1}{q \ln \tau} = \frac{\tau^q - 1}{q \ln \tau} d^q(i, j). \square
$$

## 4.4 Auxiliary LP and Iterative Rounding

The crucial ingredient of our rounding procedure is the following *auxiliary LP* ($\text{LP}_{\text{iter}}$) which we iteratively update and solve. At each step, we maintain a partition of $C$ into $C_{\text{full}}$, the set of clients which need to be fully connected, and $C_{\text{part}}$, the set of partially connected clients. For each client $j$, we also maintain a set $F_j$ of allowed connections, a *radius-level* $\ell_j$ for each client $j$ such that $D_{\ell_j}$ is the maximum connection cost for $j$, and a set $B_j$, called the *inner-ball* of $j$ which only includes facilities from $F_j$ at distance at most $D_{\ell_j-1}$ from $j$. Finally, we also maintain a subset $C^*$ of full clients such that their $F_j$ balls are disjoint, and also such that, every other full client $j \notin C^*$ is "close" (within $O(1)D_{\ell_j}$) to the set $C^*$. A small technicality is that for each facility $i \in S_0$, we also view $i$ as a *virtual client* that we shall include in $C^*$ at the beginning. This is done in order to ensure each $i \in S_0$ will be open eventually: since we split the facilities, $y^*$ obtained from Lemma 4.2 only has $\sum_{i' \text{ collocated with } i} y^*_{i'} = 1$ for every $i \in S_0$. Our algorithm will operate in a way that these virtual clients will never get removed from $C^*$, using which we will show that we open $S_0$ in the end.

Initially, $C_{\text{full}} \leftarrow \emptyset$, $C_{\text{part}} \leftarrow C$ and $C^* \leftarrow S_0$. For each client $j \in C$, the initial set $F_j$ is the one computed in Lemma 4.2, and we define $\ell_j$ to be the integer such that $D_{\ell_j} = \max_{i \in F_j} d'(i,j)$. For a virtual client $j \in S_0$, let $F_j$ be the set of facilities collocated with $j$ and thus $y^*(F_j) = 1$, and define $\ell_j = -1$ (thus $D_{\ell_j} = 0$). We remark that as we proceed with the rounding algorithm, the $F_j$'s and $\ell_j$'s will be updated but $F_j$ will only *shrink* and $\ell_j$ can only decrease. Crucially, at every step the set $F_j$ will always be contained in $\text{Ball}_F(j, D_{\ell_j})$.

$$
\min \quad \sum_{j \in C_{\text{part}}} \sum_{i \in F_j} d'^q(i,j)y_i + \sum_{j \in C_{\text{full}}} \left( \sum_{i \in B_j} d'^q(i,j)y_i + (1 - y(B_j))D^q_{\ell_j} \right) \quad \text{s.t.} \qquad (\text{LP}_{\text{iter}})
$$

$$
y(F) \leq k \qquad\qquad (8) \qquad\qquad\qquad y(F_j) \leq 1 \qquad \forall j \in C_{\text{part}} \qquad (11)
$$

$$
y(F_j) = 1 \qquad \forall j \in C^* \qquad (9) \qquad\qquad |C_{\text{full}}| + \sum_{j \in C_{\text{part}}} y(F_j) \geq m \qquad\qquad (12)
$$

$$
y(B_j) \leq 1 \qquad \forall j \in C_{\text{full}} \qquad (10) \qquad\qquad\qquad\qquad y_i \in [0,1] \qquad \forall i \in F \qquad (13)
$$

In the above LP, for a client $j \in C_{\text{part}}$, the quantity $y(F_j)$ denotes the extent that $j$ is connected, hence the constraint (11). Constraint (12) enforces that the total number of covered clients (full plus partial) is at least $m$. Then, constraint (9) ensures that clients in $C^*$ are covered fully. Since we don't enforce (9) for full clients that are not in $C^*$, we make sure that such client are close to $C^*$.

Now, we describe the objective function of ($\text{LP}_{\text{iter}}$). Notice that we use $d'$ instead of $d$ in the objective function. For any client $j \in C_{\text{part}}$, $y(F_j)$ is precisely the extent to which $j$ is connected, and so, we can simply use $\sum_{i \in F_j} d'^q(i,j)y_i$ to denote the connection cost of $j$. For a client $j \in C_{\text{full}}$, $j$ is required to be fully connected but we may not have $y(F_j) = 1$. Hence, $\sum_{i \in F_j} d'^q(i,j)y_i$ is no longer a faithful representation of its connection cost, and hence we need to express the cost in a different manner. To this end, for every $j \in C_{\text{full}}$, we require $y(B_j) \leq 1$ in (10). We guarantee that $B_j$ is always $\{i \in F_j : d'(i,j) \leq D_{\ell_j-1}\}$. Then, for connections between $j$ and any facility $i \in B_j$, we use the rounded distance $d'(i,j)$ in the objective function. For the remaining $1 - y(B_j)$ fractional connection of $j$, we use the term $D_{\ell_j}$. This gives us the objective function, and in turn completes the description of ($\text{LP}_{\text{iter}}$).

**Lemma 4.4.** *The $y^*$ computed in Section 4.2 is a feasible solution to* ($\text{LP}_{\text{iter}}$). *Moreover, the expected value of the objective function (over the random choice of $a$ defined in Section 4.3) is at most $\frac{\tau^q-1}{q\ln\tau}\tilde{U}'$.*

*Proof.* Initially we have $C^* = S_0$, $C_{\text{full}} = \emptyset$ and $C_{\text{part}} = C$. The feasibility of $y^*$ is guaranteed by Properties (4.2a)-(4.2c), and the fact that for every virtual client $j \in S_0$, we have $y^*(F_j) = 1$. Now,

15

since $C_{\text{part}} = C$ and $C_{\text{full}} = \emptyset$, the objective value of $y^*$ w.r.t (LP$_{\text{iter}}$) is $\sum_{j \in C, i \in F_j} d'(i,j) y_i^*$. Now, from Property (4.2d), we have that $\sum_{j \in C, i \in F_j} d(i,j) y_i^* \leq \tilde{U}'$. Combining this with Lemma 4.3 bounds the objective value. $\qquad\square$

---

**Algorithm 2** Iterative Rounding Algorithm

---

- **Input**: $\mathcal{I} = (F, C, d, k, m, S_0), \rho, \delta \in (0, 1/2), U \geq 0, y^* \in [0,1]^F, (F_j)_{j \in C \cup S_0}$ and $(\ell_j)_{j \in C \cup S_0}$

- **Output**: a new solution $y^*$ which is almost-integral

---

1: $C_{\text{full}} \leftarrow \emptyset, C_{\text{part}} \leftarrow C, C^* \leftarrow S_0$
2: **while** true **do**                                                                        ▷ (main loop)
3:      find an optimum vertex point solution $y^*$ to (LP$_{\text{iter}}$)
4:      **if** there exists some $j \in C_{\text{part}}$ such that $y^*(F_j) = 1$ **then**
5:           $C_{\text{part}} \leftarrow C_{\text{part}} \setminus \{j\}, C_{\text{full}} \leftarrow C_{\text{full}} \cup \{j\}, B_j \leftarrow \{i \in F_j : d'(i,j) \leq D_{\ell_j - 1}\}$ update-$C^*(j)$
6:      **else if** there exists $j \in C_{\text{full}}$ such that $y^*(B_j) = 1$ **then**
7:           $\ell_j \leftarrow \ell_j - 1, F_j \leftarrow B_j, B_j \leftarrow \{i \in F_j : d'(i,j) \leq D_{\ell_j - 1}\}$ update-$C^*(j)$
8:      **else**
9:           break
10: **return** $y^*$

---

update-$C^*(j)$:
1: **if** there exists no $j' \in C^*$ with $\ell_{j'} \leq \ell_j$ and $F_j \cap F_{j'} \neq \emptyset$ **then**
2:      remove from $C^*$ all $j'$ such that $F_j \cap F_{j'} \neq \emptyset$
3:      $C^* \leftarrow C^* \cup \{j\}$

---

We can now describe our iterative rounding algorithm, formally stated in Algorithm 2. In each iteration, we solve the LP to obtain a vertex solution $y^*$ in Step 3. If (11) is tight for some partial client $j \in C_{\text{part}}$, then we update $C_{\text{full}} \leftarrow C_{\text{full}} \cup \{j\}$ and remove $j$ from $C_{\text{part}}$. We also update $C^*$ to ensure that there is a facility in $C^*$ that is close to all full clients; this is done in the procedure update-$C^*$. Likewise if (10) is tight for some client $j \in C_{\text{full}}$, then we decrease $\ell_j$ by 1, update $F_j \leftarrow B_j$ and $B_j$ to be an even smaller set. We again call update-$C^*(j)$ to ensure $C^*$ is close to all full clients — this is needed since we decreased $\ell_j$ and the definition of "close" becomes more strict. If neither of the constraints (10) and (11) are tight for $y^*$, then we show that $y^*$ is almost integral and we return $y^*$.

## 4.5  Analysis

Our proof proceeds as follows: we first show that, if the vertex point $y^*$ computed satisfies all constraints (10) and (11) with strict inequality, then it is almost integral. Next we show that the objective value of the solutions computed is non-increasing as we iterate. Then, we show that all full clients are close to $C^*$, using which we conclude that there is one unit of fractional facilities opened within distance $O(1) D_{\ell_j}$ for all $j \in C_{\text{full}}$. This suggests that the objective function of (LP$_{\text{iter}}$) captures the actual cost well. Finally, we show that all virtual clients in $S_0$ will always be in $C^*$, implying that there will be an integrally open facility at each location in $S_0$. Combining all these leads to an almost-integral solution with bounded cost.

We begin by describing some simple invariants that are maintained in the algorithm; throughout this section, we assume that every step indicated by a number in Algorithm 2 is atomic.

**Claim 4.5.** *After each step of the algorithm, the following hold.*

1. *The sets $C_{\mathrm{full}}$ and $C_{\mathrm{part}}$ form a partition of $C$, $S_0 \subseteq C^*$ and $C^* \setminus S_0 \subseteq C_{\mathrm{full}}$.*

2. *The sets $\{F_j : j \in C^*\}$ are mutually disjoint.*

3. *If $j \in C_{\mathrm{full}}$, then $B_j = \{i \in F_j : d'(i,j) \leq D_{\ell_j - 1}\}$.*

4. *For every $j \in C$ and $i \in F_j$, we have $d'(i,j) \leq D_{\ell_j}$.*

5. *For every $j$, $D_{\ell_j} \leq \tau R_j$.*

6. *For every $j$, $\ell_j \geq -1$.*

*Proof.* $C_{\mathrm{full}}$ and $C_{\mathrm{part}}$ form a partition of $C$ since we only *move* clients from $C_{\mathrm{part}}$ to $C_{\mathrm{full}}$. Moreover, virtual clients in $S_0$ will never be removed from $C^*$ since each such client $j$ has $\ell_j = -1$. Finally, we call update-$C^*(j)$ only if $j$ is already in $C_{\mathrm{full}}$ and we only add $j$ to $C^*$ in the procedure. Thus, Property 1 holds. Also, before we add $j$ to $C^*$ in update-$C^*(j)$, we removed all $j'$ from $C^*$ such that $F_j \cap F_{j'} \neq \emptyset$, and so Property 2 holds. Property 3 holds since every time we move $j$ from $C_{\mathrm{part}}$ to $C_{\mathrm{full}}$, or update $F_j$ for some $j \in C_{\mathrm{full}}$, we define $B_j = \{i \in F_j : d'(i,j) \leq D_{\ell_j - 1}\}$. Property 4 holds at the beginning of the algorithm, and before we change $F_j$ to $B_j$ in Step 7, we have $d'(i,j) \leq D_{\ell_j - 1}$ for every $i \in B_j$. Property 5 holds initially because initially $D_{\ell_j} = \max_{i \in F_j} d'(i,j) \leq \tau \max_{i \in F_j} d(i,j) \leq \tau R_j$, and $D_{\ell_j}$ is non-increasing over time. Property 6 holds since if $\ell_j = -1$ for some $j \in C_{\mathrm{full}}$ then $B_j$ is empty thus $y^*(B_j)$ will never become 1. $\square$

We now proceed to the analysis of the algorithm, by proving the following lemmas.

**Lemma 4.6.** *If after Step 3 in Algorithm 2, none of the constraints* (10) *and* (11) *are tight for $y^*$, then there are at most $2$ strictly fractional $y^*$ variables.*

*Proof.* Assume $y^*$ has $n' \geq 3$ fractional values. We pick $|F|$ independent tight inequalities in (LP$_{\mathrm{iter}}$) that defines $y^*$. Without loss of generality, we can assume if $y_i^* \in \{0, 1\}$, then the constraint (13) for $i$ is picked. Thus, there are exactly $|F| - n'$ tight constraints from among (13) and $n'$ tight constraints from among (8), (9) or (12). Each tight constraint of form (9) should contain at least 2 fractional variables, and the tight constraint (8) should contain 2 fractional variables, that are not contained in any tight constraint of form (9). Thus, the number of tight constraints of form (8), (9) or (12) is at most $n'/2 + 1$. This is less than $n'$ if $n' \geq 3$. Thus, $n' \leq 2$. $\square$

**Lemma 4.7.** *After any step of the algorithm, the maintained solution $y^*$ is feasible to* (LP$_{\mathrm{iter}}$).

*Proof.* We only need to consider the time point after we run Step 5 or 7 during the algorithm. In Step 5, we move $j$ from $C_{\mathrm{part}}$ to $C_{\mathrm{full}}$ and define $B_j$. We were guaranteed that $y^*(F_j) = 1$ and thus after the step we have $y^*(B_j) \leq y^*(F_j) \leq 1$. Before Step 7, we were guaranteed $y^*(B_j) = 1$. Thus, after the step, we have $y^*(F_j) = 1$ and $y^*(B_j) \leq 1$. $\square$

**Lemma 4.8.** *In every iteration, the* (LP$_{\mathrm{iter}}$) *objective value of the solution $y^*$ is non-increasing.*

*Proof.* In Statement 3, we resolve the LP and thus the value of $y^*$ can only go down (due to Lemma 4.7). Thus, we only need to consider the situation in which we change the objective function, which happens in Step 5 and 7. Indeed, we show that the objective value is not affected by these operations.

1. Consider the situation where a client $j$ is moved to $C_{\text{full}}$ from $C_{\text{part}}$ in Statement 5. The cost that $j$ is contributing in the old LP is

$$\sum_{i \in F_j} d'^q(i,j) y_i^* = \sum_{i \in F_j: d'(i,j) \le D_{\ell_j - 1}} d'^q(i,j) y_i^* + \sum_{i \in F_j: d'(i,j) = D_{\ell_j}} d'^q(i,j) y_i^*$$

$$= \sum_{i \in B_j} d'^q(i,j) y_i^* + (1 - y^*(B_j)) D_{\ell_j}^q,$$

which is exactly the contribution of $j$ to the new objective function.

2. Then, consider Statement 7. To avoid confusion, let $\ell_j$, $F_j$ and $B_j$ correspond to the values before executing the statement, and $\ell_j'$, $F_j'$ and $B_j'$ correspond to the values after executing the statement. The contribution of $j$ to the old function is

$$\sum_{i \in B_j} d'^q(i,j) y_i^* + (1 - y(B_j)) D_{\ell_j}^q = \sum_{i \in B_j} d'^q(i,j) y_i^* = \sum_{i \in F_j'} d'^q(i,j) y_i^*$$

$$= \sum_{i \in F_j': d'(i,j) \le D_{\ell_j' - 1}} d'^q(i,j) y_i^* + \sum_{i \in F_j': d'(i,j) = D_{\ell_j'}} d'^q(i,j) y_i^* = \sum_{i \in B_j'} d'^q(i,j) y_i^* + \left(1 - y(B_j')\right) D_{\ell_j'}^q,$$

which is the contribution of $j$ to the new function. $\qquad\square$

**Lemma 4.9.** *At the conclusion of the algorithm, for every $j \in C_{\text{full}}$, there exists at least 1 unit of open facilities within distance $\frac{3\tau-1}{\tau-1} D_{\ell_j}$ from $j$. Formally, $\sum_{i: d(i,j) \le \frac{3\tau-1}{\tau-1} D_{\ell_j}} y_i^* \ge 1$.*

*Proof.* The following claims will serve useful in the proof of the lemma.

**Claim 4.10.** *If at any stage, $F_j \cap F_{j'} \ne \emptyset$, then we have that $d(j,j') \le D_{\ell_j} + D_{\ell_{j'}}$.*

**Claim 4.11.** *Consider a client $j$ which is added to $C^*$ at Step 3 of update-$C^*(j)$ when its radius-level $\ell_j$ was some $\ell \ge -1$. Then, there exists at least 1 unit of open facilities in the final solution $y^*$ within distance $\frac{\tau+1}{\tau-1} D_\ell$ from $j$.*

*Proof.* The proof is by induction on radius-levels. Indeed, clients with radius-level $-1$ can never be removed from $C^*$, so if such clients are added to $C^*$, then they remain in $C^*$ at the end, and so the claim is true. So suppose the claim is true up to some radius-level $\ell$-1, and consider a client $j$ which got added to $C^*$ with radius-level $\ell$ at some time in Step 3 of update-$C^*(j)$. Then, either it remains in $C^*$ till the end, in which case we are done (because its $\ell_j$ is non-increasing over time), or another client $j'$ was subsequently added to $C^*$ such that $\ell_{j'} < \ell_j$ and $F_j \cap F_{j'} \ne \emptyset$. Here, $\ell_j$ and $\ell_{j'}$ are both the radius-levels of $j$ and $j'$ at the time $j$ was removed by $j'$. Since these values are non-increasing over time, we get that $\ell_{j'} < \ell_j \le \ell$. Now, in this latter case, note that by the induction hypothesis there exists one unit of facility in the final solution $y^*$ within distance $\frac{\tau+1}{\tau-1} D_{\ell_{j'}} \le \frac{\tau+1}{\tau-1} \frac{D_{\ell_j}}{\tau} \le \frac{\tau+1}{\tau-1} \frac{D_\ell}{\tau}$ from $j'$. Moreover, we know $d(j,j') \le D_\ell + D_\ell/\tau$ from Claim 4.10. The proof then follows from the triangle inequality and $\frac{\tau+1}{\tau-1} \frac{D_\ell}{\tau} + D_\ell + \frac{D_\ell}{\tau} = \frac{\tau+1}{\tau-1} D_\ell$. $\qquad\square$

Now the proof of Lemma 4.9 is simple: consider a client $j$ which belongs to $C_{\text{full}}$ at the end of the algorithm, and consider the last time when update-$C^*(j)$ is invoked, and let $\ell_j$ denote its radius-level at this time. Note that the value $\ell_j$ has not changed subsequently.

At this point, there are two cases depending on whether $j$ was added to $C^*$ or not. If it was added, we are done by Claim 4.11. In case it was not added, then there must be a client $j'$ which belonged to $C^*$ such that

$F_{j'} \cap F_j \neq \emptyset$ and $\ell_{j'} \leq \ell_j$. Then, note that $d(j, j') \leq 2D_{\ell_j}$ from Claim 4.10, and moreover, by Claim 4.11, there is at least 1 unit of open facilities in the final solution $y^*$ within distance $\frac{\tau+1}{\tau-1}D_{\ell_{j'}} \leq \frac{\tau+1}{\tau-1}D_{\ell_j}$ from $j'$. We can then complete the proof by applying triangle inequality. □

With all the above lemmas, we can conclude this section with our main theorem.

**Theorem 4.12.** *When the algorithm terminates, it returns a valid solution $y^* \in [0,1]^F$ to (LP$_{\text{basic}}$) with at most two fractional values. $y_i^* = 1$ for every $i \in S_0$. Moreover, the cost of $y^*$ w.r.t (LP$_{\text{basic}}$), in expectation over the randomness of $a$ is at most $\alpha_q \tilde{U}'$.*

*Proof.* Firstly, at the culmination of our iterative algorithm, note that none of the constraints (10) and (11) are tight for the final $y^*$, otherwise our algorithm would not have terminated. As a result, it follows from Lemma 4.6 that there are at most two strictly fractional $y^*$ variables.

As for bounding the cost, we show this by exhibiting a setting of $x_{ij}$ for the given $y^*$ with the desired bound on cost. To this end, for clients in $C_{\text{part}}$, set $x_{ij} = y_i^*$ for all $i \in F_j$. For clients in $C_{\text{full}}$, set $x_{ij} = y_i^*$ if $i \in B_j$, and set the remaining $1 - y^*(B_j)$ fractional amount arbitrarily among open facilities within distance $\frac{3\tau-1}{\tau-1}D_{\ell_j}$ from $j$ — Lemma 4.9 guarantees the existence of such open facilities.

We now verify that such a solution satisfies the constraints of (LP$_{\text{basic}}$). Indeed, it is easy to see that $\sum_i y_i^* \leq k$, and also that $x_{i,j} \leq y_i^*$. As for the coverage constraint, note that for full clients $j$, we have $\sum_i x_{i,j} = 1$ and for partial clients $j$, we have $\sum_i x_{i,j} = \sum_{i \in F_j} y_i^* = y^*(F_j) \leq 1$. Finally the total coverage $\sum_{j,i} x_{i,j} = |C_{\text{full}}| + \sum_{j \in C_{\text{part}}} y^*(F_j) \geq m$ since $y^*$ satisfies (12). Also, every virtual client $j \in S_0$ is always in $C^*$ and we always have $F_j$ is the set of facilities collocated with $j$. Eventually we have $y^*(F_j) = 1$. W.l.o.g we can assume $y_i^* = 1$ for every $i \in S_0$.

To complete the proof, we compute the cost of the solution we have constructed. Indeed, the contribution of a partial client $j$ in our solution is at most their contribution to the objective value of $y^*$ for (LP$_{\text{iter}}$). This is because $\sum_i x_{i,j} d(i,j) \leq \sum_i x_{i,j} d'(i,j) = \sum_{i \in F_j} d'(i,j) y_i^*$. Now consider a full client $j \in C_{\text{full}}$. By our setting we have $x_{i,j} = y_i^*$ for $i \in B_j$ and hence $\sum_{i \in B_j} x_{i,j} d(i,j) \leq \sum_{i \in B_j} x_{i,j} d'(i,j) = \sum_{i \in B_j} y_i^* d'(i,j)$. Hence, the contribution to the objective value of assignment within $B_j$ for full clients has a one-to-one correspondence with their contribution in the auxiliary LP (LP$_{\text{iter}}$). It remains to bound the connection cost to facilities outside $B_j$. For bounding this, note that the extend of outside connections for each client $j$ is exactly $1 - y^*(B_j)$, and from Lemma 4.9, we get that all such connections are at a distance at most $\frac{3\tau-1}{\tau-1}D_{\ell_j}$ from $j$, incurring a cost of $\frac{(3\tau-1)^q}{(\tau-1)^q}D_{\ell_j}^q$. In contrast, these connections contribute a total of $D_{\ell_j}^q(1 - y^*(B_j))$ to the objective value of (LP$_{\text{iter}}$). Thus, the cost of our solution w.r.t (LP$_{\text{basic}}$) is at most $\frac{(3\tau-1)^q}{(\tau-1)^q}$ times the cost of $y^*$ w.r.t (LP$_{\text{iter}}$), which in expectation is at most $\frac{\tau^q-1}{q \ln \tau}$ times $\tilde{U}'$, by Lemma 4.4 and 4.8. Hence, the cost of our constructed almost-integral solution w.r.t. (LP$_{\text{basic}}$) is at most $\frac{(3\tau-1)^q}{(\tau-1)^q} \cdot \frac{\tau^q-1}{q \ln \tau} \cdot \tilde{U}'$ in expectation. By the way we choose $\tau$, this is exactly $\alpha_q \tilde{U}'$. □

*Pseuso-approximation Setting:* for this setting, we simply change the two possible fractionally open facilities to be integrally open, and hence we can obtain an $\alpha_q$-approximation for RkMed/RkMeans with $k+1$ open facilities.

## 5 Opening exactly $k$ facilities

In this section, we show how to convert the almost-integral $y^*$ computed by the iterative rounding Algorithm 2 into a fully integral solution with bounded cost. Indeed, this last step of converting an almost-integral solution

into a fully-integral one is why we needed the three steps of pre-processing and the stronger LP relaxation in the first place. From Theorem 4.12, we know that the final $y^*$ computed has at most two fractional values and moreover, the expected cost of $y^*$ to the (LP$_{\text{basic}}$) is at most $\alpha_q \tilde{U}'$.

**Claim 5.1.** *The solution $y^*$ either is already fully integral, or has exactly two fractional values, say, $y_{i_1}$ and $y_{i_2}$. Moreover it holds that $y_{i_1}^* + y_{i_2}^* = 1$.*

*Proof.* WLOG we can assume $y^*(F) = k$, as if $y^*(F) < k$, then we could open more facilities in $y^*$. This along with the fact that $y^*$ has at most two fractional values establishes the proof. $\square$

Now, in case $y^*$ is already integral, Theorem 4.12 already gives a solution of cost $\alpha_q \tilde{U}'$. Hence, we focus on the other case, where, by Claim 5.1, there are exactly two fractional facilities, say $i_1$ and $i_2$. We assume that we are given as input the $y^*$ output by Algorithm 2, along with the sets $F_j$ for $j \in C$, radius-levels $\ell_j$, and the partition of $C$ into $C_{\text{part}}$ and $C_{\text{full}}$ maintained by Algorithm 2 when it terminated.

Let $C_1 = \{j \in C_{\text{part}} : i_1 \in F_j \text{ and } i_2 \notin F_j\}$, and similarly let $C_2 = \{j \in C_{\text{part}} : i_1 \notin F_j \text{ and } i_2 \in F_j\}$. By renaming $i_1$ and $i_2$ if necessary, we assume $|C_1| \geq |C_2|$. Our final solution $\hat{y}$ will be defined as: $\hat{y}_{i_1} = 1$, $\hat{y}_{i_2} = 0$, and $\hat{y}_i = y_i^*$ for $i \notin \{i_1, i_2\}$. Our algorithm will connect $C_{\text{full}} \cup C_1$ to their respective nearest facilities in $\hat{y}$.

**Lemma 5.2.** *The solution described above connects at least $m$ clients.*

*Proof.* This is easy to see, as it covers all full clients, and the only partial clients which are connected in the original LP solution $y^*$ output by Algorithm 2 are those in $C_1$ and $C_2$, contributing a total of $y_{i_1}|C_1| + y_{i_2}|C_2|$ to constraint (12). Since $|C_1| \geq |C_2|$ and $y_{i_1} + y_{i_2} = 1$, we have $|C_1| + |C_{\text{full}}| \geq y_{i_1}|C_1| + y_{i_2}|C_2| + |C_{\text{full}}| \geq m$. $\square$

**Lemma 5.3.** *The total connection cost of all clients in $C_1$ is at most $2\rho\tilde{U}$.*

*Proof.* Notice that $i_1$ is open in the solution $\hat{y}$. The facility $i_1$ is not collocated with $S_0$ since there are exactly one integral facility open at each location in $S_0$. The connection cost is thus at most

$$\sum_{j \in C_1} d^q(i_1, j) \leq \sum_{j: i_1 \in F_j} d^q(i_1, j) \leq 2\rho\tilde{U}.$$

The first inequality holds, since if $j$ was only partially connected to $i_1$, then $i_1 \in F_j$. Also, since $F_j$ can only shrink during the algorithm and by Property (4.2f), the second inequality holds. $\square$

It remains to consider what we need to change in Theorem 4.12 for the analysis of connection cost for full clients, when we change $y^*$ to $\hat{y}$. Notice that for every $j \in C_{\text{full}}$, we have one unit of open facility in $\text{Ball}_F(j, \frac{3\tau-1}{\tau-1} D_{\ell_j})$ in $\hat{y}$. This holds due to Lemma 4.9, and the fact that $y_{i_2}^*$ is strictly smaller than 1, meaning that if $i_2$ is in such a ball, then so is $i_1$, or some other facility $i_3$ with $y_{i_3}^* = \hat{y}_{i_3} = 1$.

Thus, we only need to focus on the set $J = \{j \in C_{\text{full}} : i_2 \in B_j\}$ of clients and consider the total connection cost of these clients in the new solution. Let $\delta' = \Theta(\delta)$ be a number such that $\delta' \leq \frac{\delta}{4+3\delta} \cdot \frac{(\tau-1)(1-\delta')}{\tau(3\tau-1)}$; setting $\delta' = \frac{\delta}{6} \cdot \frac{\tau-1}{2\tau(3\tau-1)}$ satisfies the property. Let $i^*$ be the nearest open facility to $i_2$ in the solution $\hat{y}$ and let $t' = d(i_2, i^*)$. Let $J_1 = \{j \in J : d(j, i_2) \geq \delta' t'\}$ and $J_2 = J \setminus J_1 = \{j \in J : d(j, i_2) < \delta' t'\}$.

For every $j \in J_1$, we have $d(j, i^*) \leq d(j, i_2) + d(i_2, i^*) \leq d(j, i_2) + \frac{d(j, i_2)}{\delta'} = (1 + 1/\delta')d(j, i_2)$. Thus,

$$\sum_{j \in J_1} d^q(j, i^*) \leq \left(1 + \frac{1}{\delta'}\right)^q \sum_{j \in J_1} d^q(j, i_2) \leq \left(1 + \frac{1}{\delta'}\right)^q \cdot 2\rho\tilde{U} = O\left(\frac{\rho}{\delta^q}\right) U.$$

20

To see the second inequality, we know that for every $j \in J_1 \subseteq J$, we have $i_2 \in B_j \subseteq F_j$ at the end of Algorithm 2. Thus, $i_2 \in F_j$ at the beginning, and hence by Property (4.2f), we have $\sum_{j \in J_1} d^q(j, i_2) \leq 2\rho\tilde{U}$.

For every $j \in J_2$, we have $t' - d(i_2, j)$ is at most the distance from $j$ to its the nearest open facility in $\hat{y}$, which is at most $\frac{3\tau - 1}{\tau - 1} D_{\ell_j}$. Thus, $R_j \geq \frac{D_{\ell_j}}{\tau} \geq \frac{\tau - 1}{\tau(3\tau - 1)}(t' - d(i_2, j)) \geq \frac{\tau - 1}{\tau(3\tau - 1)}(1 - \delta')t'$.

Let $t = \frac{\tau - 1}{\tau(3\tau - 1)}(1 - \delta)t'$. Then for every $j \in J_2$, we have $R_j \geq t$. Also, by our choice of $\delta'$, we have $\frac{\delta t}{4 + 3\delta} \geq \delta' t'$. So,

$$|J_2| \leq \left| \left\{ j \in \mathrm{Ball}_C \left( i_2, \frac{\delta t}{4 + 3\delta} \right) : R_j \geq t \right\} \right| \leq \frac{\rho(1 + 3\delta/4)^q}{(1 - \delta/4)^q} \cdot \frac{U}{t^q} = O(\rho) \cdot \frac{U}{t^q}.$$

The second inequality is by Property (3.3a'). So,

$$\sum_{j \in J_2} d^q(j, i^*) \leq |J_2|((1 + \delta')t')^q \leq O(1) \cdot t'^q \cdot O(\rho) \cdot \frac{U}{t^q} = O(\rho)U.$$

Thus, overall, we have $\sum_{j \in J} d^q(j, i^*) \leq O(\rho/\delta^q)U$. This shows that the additional cost incurred due to changing the almost-integral solution to an integral one is at most $O(\rho/\delta^q)U$, finishing the proof of Theorem 3.4.

# 6  Improved Approximation Algorithm for Matroid Median

The Matroid Median problem (MatMed) is a generalization of $k$ median where we require the set of open facilities to be an independent set of a given matroid. The $k$ median problem is a special case when the underlying matroid is a uniform matroid of rank $k$. We now show how we can use our framework to obtain an $\alpha_1$ approximation for MatMed.

Unlike RkMed, the natural LP relaxation for MatMed only has a small integrality gap, and so we can skip the pre-processing steps and jump directly to the iterative rounding framework. To give a description in the simplest way, we just assume $U = \infty$, $U' = \tilde{U}'$ is a guessed upper bound on the cost of the optimum solution, $R_j = \infty$ for every $j \in C$ and $m = n$, at the beginning of Section 4. The only change is that $y(F) \leq k$ in (LP$_{\mathrm{strong}}$) and (LP$_{\mathrm{iter}}$) need to be replaced with the matroid polytope constraints $\sum_{v \in S} y_v \leq r_{\mathcal{M}}(S), \forall S \subseteq F$. Eventually, all clients $j$ will be moved to $C_{\mathrm{full}}$ and thus the Constraint (12) becomes redundant. The final $y^*$ is a vertex point in the intersection of a partition matroid polytope and the given matroid polytope. Thus $y^*$ is integral. This leads to an integral solution with expected cost $\alpha_1 U'$, leading to an $\alpha_1$-approximation for MatMed.

# 7  Improved approximation algorithm for Knapsack Median

In the knapsack median problem (KnapMed), we are given a knapsack constraint on the facilities instead of a bound on number of facilities opened as in $k$ median. Formally, for every facility $i$, there is a non-negative weight $w_i$ associated with it, and we have an upper bound $W$ on the total weight of facilities. The problem is to choose a vector $y \in \{0, 1\}^n$ that satisfies $\sum_{i \in F} w_i y_i \leq W$ and we minimize the total connection cost of all the clients. The classical $k$ median is a special case of knapsack median, when all the weights $w_i$ are equal to 1 and $W$ is equal to $k$. Due to the knapsack constraint, the natural LP for KnapMed has an unbounded integrality gap. We perform the similar preprocessing as for RkMed (Section 3), except now we use a simple way to give the upper bound vector $(R_j)$, instead of applying Theorem 3.3.

**Definition 7.1** (Sparse Instances for KnapMed). *Suppose we are given an extended* KnapMed *instance* $\mathcal{I} = (F, C, d, w, B, S_0)$, *and parameters* $\rho, \delta \in (0, 1/2)$ *and* $U \geq 0$. *Let* $S^*$ *be a solution to* $\mathcal{I}$ *with cost at most* $U$, *and* $(\kappa^*, c^*)$ *be the nearest-facility-vector-pair for* $S^*$. *Then we say the instance* $\mathcal{I}$ *is* $(\rho, \delta, U)$-*sparse w.r.t the solution* $S^*$, *if*

*(7.1a) for every* $i \in S^* \setminus S_0$, *we have* $\sum_{j \in C: \kappa_j^* = i} c_j^* \leq \rho U$,

*(7.1b) for every* $p \in F \cup C$, *we have* $\left| \text{Ball}_C(p, \delta c_p^*) \right| \cdot c_p^* \leq \rho U$.

Using a similar argument as in Theorem 3.2, given a KnapMed instance $\mathcal{I}$, we can produce a family of $n^{O(1/\rho)}$ extended KnapMed instances, one of which, say $\mathcal{I}' = (F, C' \subseteq C, d, w, B, S_0 \subseteq S^*)$, is $(\rho, \delta, U)$-sparse w.r.t $S^*$, the optimum solution of $\mathcal{I}$. Moreover, $\frac{1-\delta}{1+\delta} \sum_{j \in C \setminus C'} d(j, S_0) + \sum_{j \in C'} d(j, S^*) \leq U$. Therefore it suffices for us to prove the following theorem:

**Theorem 7.2.** *Let* $\mathcal{I}' = (F, C', d, w, B, S_0)$ *be a* $(\rho, \delta, U)$-*sparse instance w.r.t some solution* $S^*$, *for some* $\rho, \delta \in (0, 1/2)$ *and* $U \geq 0$. *Let* $U'$ *be the cost of* $S^*$ *for* $\mathcal{I}'$. *There is an efficient randomized algorithm that computes a solution to* $\mathcal{I}'$ *with expected cost at most* $\alpha_1 U' + O(\rho/\delta)U$.

Again, we shall replace $\mathcal{I}'$ with $\mathcal{I}$ and $C'$ with $C$. We use a very simple way to define $R_j$: for every $j$, let $R_j$ be the maximum number $R$ such that $|\text{Ball}_C(j, \delta R)| \cdot R \leq \rho U$. Let $(\kappa^*, c^*)$ be the nearest-facility-vector-pair for $S^*$. By Property (7.1b), we have, $c_j^* \leq R_j$. We then formulate a stronger LP relaxation:

$$\min \quad \sum_{i \in F, j \in C} x_{i,j} d(i, j) \quad \text{s.t.} \quad (\text{LP}_{\text{k-strong}})$$

$$y(F) \leq k \quad (14) \qquad x_{i,j} = 0 \qquad \forall i, j \text{ s.t } d(i,j) > R_j \quad (18)$$

$$x_{i,j} \leq y_i \quad \forall i,j \quad (15) \qquad x_{i,j} = 0 \qquad \forall i \in F \setminus S_0, j \text{ s.t } d(i,j) > \rho U \quad (19)$$

$$\sum_{i \in F} x_{i,j} = 1 \quad \forall j \quad (16) \qquad \sum_i d(i,j) x_{i,j} \leq \rho U y_i \qquad \forall i \in F \setminus S_0 \quad (20)$$

$$y_i = 1 \quad \forall i \in S_0 \quad (17)$$

There is a solution $(x, y)$ to the above LP with cost at most $U'$, as the indicator vector $(x, y)$ for solution $S^*$ satisfies all the above constraints. After this LP, we eliminate the $x_{i,j}$ variables akin to Section 4.2, then perform a randomized discretization of distances identical to the one in Section 4.3, and obtain a feasible solution to the following auxiliary LP. Like in Section 4, we begin with $C_{\text{part}} = C, C_{\text{full}} = C^* = \emptyset$.

Then we formulate an LP for the iterative rounding. Now we can require $y(F_j) = 1$ for every $j \in C_{\text{part}}$ and thus the coverage constraint is not needed:

$$\min \quad \sum_{j \in C_{\text{part}}} \sum_{i \in F_j} d'(i,j) y_i \quad + \quad \sum_{j \in C_{\text{full}}} \left( \sum_{i \in B_j} d'(i,j) y_i + (1 - y(B_j)) d_{\ell_j} \right) \quad \text{s.t.} \quad (\text{LP}_{\text{k-iter}})$$

$$\sum_i w_i y_i \leq W \quad (21) \qquad y(B_j) \leq 1 \qquad \forall j \in C_{\text{full}} \quad (23)$$

$$y(F_j) = 1 \qquad \forall j \in C^* \cup C_{\text{part}} \quad (22) \qquad y_i \in [0,1] \qquad \forall i \in F \quad (24)$$

We then run the same Algorithm 2 except we solve $(\text{LP}_{\text{k-iter}})$ each time instead of $(\text{LP}_{\text{iter}})$. Eventually every client will become a full client, i.e., $C_{\text{full}} = C$ and $C_{\text{part}} = \emptyset$. At the end of iterative rounding algorithm, the only tight constraints are (21), (22) and (24). Now, we use the fact that the intersection of a laminar family of

constraints and a knapsack polytope is an almost-integral polytope with at most two fractional variables. Moreover, following the same analysis as the proof of Theorem 4.12, we can bound the expected connection cost to be at most $\frac{3\tau-1}{\ln\tau}U' = \alpha_1 U'$, where the expectation is over the random choice for $a$ in discretization step.

In order to get an integral solution $\hat{y}$, we let $\hat{y} = y$ and consider following two cases. If there is exactly one strictly fractional facility, we call it $i_2$ and close it in the solution $\hat{y}$. If there are two strictly fractional facilities $i_1$ and $i_2$, then we must have $y^*_{i_1} + y^*_{i_2} = 1$. Let $i_1$ be the facility with smaller weight and $i_2$ be the one with bigger weight. We open $i_1$ and close $i_2$ in $\hat{y}$. Thus, we ensure that the knapsack constraint is not violated in the integral solution.

Unlike RkMed, there is only one type of cost that we incur during this process since all clients are fully connected. The clients that were using $i_2$ partially now need to be connected to their nearest open facility. We bound the increase of connection cost of these clients. Again, it suffices to focus on the set $J = \{j \in C_{\text{full}} : i_2 \in B_j\}$ of clients and consider the total connection cost of these clients in the new solution $\hat{y}$. In the solution $\hat{y}$, for every client $j \in C_{\text{full}} = C$, there is always a open facility with distance $D_{\ell_j}$ from $j$.

Let $\delta' = \Theta(\delta)$ be a number such that $\delta' \leq \frac{(\tau-1)(1-\delta')}{2\tau(3\tau-1)}\delta$; setting $\delta' = \frac{\tau-1}{4\tau(3\tau-1)}\delta$ satisfies the property. Let $i^*$ be the nearest open facility to $i_2$ in the solution $\hat{y}$ and let $t = d(i_2, i^*)$. Let $J_1 = \{j \in J : d(j, i_2) \geq \delta't\}$ and $J_2 = J \setminus J_1 = \{j \in J : d(j, i_2) < \delta't\}$.

For every $j \in J_1$, we have $d(j, i^*) \leq d(j, i_2) + d(i_2, i^*) \leq d(j, i_2) + \frac{d(j,i_2)}{\delta'} = (1 + 1/\delta')d(j, i_2)$. Thus,

$$\sum_{j \in J_1} d(j, i^*) \leq \left(1 + \frac{1}{\delta'}\right)\sum_{j \in J_1} d(j, i_2) \leq \left(1 + \frac{1}{\delta'}\right) \cdot 2\rho U = O\left(\frac{\rho}{\delta}\right)U.$$

To see the second inequality, we know that for every $j \in J_1 \subseteq J$, we have $i_2 \in B_j \subseteq F_j$ in the end of Algorithm 2. Thus, $i_2 \in F_j$ after the splitting. By the property that each star has small cost, we have $\sum_{j \in J_1} d(j, j_2) \leq 2\rho U$.

Assuming $J_2 \neq \emptyset$. Fix any $j^* \in J_2$. $t - \delta't \leq t - d(j^*, i_2) = d(i_2, i^*) - d(j^*, i_2) \leq d(j^*, i^*) \leq \frac{3\tau-1}{\tau-1}D_{\ell_j} \leq \frac{(3\tau-1)\tau}{\tau-1}R_j$. Thus $R_j \geq \frac{(1-\delta')(\tau-1)}{\tau(3\tau-1)}t$. By our choice of $\delta'$, we have $\delta R_j \geq 2\delta't$. So,

$$|J_2| \leq |\text{Ball}_C(j^*, \delta R_j)| \leq \frac{\rho U}{R_j}.$$

So,

$$\sum_{j \in J_2} d(j, i^*) \leq |J_2|(1 + \delta')t \leq \frac{\rho U}{R_j} \cdot O(1)R_j = O(\rho)U.$$

Thus, overall, we have $\sum_{j \in J} d(j, i^*) \leq O(\rho/\delta)U$. This shows that the additional cost incurred due to changing the almost-integral solution to an integral one is at most $O(\rho/\delta)U$, finishing the proof of Theorem 7.2.

# References

[1] Sara Ahmadian, Ashkan Norouzi-Fard, Ola Svensson, and Justin Ward. Better guarantees for k-means and euclidean k-median by primal-dual algorithms. *Proceedings, IEEE Symposium on Foundations of Computer Science (FOCS)*, abs/1612.07925, 2017.

[2] Sanjeev Arora, Prabhakar Raghavan, and Satish Rao. Approximation schemes for euclidean k-medians and related problems. In *Proceedings of STOC*, STOC '98, pages 106–113, New York, NY, USA, 1998. ACM.

[3] David Arthur, Bodo Manthey, and Heiko Röglin. Smoothed analysis of the k-means method. *J. ACM*, 58(5):19:1–19:31, 2011.

[4] David Arthur and Sergei Vassilvitskii. K-means++: The advantages of careful seeding. In *Proceedings of ACM-SIAM SODA 2007*.

[5] Vijay Arya, Naveen Garg, Rohit Khandekar, Adam Meyerson, Kamesh Munagala, and Vinayaka Pandit. Local search heuristic for k-median and facility location problems. In *Proceedings of STOC 2001*.

[6] Pranjal Awasthi, Avrim Blum, and Or Sheffet. Stability yields a PTAS for k-median and k-means clustering. In *Proceedings of FOCS 2010*, pages 309–318. IEEE Computer Society, 2010.

[7] Maria-Florina Balcan, Avrim Blum, and Anupam Gupta. Clustering under approximation stability. *J. ACM*, 60(2):8:1–8:34, 2013.

[8] Jaroslaw Byrka. An optimal bifactor approximation algorithm for the metric uncapacitated facility location problem. In *APPROX/RANDOM 2007, Princeton, NJ, USA, Proceedings*, pages 29–43, 2007.

[9] Jaroslaw Byrka, Thomas Pensyl, Bartosz Rybicki, Joachim Spoerhase, Aravind Srinivasan, and Khoa Trinh. An improved approximation algorithm for knapsack median using sparsification. In *Proceedings of ESA 2015*, pages 275–287, 2015.

[10] Jaroslaw Byrka, Thomas Pensyl, Bartosz Rybicki, Aravind Srinivasan, and Khoa Trinh. An improved approximation for k-median and positive correlation in budgeted optimization. *ACM Trans. Algorithms*, 13(2):23:1–23:31, March 2017.

[11] M. Charikar, S. Guha, D. Shmoys, and E. Tardos. A constant-factor approximation algorithm for the k-median problem. *ACM Symp. on Theory of Computing (STOC)*, 1999.

[12] M. Charikar, S. Khuller, D. M. Mount, and G. Narasimhan. Algorithms for facility location problems with outliers. *Proceedings, ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2001.

[13] Moses Charikar and Sudipto Guha. Improved combinatorial algorithms for the facility location and k-median problems. In *Proceedings of FOCS 1999*.

[14] Moses Charikar and Shi Li. A dependent lp-rounding approach for the k-median problem. In *Proceedings of ICALP 2012*.

[15] Sanjay Chawla and Aristides Gionis. k-means–: A unified approach to clustering and outlier detection. In *Proceedings of the 13th SIAM International Conference on Data Mining, May 2-4, 2013. Austin, Texas, USA.*, pages 189–197, 2013.

[16] Ke Chen. A constant factor approximation algorithm for k-median clustering with outliers. In *Proceedings of ACM-SIAM SODA 2008*.

[17] Fabián A. Chudak and David B. Shmoys. Improved approximation algorithms for the uncapacitated facility location problem. *SIAM J. Comput.*, 33(1):1–25, 2003.

[18] V. Cohen-Addad and C. Schwiegelshohn. On the Local Structure of Stable Clustering Instances. *Proceedings, IEEE Symposium on Foundations of Computer Science (FOCS)*, October 2017.

[19] Vincent Cohen-Addad, Philip N. Klein, and Claire Mathieu. The power of local search for clustering. *Proceedings, IEEE Symposium on Foundations of Computer Science (FOCS)*, abs/1603.09535, 2016.

[20] Zachary Friggstad, Kamyar Khodamoradi, Mohsen Rezapour, and Mohammad R. Salavatipour. Approximation schemes for clustering with outliers. *Proceedings, ACM-SIAM Symposium on Discrete Algorithms (SODA)*, abs/1707.04295, 2018.

[21] Zachary Friggstad, Mohsen Rezapour, and Mohammad R. Salavatipour. Local search yields a PTAS for k-means in doubling metrics. *Proceedings, IEEE Symposium on Foundations of Computer Science (FOCS)*, abs/1603.08976, 2016.

[22] Sudipto Guha and Samir Khuller. Greedy strikes back: Improved facility location algorithms. In *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '98, pages 649–657, Philadelphia, PA, USA, 1998. Society for Industrial and Applied Mathematics.

[23] Shalmoli Gupta, Ravi Kumar, Kefu Lu, Benjamin Moseley, and Sergei Vassilvitskii. Local search methods for k-means with outliers. *Proceedings, International Conference on Very Large Data Bases (VLDB)*, 10(7):757–768, March 2017.

[24] M. Hajiaghayi, R. Khandekar, and G. Kortsarz. Local search algorithms for the red-blue median problem. *Algorithmica*, 63(4):795–814, Aug 2012.

[25] K. Jain and V. V. Vazirani. Approximation algorithms for metric facility location and k-median problems using the primal-dual schema and lagrangian relaxation. *J. ACM*, 48(2):274 – 296, 2001.

[26] Kamal Jain, Mohammad Mahdian, Evangelos Markakis, Amin Saberi, and Vijay V. Vazirani. Greedy facility location algorithms analyzed using dual fitting with factor-revealing lp. *J. ACM*, 50(6):795–824, November 2003.

[27] Kamal Jain, Mohammad Mahdian, and Amin Saberi. A new greedy approach for facility location problems. In *Proceedings of STOC 2002*.

[28] Madhukar R. Korupolu, C. Greg Plaxton, and Rajmohan Rajaraman. Analysis of a local search heuristic for facility location problems. In *Proceedings of ACM-SIAM SODA 1998*, pages 1–10.

[29] Ravishankar Krishnaswamy, Amit Kumar, Viswanath Nagarajan, Yogish Sabharwal, and Barna Saha. The matroid median problem. In *Proceedings of ACM-SIAM SODA 2011*.

[30] Amit Kumar. Constant factor approximation algorithm for the knapsack median problem. In *Proceedings of the Twenty-third Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '12, pages 824–832, Philadelphia, PA, USA, 2012. Society for Industrial and Applied Mathematics.

[31] Amit Kumar and Ravindran Kannan. Clustering with spectral norm and the k-means algorithm. In *Proceedings of FOCS 2010*.

[32] Euiwoong Lee, Melanie Schmidt, and John Wright. Improved and simplified inapproximability for k-means. *Inf. Process. Lett.*, 120:40–43, 2017.

[33] S. Li and O. Svensson. Approximating k-median via pseudo-approximation. *ACM Symp. on Theory of Computing (STOC)*, 2013.

[34] Shi Li. *A 1.488 Approximation Algorithm for the Uncapacitated Facility Location Problem*, pages 77–88. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011.

[35] Jyh-Han Lin and Jeffrey Scott Vitter. Approximation algorithms for geometric median problems. *Inf. Process. Lett.*, 44(5):245–249.

[36] S. Lloyd. Least squares quantization in pcm. *IEEE Trans. Inf. Theor.*, 28(2):129–137, September 2006.

[37] Mohammad Mahdian, Yinyu Ye, and Jiawei Zhang. Approximation algorithms for metric facility location problems. *SIAM J. Comput.*, 36(2):411–432, 2006.

[38] Rafail Ostrovsky, Yuval Rabani, Leonard J. Schulman, and Chaitanya Swamy. The effectiveness of lloyd-type methods for the k-means problem. *J. ACM*, 59(6):28:1–28:22, 2012.

[39] Lionel Ott, Linsey Pang, Fabio T Ramos, and Sanjay Chawla. On integrated clustering and outlier detection. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 1359–1367. 2014.

[40] N. Rujeerapaiboon, K. Schindler, D. Kuhn, and W. Wiesemann. Size Matters: Cardinality-Constrained Clustering and Outlier Detection via Conic Optimization. *ArXiv e-prints*, May 2017.

[41] David B. Shmoys, Éva Tardos, and Karen Aardal. Approximation algorithms for facility location problems (extended abstract). In *Proceedings of STOC 1997*.

[42] Chaitanya Swamy. Improved approximation algorithms for matroid and knapsack median problems and applications. *ACM Trans. Algorithms*, 12(4):49:1–49:22, August 2016.

[43] David P. Williamson and David B. Shmoys. *The Design of Approximation Algorithms*. Cambridge University Press, New York, NY, USA, 1st edition, 2011.