# Local Search Yields a PTAS for $k$-Means in Doubling Metrics

Zachary Friggstad
*Department of Computing Science*
*University of Alberta*
*Email:zachary@ualberta.ca*

Mohsen Rezapour
*Department of Computing Science*
*University of Alberta*
*Email:rezapour@ualberta.ca*

Mohammad R. Salavatipour
*Department of Computing Science*
*University of Alberta*
*Email:mrs@ualberta.ca*

*Abstract*—The most well known and ubiquitous clustering problem encountered in nearly every branch of science is undoubtedly $k$-MEANS: given a set of data points and a parameter $k$, select $k$ centres and partition the data points into $k$ clusters around these centres so that the sum of squares of distances of the points to their cluster centre is minimized. Typically these data points lie in Euclidean space $\mathbb{R}^d$ for some $d \geq 2$. $k$-MEANS and the first algorithms for it were introduced in the 1950's. Over the last six decades, hundreds of papers have studied this problem and different algorithms have been proposed for it. The most commonly used algorithm in practice is known as Lloyd-Forgy, which is also referred to as "the" $k$-MEANS algorithm, and various extensions of it often work very well in practice. However, they may produce solutions whose cost is arbitrarily large compared to the optimum solution. Kanungo et al. [2004] analyzed a very simple local search heuristic to get a polynomial-time algorithm with approximation ratio $9+\epsilon$ for any fixed $\epsilon > 0$ for $k$-MEANS in Euclidean space. Finding an algorithm with a better worst-case approximation guarantee has remained one of the biggest open questions in this area, in particular whether one can get a true PTAS for fixed dimension Euclidean space.

We settle this problem by showing that a simple local search algorithm provides a PTAS for $k$-MEANS for $\mathbb{R}^d$ for any fixed $d$. More precisely, for any error parameter $\epsilon > 0$, the local search algorithm that considers swaps of up to $\rho = d^{O(d)} \cdot \epsilon^{-O(d/\epsilon)}$ centres at a time will produce a solution using *exactly* $k$ centres whose cost is at most a $(1+\epsilon)$-factor greater than the optimum solution. Our analysis extends very easily to the more general settings where we want to minimize the sum of $q$'th powers of the distances between data points and their cluster centres (instead of sum of squares of distances as in $k$-MEANS) for any fixed $q \geq 1$ and where the metric may not be Euclidean but still has fixed doubling dimension.

*Keywords*-clustering, $k$-means, approximation, Euclidean, Polynomial Time Approximation Scheme, doubling metric

## I. INTRODUCTION

With advances in obtaining and storing data, one of the emerging challenges of our age is data analysis. It is hard to find a scientific research project which does not involve some form of methodology to process, understand, and summarize data. Two problems often encounterd in data analysis are classification and clustering. Classification (which is an instance of supervised learning) is the task of predicting the label of a new data point after being trained with a set of labeled data points (called training set). Clustering (which is an instance of unsupervised learning) is the task of grouping a given set of objects or data points into clusters/groups such that the data points that are more similar fall into the same cluster while data points (objects) that do not seem similar are in different clusters. Some of the main purposes of clustering are to understand the underlying structure and relation between objects and find a compact representation of data points.

Clustering and different methods to achieve it have been studied since 1950's in different branches of science. Perhaps the most widely used clustering model is the $k$-MEANS clustering: Given a set $\mathcal{X}$ of $n$ *data points* in $d$-dimensional Euclidean space $\mathbb{R}^d$, and an integer $k$, find a set of $k$ points $c_1, \ldots, c_k \in \mathbb{R}^d$ to act as as *centres* that minimize the sum of squared distances of each data point to its nearest centre. In other words, we would like to partition $\mathcal{X}$ into $k$ cluster sets, $\{C_1, \ldots, C_k\}$ and find a centre $c_i$ for each $C_i$ to minimize $\sum_{i=1}^{k} \sum_{x \in C_i} ||x - c_i||_2^2$. Here, $||x - c_i||_2$ is the standard Euclidean distance in $\mathbb{R}^d$ between points $x$ and $c_i$. The sum is called the cost of the clustering. Typically, the centres $c_i$ are selected to be the centroid (mean) of the cluster $C_i$. In other situations the centres must be from the data points themselves (i.e. $c_i \in C_i$) or from a given set $\mathcal{C}$. This latter version is referred to as discrete $k$-MEANS clustering. Although in most application of $k$-MEANS the data points are in some Euclidean space, the discrete variant can be defined in general metrics. The $k$-MEANS clustering problem is known to be an NP-hard problem even for $k = 2$ or when $d = 2$ [1], [40], [21], [45]. Clustering, in particular the $k$-MEANS clustering problem as the most popular model for it, has found numerous applications in very different areas (see [31]).

### A. Previous work

The most widely used algorithm for $k$-MEANS (which is also sometimes referred to as "the" $k$-means algorithm) is a simple heuristic introduced by Lloyd in 1957 [39]. Although this algorithm works well in practice it is known that it has unbounded approximation ratio (see [32]). Various modifications and extensions of this algorithm have been produced and studied (see [31]), but none of them are known to have a bounded approximation ratio in the general setting. Arthur and Vassilvitskii [6] show that Lloyd's method with properly chosen initial centres will be an $O(\log k)$-approximation.

Ostrovsky et al. [43] show that under some assumptions about the data points the approximation ratio is bounded by a constant. The problem of finding an efficient algorithm for $k$-MEANS with a proven theoretical bound on the cost of the solution returned is probably one of the most well studied problems in the whole field of clustering with hundreds of research papers devoted to this.

Arthur and Vassilvitskii [5], Dastupta and Gupta [19], and Har-Peled and Sadri [29] study convergence rate of Lloyd's algorithm. In particular [5] show that it can be super-polynomial. More recently, Vattani [46] shows that it can take exponential time even in two dimensions. Arthur et al. [4] proved that Lloyd's algorithm has polynomial-time smoothed complexity. Kumar and Kannan [34], Ostrovsky et al. [43], and Awasthi et al. [9] gave empirical and theoretical evidence for when and why the known heuristics work well in practice. For instance [9] show that when the size of an optimum $(k-1)$-MEANS is sufficiently larger than the cost of $k$-MEANS then one can get a near optimum solution to $k$-MEANS using a variant of Lloyd's algorithm.

The $k$-MEANS problem is known to be NP-hard [1], [21], [40]. In fact, the $k$-MEANS problem is NP-hard if $d$ is arbitrary even for $k = 2$ [1], [21]. Also, if $k$ is arbitrary the problem is NP-hard even for $d = 2$ [40], [45]. However, the $k$-MEANS problem can be solved in polynomial time by the algorithm of [30] when both $k$ and $d$ are constant.

Matoušek [41] gave a PTAS for fixed $k, d$ with running time $O(n(\log n)^k \epsilon^{-2k^2 d})$. Since then several other PTASs have been proposed for variant settings of parameters but all need $k$ to be constant [42], [11], [20], [28], [35], [36], [23], [27]. Recently, Bandyapadhyay and Varadarajan [12] presented a pseudo-approximation for $k$-MEANS: their algorithm finds a solution whose cost is at most $1 + \epsilon$ times of the optimum but might use up to $(1 + \epsilon) \cdot k$ clusters. The result of Matoušek [41] also shows that one can select a set $\mathcal{C}$ of "candidate" centers in $\mathbb{R}^d$ from which the $k$ centres should be chosen from with a loss of at most $(1 + \epsilon)$ and this set can be computed in time $O(n\epsilon^{-d} \log(1/\epsilon))$. This reduces the $k$-means problem to the discrete setting where along with $\mathcal{X}$ we have a set $\mathcal{C}$ of candidate centres and we have to select $k$ centres from $\mathcal{C}$.

Kanungo et al. [32] proved that a simple local search heuristic yields an algorithm with approximation ratio $9 + \epsilon$ for $\mathbb{R}^d$. This remains the best known approximation algorithm with polynomial running time for $\mathbb{R}^d$. For general metrics, Gupta and Tangwongsan [26] proved that local search is a $(25 + \epsilon)$-approximation. It was an open problem for a long time whether $k$-MEANS is APX-hard or not in Euclidean metrics. This was recently answered positively by Awasthi et al. [10] where they showed that the problem is APX-hard in $\mathbb{R}^d$ but the dimension $d$ used in the proof is $\Omega(\log n)$. Blomer et al. [13] have a nice survey of theoretical analysis of different $k$-means algorithms.

$k$-MEDIAN: Another very well studied problem that is also closely related to $k$-MEANS is $k$-MEDIAN. The only difference is that the goal (objective function) in $k$-MEDIAN is to minimize the sum of distances, instead of sum of square of distances as in $k$-MEANS, i.e. minimize $\sum_{i=1}^{k} \sum_{x \in C_i} \delta(x, c_i)$ where $\delta(x, c_i)$ is the distance between $x$ and $c_i$.

The simple local search (which swaps in and out a constant number of centres in each iteration) is known to give a $3 + \epsilon$ approximation for general metrics by Arya et al. [7], [8]. The current best approximation uses different techniques and has an approximation ratio of $2.611 + \epsilon$ [38], [14]. The local search $(9 + \epsilon)$-approximation (for $k$-MEANS) in [32] can be seen as an extension of the analysis in [7] for $k$-MEDIAN. One reason that analysis of $k$-MEANS is more difficult is that the squares of distances do not necessarily satisfy the triangle inequality. For instances of $k$-MEDIAN on Euclidean metrics, Arora et al. [3], building on the framework of Arora [2], gave the first PTAS. Kolliopoulos and Rao [33] improved the time complexity.

### B. Our result and technique

Although a PTAS for $k$-MEDIAN in fixed dimension Euclidean space has been known for almost two decades, getting a PTAS for $k$-MEANS in fixed dimension Euclidean space has remained an open problem. We provide a PTAS for this setting. We focus on the discrete case where we have to select a set of $k$ centres from a given set $\mathcal{C}$. Our main result is to show that a simple local search heuristic that swaps up to $\rho = d^{O(d)} \cdot \epsilon^{-O(d/\epsilon)}$ centres at a time and assigns each point to the nearest centre is a PTAS for $k$-MEANS in $\mathbb{R}^d$.

Recall that the *doubling dimension* of a metric space is the smallest $\tau$ such that any ball of radius $2r$ around a point can be covered by at most $2^\tau$ balls of radius $r$. A *doubling metric* is one with constant doubling dimension (as in [44]). Our analysis implies a PTAS for more general settings where the data points are in a metric space with constant doubling dimension (described below) and when the objective function is to minimize the sum of $q$'th power of the distances for some fixed $q \geq 1$. Let $\rho(\epsilon, d) := d^{O(d)} \cdot \epsilon^{O(-d/\epsilon)}$. We will articulate the absolute constants suppressed by the $O(\cdot)$ notation later on in our analysis.

*Theorem 1:* The local search algorithm that swaps up to $\rho(\epsilon, d)$ centres at a time is a $(1 + \epsilon)$-approximation for $k$-MEANS in metrics with doubling dimension $d$.

Now consider the generalization where the objective function measures the sum of $q$'th power of the distances between points and their assigned cluster centre, we call this $\ell_q^q$-NORM $k$-CLUSTERING Here, we are given the points $\mathcal{X}$ in a metric space $\delta(\cdot, \cdot)$ along with a set $\mathcal{C}$ of potential centres. We are to select $k$ centres from $\mathcal{C}$ and partition the points into $k$ cluster sets $C_1, \ldots, C_k$ with each $C_i$ having a centre $c_i$ so as to minimize $\sum_{i=1}^{k} \sum_{x \in C_i} \delta(x, c_i)^q$. Note

that the case of $q = 2$ is the $k$-MEANS problem and $q = 1$ is $k$-MEDIAN.

Our analysis extends to provide a PTAS for this setting when $q$ is fixed. That is, Theorem 1 holds for $\ell_q^q$-NORM $k$-CLUSTERING for constants $q \geq 1$, except that we require that the local search procedure be the $\rho'$-swap heuristic where $\rho' = d^{O(d)} \cdot (2^q/\epsilon)^{O(2^q \cdot d/\epsilon)}$.

Note that even for the case of $k$-MEDIAN, this is the first PTAS for metrics with constant doubling dimension. Also, while a PTAS was known for $k$-MEDIAN for constant-dimensional Euclidean metrics, determining if local search provided such a PTAS was an open problem.

The analysis of Theorem 1 is easily adapted to prove that local search yields a PTAS in doubling metrics for UNCAPACITATED FACILITY LOCATION with *non-uniform* opening costs and the common generalization of $k$-MEDIAN and UNCAPACITATED FACILITY LOCATION where facilities have opening costs and we are only allowed to open $k$ such facilities. Prior to this, [18] showed that local search yields a PTAS for UNCAPACITATED FACILITY LOCATION with uniform opening costs in constant-dimensional Euclidean metrics. A PTAS for Euclidean UNCAPACITATED FACILITY LOCATION was first given by Arora et al. [3].

*Theorem 2:* The local search algorithm that adds and subtracts up to $\rho(\epsilon, d)$ facilities is a $(1+\epsilon)$-approximation for UNCAPACITATED FACILITY LOCATION with non-uniform opening costs in metrics with doubling dimension $d$. Similarly, the local search algorithm that swaps $\rho(\epsilon, d)$ facilities is a $(1 + \epsilon)$-approximation for the generalization of UNCAPACITATED FACILITY LOCATION where we must open precisely $k$ facilities.

As mentioned earlier, Awasthi et al. [10] proved that $k$-MEANS is APX-hard for $d = \Omega(\log n)$ and they left the approximability of $k$-MEANS for lower dimensions as an open problem. A consequence of our algorithm is that one can get a $(1 + \epsilon)$-approximation for $k$-MEANS that runs in sub-exponential time for values of $d$ up to $O(\log n / \log \log n)$. More specifically, for any given $0 < \epsilon < 1$ and $d = \sigma \log n / \log \log n$ for sufficiently small absolute constant $\sigma$ we get a $(1 + \epsilon)$-approximation for $k$-MEANS that runs in time $O(2^{n^\kappa})$, for some constant $\kappa = \kappa(\sigma) < 1$; for $d = O(\log \log n / \log \log \log n)$ we get a *quasi-polytime approximation scheme* (QPTAS). Therefore, our result in a sense shows that the requirement of [10] of $d = \Omega(\log n)$ to prove APX-hardness of $k$-MEANS is almost tight unless $\mathrm{NP} \subseteq DTIME(2^{n^\sigma})$.

The notion of coresets and using them for finding faster algorithms for $k$-means has been studied extensively (e.g. [28], [27], [15], [23], [24] and references there). In the full version, we discuss how to use these ideas to improve the running time of the local search algorithm.

**Note:** Shortly after we announced our result [25], Cohen-Addad, Klein, and Mathieu [17], [16] announced similar results for Euclidean and minor-free metrics using the local search method.

### C. Proof Outline

The general framework for analysis of our local search algorithms for is similar to that of $k$-MEDIAN and $k$-MEANS in [8], [32], [26]. Let $\mathcal{S}$ and $\mathcal{O}$ be a local optimum and a global optimum solution, respectively. We carefully identify a set $Q$ of potential swaps between local and global optimum. In each such swap, the cost of assigning a data point $x$ to the nearest centre after a swap is bounded with respect to the local and global cost assignment. In other words, if $\mathcal{B}(\mathcal{S})$ is the set of solutions that can be obtained by performing swaps from $Q$, the main task is to show that $\sum_{\mathcal{S}' \in \mathcal{B}(\mathcal{S})}(\mathrm{cost}(\mathcal{S}') - \mathrm{cost}(\mathcal{S})) \leq (1 + O(\epsilon) \cdot \mathrm{cost}(\mathcal{O}) - (1 - O(\epsilon)\mathrm{cost}(\mathcal{S})$. Given that $0 \leq \mathrm{cost}(\mathcal{S}') - \mathrm{cost}(\mathcal{S})$ for all $\mathcal{S}' \in \mathcal{B}(\mathcal{S})$ (because $\mathcal{S}$ is a local optimum), $\mathrm{cost}(\mathcal{S}) \leq \frac{1+O(\epsilon)}{1-O(\epsilon)} \cdot \mathrm{cost}(\mathcal{O})$.

Our analysis has many more ingredients and several intermediate steps to get us what we want. Note that the following only describes steps used in the analysis of the local search algorithm; we do not perform any of the steps described below in the algorithm itself. Let us define $\mathcal{S}$ and $\mathcal{O}$ as before. First, we do a filtering over $\mathcal{S}$ and $\mathcal{O}$ to obtain subsets $\overline{\mathcal{S}} \subseteq \mathcal{S}$ and $\overline{\mathcal{O}} \subseteq \mathcal{O}$ such that every centre in $\mathcal{S} - \overline{\mathcal{S}}$ (in $\mathcal{O} - \overline{\mathcal{O}}$) is "close" to a centre in $\overline{\mathcal{S}}$ (in $\overline{\mathcal{O}}$) while these filtered centres are far apart. We define a "net" around each centre $i \in \overline{\mathcal{S}}$ which captures a collection of other filtered centres in $\overline{\mathcal{O}}$ that are relatively close to $i$. The idea of the net is that if we choose to close $i$ (in a test swap) then the data points that were to be assigned to $i$ will be assigned to a nearby centre in the net of $i$. Since the metric is a constant-dimensional Euclidean metric (or, more generally, a doubling metric), we can choose these nets to have constant size.

For each $j$ assigned to $i$ in $\mathcal{S}$, if the centre $i^*$ that $j$ is assigned to in the optimum solution lies somewhat close to $i$ then we can reassign $j$ to a facility in the net around $i$ that is close to $i^*$. In this case, the reassignment cost for $j$ will be close to $c_j^* - c_j$. Otherwise, if $i^*$ lies far from $i$ then we can reassign $j$ to a facility near $i$ in the net around $i$ and the reassignment cost will only be $O(\epsilon) \cdot (c_j^* + c_j)$ and we will generate the $c_j^* - c_j$ term for the local search analysis when $i^*$ is opened in another different swap.

One main part of our proof is to show that there exists a suitable randomized partitioning of $\mathcal{S} \cup \mathcal{O}$ such that each part has small size (which will be a function of only $\epsilon$ and $d$ and, ultimately, determines the size of the swaps in our local search procedure), and for any pair $(i, i^*) \in \overline{\mathcal{S}} \times \overline{\mathcal{O}}$ where $i^*$ lies in the net of $i$ we have $\Pr[i, i^* \text{ lie in the same part}] \geq 1 - \epsilon$. This randomized partitioning is the only part of the proof that we rely on properties of doubling metrics (or $\mathbb{R}^d$). For those small portion of facilities that their net is "cut" by our partitioning, we show that when we close those centres in our test swaps then the reassignment cost of a point $j$ that was assigned to them is only $O(1)$ times more than the sum

of their assignment costs in $\mathcal{O}$ and $\mathcal{S}$. Given that this only happens with small probability (due to our random partition scheme), this term is negligible in the final analysis. So we get an overall bound of $1 + O(\epsilon)$ on the ratio of cost of $\mathcal{S}$ over $\mathcal{O}$.

**Outline of the paper:** We start with some basic definitions and notation in Section II. In Sections III and IV we show that the local search algorithm with an appropriate number of swaps yields a PTAS for $k$-MEANS in $\mathbb{R}^d$. The extension to doubling metrics and to the setting where we measure the $\ell_q^q$-norm of the solution for any constant $q \geq 1$, the proof of Theorem 2, as well as many proofs of the supporting results used to prove Theorem 1 can be found in the full version [25].

## II. NOTATION AND PRELIMINARIES

Recall that in the $k$-MEANS problem we are given a set $\mathcal{X}$ of $n$ points in $\mathbb{R}^d$ and an integer $k \geq 1$; we have to find $k$ centres $c_1, \ldots, c_k \in \mathbb{R}^d$ so as to minimize the sum of squares of distances of each point to the nearest centre. As mentioned earlier, by using the result of [41], at a loss of $(1+\epsilon)$ factor we can assume we have a set $\mathcal{C}$ of "candidate" centres from which the centres can be chosen from. This set can be computed in time $O(n\epsilon^{-d}\log(1/\epsilon))$ and $|\mathcal{C}| = O(n\epsilon^{-d}\log(1/\epsilon))$. Therefore, we can reduce the problem to the discrete case. Formally, suppose we are given a set $\mathcal{C}$ of points (as possible cluster centres) along with $\mathcal{X}$ and we have to select the $k$ centres from $\mathcal{C}$. Furthermore, we assume the points are given in a metric space $(V, \delta)$ (not necessarily $\mathbb{R}^d$). For any two points $p, q \in V$, $\delta(p, q)$ denotes the distance between them: for the case of the metric being $\mathbb{R}^d$, then $\delta(p, q) = (\sum_{\ell=1}^d |p_\ell - q_\ell|^2)^{\frac{1}{2}}$.

We usually refer to a potential centre in $\mathcal{C}$ by a simple index $i$ and a point in $\mathcal{X}$ by a simple index $j$ (or slight variants like $i^*$ or $\overline{i'}$). This is to emphasize that we do not need to talk about specific coordinates of points in Euclidean space. In fact, only once in our proof do we rely on the particular embedding of the points in Euclidean space. So, for any set $S \subseteq \mathcal{C}$ and any $j \in \mathcal{X}$, let $\delta(j, S) = \min_{i \in S} \delta(j, i)$. We also define $\text{cost}(S) = \sum_{j \in \mathcal{X}} \delta(j, S)^2$.

Our goal in (discrete) $k$-MEANS is to find a set of centres $S \subseteq \mathcal{C}$ of size $k$ to minimize $\text{cost}(S)$. Note that once we fix the set of centres we can find a partitioning of $\mathcal{X}$ that realizes $\sum_{j \in \mathcal{X}} \delta(j, S)^2$ by assigning each $j \in \mathcal{X}$ to the nearest centre in $S$, breaking ties arbitrarily.

As mentioned, for ease of exposition we focus on $k$-MEANS on $\mathbb{R}^d$ The following simple $\rho$-swap local search heuristic (Algorithm 1) is essentially the same one considered in [32].

Throughout we assume that $\epsilon > 0$ is sufficiently small. Recall that we defined $\rho(\epsilon, d) = d^{O(d)} \cdot \epsilon^{O(d/\epsilon)}$, where the constants will be specified later and consider the local search algorithm with $\rho = \rho(\epsilon, d)$ swaps. By a standard argument (as in [7], [32]) one can show that replacing the condition of

the while loop with $\text{cost}((\mathcal{S} - Q) \cup P) \leq (1 - \frac{\epsilon}{k}) \cdot \text{cost}(\mathcal{S})$, the algorithm terminates in polynomial time and we only loos a $(1 + \epsilon)$ factor in the approximation ratio. For ease of exposition, we ignore this factor $1 + \epsilon$ loss, and consider the solution $\mathcal{S}$ returned by Algorithm 1. Recall that we use $\mathcal{O}$ to denote the global optimum solution. For $j \in \mathcal{X}$, let $c_j^* = \delta(j, \mathcal{O})^2$ and $c_j = \delta(j, \mathcal{S})^2$, so $\text{cost}(\mathcal{O}) = \sum_{j \in \mathcal{X}} c_j^*$ and $\text{cost}(\mathcal{S}) = \sum_{j \in \mathcal{X}} c_j$. We also denote the centre in $\mathcal{O}$ nearest to $j$ by $\sigma^*(j)$ and the centre in $\mathcal{S}$ nearest to $j$ by $\sigma(j)$. Define $\phi : \mathcal{O} \cup \mathcal{S} \to \mathcal{O} \cup \mathcal{S}$ to be the function that assigns $i^* \in \mathcal{O}$ to its nearest centre in $\mathcal{S}$ and assigns $i \in \mathcal{S}$ to its nearest centre in $\mathcal{O}$. For any two sets $S, T \subseteq \mathcal{O} \cup \mathcal{S}$, we let $S \triangle T = (S \cup T) - (S \cap T)$. We assume $\mathcal{O} \cap \mathcal{S} = \emptyset$. This is without loss of generality because we could duplicate each location in $\mathcal{C}$ and say $\mathcal{O}$ uses the originals and $\mathcal{S}$ the duplicates. It is easy to check that $\mathcal{S}$ would still be a locally optimum solution in this instance. We can also assume that these are the only possible colocated facilities, so $\delta(i, i') > 0$ for distinct $i, i' \in \mathcal{O}$ or distinct $i, i' \in \mathcal{S}$. Finally, we will assume $\epsilon$ is sufficiently small (independent of all other parameters, including $d$) so that all of our bounds hold.

## III. LOCAL SEARCH ANALYSIS FOR $\mathbb{R}^d$

In this section we focus on $\mathbb{R}^d$ (for fixed $d \geq 2$) and define $\rho(\epsilon, d) = 32 \cdot (2d)^{8d} \cdot \epsilon^{-36 \cdot d/\epsilon}$. Our goal in this section is to prove that the $\rho$-swap local search with $\rho = \rho(\epsilon, d)$ is a PTAS for $k$-MEANS in $\mathbb{R}^d$.

*Theorem 3:* Let $\mathcal{S}$ be a locally-optimum solution with respect to the $\rho(\epsilon, d)$-swap local search heuristic when the points lie in $\mathbb{R}^d$. Then $\text{cost}(\mathcal{S}) \leq (1 + O(\epsilon)) \cdot \text{cost}(\mathcal{O})$.

To prove this, we will construct a set of test swaps that yield various inequalities which, when combined, provide the desired bound on $\text{cost}(\mathcal{S})$. That is, we will partition $\mathcal{O} \cup \mathcal{S}$ into sets where $|P \cap \mathcal{O}| = |P \cap \mathcal{S}| \leq \rho(\epsilon, d)$ for each part $P$. For each such set $P$, $0 \leq \text{cost}(\mathcal{S} \triangle P) - \text{cost}(\mathcal{S})$ because $\mathcal{S}$ is a locally optimum solution. We will provide an explicit upper bound on this cost change that will reveal enough information to easily conclude $\text{cost}(\mathcal{S}) \leq (1 + O(\epsilon)) \cdot \text{cost}(\mathcal{O})$. For example, for a point $j \in \mathcal{X}$ if $\sigma^*(j) \in P$ then the change in $j$'s assignment cost is at most $c_j^* - c_j$ because we could assign $j$ from $\sigma(j)$ to $\sigma^*(j)$. The problem is that points $j$ with $\sigma(j) \in P$ but $\sigma^*(j) \notin P$ must go somewhere else; most of our effort is ensuring that the test swaps are carefully chosen so such reassignment cost increases are very small.

First we need to describe the partition of $\mathcal{O} \cup \mathcal{S}$. This is a fairly elaborate scheme that involves several steps. As mentioned earlier, the actual algorithm for $k$-MEANS is the simple local search we described and the algorithms we describe below to get this partitioning scheme are only for the purpose of proof and analysis of the local search algorithm.

*Definition 1:* For $i^* \in \mathcal{O}$ let $D_{i^*} := \delta(i^*, \mathcal{S}) = \delta(i^*, \phi(i^*))$. For $i \in \mathcal{S}$ let $D_i := \delta(i, \mathcal{O}) = \delta(i, \phi(i))$.

---
**Algorithm 1** $\rho$-Swap Local Search
---
Let $\mathcal{S}$ be an arbitrary set of $k$ centres from $\mathcal{C}$
**while** $\exists$ sets $P \subseteq \mathcal{C} - \mathcal{S}$, $Q \subseteq \mathcal{S}$ with $|P| = |Q| \le \rho$ s.t. $\mathrm{cost}((\mathcal{S} - Q) \cup P) < \mathrm{cost}(\mathcal{S})$ **do**
　　$\mathcal{S} \leftarrow (\mathcal{S} - Q) \cup P$
**return** $\mathcal{S}$

---

The first thing is to sparsify $\mathcal{O}$ and $\mathcal{S}$ using a simple filtering step. In particular,

*Lemma 1:* There exists $\overline{\mathcal{S}} \subseteq \mathcal{S}$, $\overline{\mathcal{O}} \subseteq \mathcal{O}$ and a *proxy* function $\eta : \mathcal{S} \cup \mathcal{O} \to \overline{\mathcal{S}} \cup \overline{\mathcal{O}}$ mapping $\mathcal{S}$ into $\overline{\mathcal{S}}$ and $\mathcal{O}$ into $\overline{\mathcal{O}}$ with the following properties: a) $\eta(i) = i$ for each $i \in \overline{\mathcal{S}} \cup \overline{\mathcal{O}}$; b) $\delta(i, \eta(i)) \le \epsilon \cdot D_i$ for each $i \in \mathcal{O} \cup \mathcal{S}$; c) $\delta(i, i') \ge \epsilon \cdot \max\{D_i, D_{i'}\}$ for distinct $i, i' \in \overline{\mathcal{S}} \cup \overline{\mathcal{O}}$, d) $D_{\eta(i)} \le D_i$ for each $i \in \mathcal{S} \cup \mathcal{O}$.

The proof follows a simple modification of the standard clustering technique that is used in many UNCAPACITATED FACILITY LOCATION and $k$-MEDIAN LP rounding algorithms. Details are in the full version of this paper.

Next we define mappings similar to $\phi, \sigma, \sigma^*$ except they only concern centres that were not filtered out. Let $\overline{\phi} : \overline{\mathcal{O}} \cup \overline{\mathcal{S}} \to \overline{\mathcal{O}} \cup \overline{\mathcal{S}}$ map each $i \in \overline{\mathcal{O}}$ to its nearest location in $\overline{\mathcal{S}}$ and vice versa; $\overline{\sigma}^* : \mathcal{X} \to \overline{\mathcal{O}}$ defined by $\overline{\sigma}^*(j) = \eta(\sigma^*(j))$, and $\overline{\sigma} : \mathcal{X} \to \overline{\mathcal{S}}$ defined by $\overline{\sigma}(j) = \eta(\sigma(j))$. Finally, for each $i \in \overline{\phi}(\overline{\mathcal{O}})$, let $\mathrm{cent}(i)$ be the centre in $\phi^{-1}(i)$ that is closest to $i$, breaking ties arbitrarily. Note $\overline{\sigma}(j)$ may not necessarily be the centre in $\overline{\mathcal{S}}$ that is closest to $j$. Also note that if one considers a bipartite graph with parts $\overline{\mathcal{O}}$ and $\overline{\mathcal{S}}$, then $\overline{\phi}$ maps centres from one side to the other.

*Lemma 2:* For each $i' \in \overline{\mathcal{O}} \cup \overline{\mathcal{S}}$, $D_{i'} \le \delta(i', \overline{\phi}(i')) \le (1 + \epsilon) \cdot D_{i'}$.

Finally, the last definition in this section identifies pairs of centres that we would like to have in the same part of the partition we construct. Finally we define

$$\mathcal{T} := \{(\mathrm{cent}(i), i) : i \in \overline{\phi}(\overline{\mathcal{O}}) \text{ and } \epsilon \cdot \delta(\mathrm{cent}(i), i) \le D_i\}$$
$$\mathcal{N} := \{(i^*, i) \in \overline{\mathcal{O}} \times \overline{\mathcal{S}} : \delta(i, i^*) \le \epsilon^{-1} \cdot D_i \text{ and } D_{i^*} \ge \epsilon \cdot D_i\}.$$

For each $i \in \overline{\mathcal{S}}$, the set $\{i^* : (i^*, i) \in \mathcal{N}\}$ is the "net" for centre $i$ that was discussed in the proof outline in Section I-C. Ultimately we will require that pairs in $\mathcal{T}$ are not separated by the partition. Our requirement for $\mathcal{N}$ is not quite as strong. The partition is constructed randomly and it will be sufficient to have each pair in $\mathcal{N}$ being separated by the partition with probability at most $\epsilon$.

The following lemma says that if at least one centre of each pair in $\mathcal{T}$ is open after a swap, then every centre in $\overline{\mathcal{O}} \cup \overline{\mathcal{S}}$ is somewhat close to some open centre. The bound is a bit big, but it will be multiplied by $O(\epsilon)$ whenever it is used in the local search analysis. The main idea of the proof is that for each $i' \in \mathcal{O} \cup \mathcal{S}$ we consider the sequence $i' = i_0, i_1, i_2, \dots$ with $i_{a+1} = \overline{\phi}(i_a)$ for each $a \ge 0$. Every second step decreases in length geometrically as long as no point of a pair in $\mathcal{T}$ is hit, which also shows it eventually reaches a point $i_{a'}$ with either $i_{a'} \in A$ or $\mathrm{cent}(i_{a'}) \in A$.

*Lemma 3:* Let $A \subseteq \overline{\mathcal{O}} \cup \overline{\mathcal{S}}$ be such that $A \cap \{\mathrm{cent}(i), i\} \ne \emptyset$ for each $(\mathrm{cent}(i), i) \in \mathcal{T}$. Then $\delta(i', A) \le 5 \cdot D_{i'}$ for any $i' \in \mathcal{O} \cup \mathcal{S}$.

### A. Good Partitioning of $\mathcal{O} \cup \mathcal{S}$ and Proof of Theorem 3

The main tool used in our analysis is the existence of the following randomized partitioning scheme.

*Theorem 4:* There is a randomized algorithm that samples a partition $\pi$ of $\mathcal{O} \cup \mathcal{S}$ such that:

- For each part $P \in \pi$, $|P \cap \mathcal{O}| = |P \cap \mathcal{S}| \le \rho(\epsilon, d)$.
- For each part $P \in \pi$, $\mathcal{S} \triangle P$ includes at least one centre from every pair in $\mathcal{T}$.
- For each $(i^*, i) \in \mathcal{N}$,
  $\Pr[i, i^* \text{ lie in different parts of } \pi] \le \epsilon$.

We prove this theorem in Section IV. For now, we will complete the analysis of the local search algorithm using this partitioning scheme. Note that in the following we do not use the geometry of the metric (i.e. all arguments hold for general metrics); it is only in the proof of Theorem 4 that we use properties of $\mathbb{R}^d$. The following simple lemma gives a way to handle the fact that the triangle inequality does not hold with squares of the distances.

*Lemma 4:* For any real numbers $x, y$ we have $(x+y)^2 \le 2(x^2 + y^2)$.

*Lemma 5:* For each point $j \in \mathcal{X}$, $D_{\overline{\sigma}(j)} \le D_{\sigma(j)} \le \delta(j, \sigma(j)) + \delta(j, \sigma^*(j))$. Similarly, $D_{\overline{\sigma}^*(j)} \le D_{\sigma^*(j)} \le \delta(j, \sigma(j)) + \delta(j, \sigma^*(j))$.

*Proof:* We prove the first statement, the second is nearly identical. That $D_{\overline{\sigma}(j)} \le D_{\sigma(j)}$ follows from d) in Lemma 1. For the other inequality, note $D_{\sigma(j)} = \delta(\sigma(j), \phi(\sigma(j))) \le \delta(\sigma(j), \sigma^*(j)) \le \delta(j, \sigma(j)) + \delta(j, \sigma^*(j))$. ∎

### Proof of Theorem 3.

Let $\pi$ be a partition sampled by the algorithm from Theorem 4. For each point $j \in \mathcal{X}$ and each part $P$ of $\pi$, let $\Delta_j^P := \delta(j, \mathcal{S} \triangle P)^2 - \delta(j, \mathcal{S})^2$ denote the change in assignment cost for the point after swapping $P$. Local optimality of $\mathcal{S}$ and $|P \cap \mathcal{S}| = |P \cap \mathcal{O}| \le \rho(\epsilon, d)$ means $0 \le \sum_j \Delta_j^P$ for any part $P$. Classify each point $j \in \mathcal{X}$ as follows:

- **Lucky**: $\sigma(j)$ and $\overline{\sigma}(j)$ do not lie in the same part of $\pi$.

- **Long**: $j$ is not lucky but $\delta(\overline{\sigma}(j), \overline{\sigma}^*(j)) > \epsilon^{-1} \cdot D_{\overline{\sigma}(j)}$.

- **Bad**: $j$ is not lucky or long and $(\overline{\sigma}^*(j), \overline{\sigma}(j)) \in \mathcal{N}$ yet $\overline{\sigma}(j), \overline{\sigma}^*(j)$ lie in different parts of $\pi$.

- **Good**: $j$ is neither lucky, long, nor bad.

We now place a bound on $\sum_{P \in \pi} \Delta_j^P$ for each point $j \in \mathcal{X}$. Note that each centre in $\mathcal{S}$ is swapped out exactly once over all swaps $P$ and each centre in $\mathcal{O}$ is swapped in exactly once. With this in mind, consider the following cases for a point $j \in \mathcal{X}$. In the coming arguments, we let $\delta_j := \delta(j, \sigma(j))$ and $\delta_j^* := \delta(j, \sigma^*(j))$ for brevity. Note $c_j = \delta_j^2$ and $c_j^* = \delta_j^{*2}$.

In all cases for $j$ except when $j$ is bad, the main idea is that we can bound the distance from $j$ to some point in $\mathcal{S} \triangle P$ by first moving it to either $\sigma(j)$ or $\sigma^*(j)$ and then moving it a distance of $O(\epsilon) \cdot (\delta_j + \delta_j^*)$ to reach an open facility. Considering that we reassigned $j$ from $\sigma(j)$, the reassignment cost will be $(\delta_j + O(\epsilon) \cdot (\delta_j + \delta_j^*))^2 - c_j = O(\epsilon) \cdot (\delta_j + \delta_j^*)$ or $(\delta_j^* + O(\epsilon) \cdot (\delta_j + \delta_j^*))^2 - c_j = (1 + O(\epsilon)) \cdot c_j^* - (1 - O(\epsilon)) \cdot c_j$. Full details for the first case are provided here, more details for the remaining cases are in the full version of this paper.

● **Case: $j$ is lucky**
For the part $P \in \pi$ with $\sigma^*(j) \in P$, we have $\Delta_j^P \leq c_j^* - c_j$ as we could move $j$ from $\sigma(j)$ to $\sigma^*(j)$. If $\sigma(j)$ is swapped out in a different swap $P'$, we move $j$ to $\overline{\sigma}(j)$ (which remains open because $j$ is lucky) and bound $\Delta_j^{P'}$ by:

$$\begin{aligned}
\Delta_j^{P'} &\leq \delta(j, \overline{\sigma}(j))^2 - \delta_j^2 \\
&\leq (\delta_j + \delta(\sigma(j), \overline{\sigma}(j)))^2 - c_j \\
&\leq (\delta_j + \epsilon \cdot D_{\sigma(j)})^2 - c_j \quad \text{(Lemma 1.b)} \\
&= 2\epsilon \cdot \delta_j \cdot D_{\sigma(j)} + \epsilon^2 \cdot D_{\sigma(j)}^2 \\
&\leq 2\epsilon \cdot \delta_j \cdot (\delta_j + \delta_j^*) + \epsilon^2 \cdot (\delta_j + \delta_j^*)^2 \quad \text{(Lemma 5)} \\
&\leq 2\epsilon \cdot (\delta_j + \delta_j^*)^2 + \epsilon^2 \cdot (\delta_j + \delta_j^*)^2 \\
&\leq 4\epsilon \cdot (c_j^* + c_j) + 2\epsilon^2 \cdot (c_j^* + c_j) \quad \text{(Lemma 4)} \\
&\leq 6\epsilon \cdot (c_j^* + c_j). \quad \text{(for $\epsilon$ sufficiently small)}
\end{aligned}$$

For every other swap $P''$, we have that $\sigma(j)$ remains open after the swap so $\Delta_j^{P''} \leq 0$ as we could just leave $j$ at $\sigma(j)$. In total, we have $\sum_{P \in \pi} \Delta_j^P \leq c_j^* - c_j + 6\epsilon \cdot (c_j^* + c_j)$.

● **Case: $j$ is long**
Again, for $P \in \pi$ with $\sigma^*(j) \in P$ we get $\Delta_j^P \leq c_j^* - c_j$. If $\sigma(j)$ is swapped out in a different swap $P'$, then we bound $\Delta_j^{P'}$ by moving $j$ from $\sigma(j)$ to the open centre nearest to $\sigma(j)$. Note that $\mathcal{S} \triangle P'$ contains at least one centre from every pair in $\mathcal{T}$, so we can bound this distance using Lemma 3. Using this, we bound $\Delta_j^{P'}$ as follows.

$$\begin{aligned}
\Delta_j^{P'} &\leq (\delta_j + \delta(\sigma(j), \mathcal{S} \triangle P))^2 - c_j \\
&\leq (\delta_j + 5D_{\sigma(j)})^2 - c_j \quad \text{(Lemma 3)}
\end{aligned}$$

We can bound $D_{\sigma(j)}$ by $2\epsilon(1 + \epsilon)(\delta_j + \delta_j^*)$, which leads to $\Delta_j^{P'} \leq 46\epsilon \cdot (c_j + c_j^*)$. In every other swap we could leave $j$ at $\sigma(j)$. Thus, for a long point $j$ we have $\sum_{P \in \pi} \Delta_j^P \leq c_j^* - c_j + 46\epsilon \cdot (c_j^* + c_j)$.

● **Case: $j$ is bad**
We only move $j$ in the swap $P$ when $\sigma(j)$ is closed. In this case, we assign $j$ in the same way as if it was long.

Our bound is weaker here and introduces significant positive dependence on $c_j$. This will eventually be compensated by the fact that $j$ is bad with probability at most $\epsilon$ over the random choice of $\pi$. For now, we just provide the reassignment cost bound for bad $j$.

$$\begin{aligned}
\Delta_j^P &\leq (\delta_j + \delta(\sigma(j), \mathcal{S} \triangle P))^2 - c_j \\
&\leq (\delta_j + 5D_{\sigma(j)})^2 \quad \text{(Lemma 3)} \\
&\leq (6\delta_j + 5\delta_j^*)^2 \leq 72 \cdot (c_j^* + c_j),
\end{aligned}$$

where we used Lemma 3 for the 2nd and Lemma 4 for the 3rd inequality. Therefore, $\sum_{P \in \pi} \Delta_j^P \leq 72 \cdot (c_j^* + c_j)$.

● **Case: $j$ is good**
This breaks into two subcases. We know $\delta(\overline{\sigma}^*(j), \overline{\sigma}(j)) \leq \epsilon^{-1} \cdot D_{\overline{\sigma}(j)}$ because $j$ is not long. In one subcase, $D_{\overline{\sigma}^*(j)} \geq \epsilon \cdot D_{\overline{\sigma}(j)}$ so $(\overline{\sigma}^*(j), \overline{\sigma}(j)) \in \mathcal{N}$. Since $j$ is not bad and not lucky, we have $\sigma(j), \overline{\sigma}(j), \overline{\sigma}^*(j) \in P$ for some common part $P \in \pi$. In the other subcase, $D_{\overline{\sigma}^*(j)} < \epsilon \cdot D_{\overline{\sigma}(j)}$. Note in this case we still have $\sigma(j), \overline{\sigma}(j) \in P$ for some common part $P$ because $j$ is not lucky.

**Subcase: $D_{\overline{\sigma}^*(j)} \geq \epsilon \cdot D_{\overline{\sigma}(j)}$**
The only time we move $j$ is when $\sigma(j)$ is closed. As observed in the previous paragraph, this happens in the same swap when $\overline{\sigma}^*(j)$ is opened, so send $j$ to $\overline{\sigma}^*(j)$.

$$\begin{aligned}
\Delta_j^P &\leq (\delta_j^* + \delta(\sigma^*(j), \overline{\sigma}^*(j)))^2 - c_j \\
&\leq c_j^* - c_j + 5\epsilon \cdot (c_j^* + c_j).
\end{aligned}$$

**Subcase: $D_{\overline{\sigma}^*(j)} < \epsilon \cdot D_{\overline{\sigma}(j)}$**
Again, the only time we move $j$ is when $\sigma(j)$ is closed. We reassign $j$ by first moving it to $\sigma^*(j)$ and then using Lemma 3 to further bound the cost. We bound the cost change for this reassignment as follows. Recall that $D_{\sigma^*(j)} \leq (1 + \epsilon)D_{\overline{\sigma}^*(j)}$ which in turn is bounded by $\epsilon(1 + \epsilon)D_{\overline{\sigma}(j)}$ in this subcase (by the assumption of subcase). Thus,

$$\begin{aligned}
\Delta_j^P &\leq (\delta_j^* + \delta(\sigma^*(j), \mathcal{S} \triangle P))^2 - c_j \\
&\leq c_j^* - c_j + 24\epsilon \cdot (c_j^* + c_j).
\end{aligned}$$

Considering both subcases we can say that for any good point $j$ that $\sum_{P \in \pi} \Delta_j^P \leq c_j^* - c_j + 24\epsilon \cdot (c^* + c)$. Aggregating these bounds and remembering $0 \leq \sum_{j \in \mathcal{X}} \Delta_j^P$ for each $P \in \pi$ because $\mathcal{S}$ is a locally optimum solution, we have

$$\begin{aligned}
0 &\leq \sum_{j \in \mathcal{X}} \sum_{P \in \pi} \Delta_j^P \\
&\leq \sum_{\substack{j \in \mathcal{X} \\ j \text{ not bad}}} \left[(1 + 46\epsilon)c_j^* - (1 - 46\epsilon)c_j\right] + \sum_{\substack{j \in \mathcal{X} \\ j \text{ bad}}} 72(c_j^* + c_j).
\end{aligned}$$

The last step is to average this inequality over the random choice of $\pi$. Note that any point $j \in \mathcal{X}$ is bad with probability at most $\epsilon$ by the guarantee in Theorem 4 and

the definition of *bad*. Thus, we see

$$
\begin{aligned}
0 \;\leq\; & \mathbf{E}_\pi\left[\sum_{j\in\mathcal{X}}\sum_{P\in\pi}\Delta_j^P\right] \\
\leq\; & \sum_{j\in\mathcal{X}}\Pr[j\text{ not bad}]\cdot\left[(1+46\epsilon)c_j^*-(1-46\epsilon)c_j\right]+ \\
& \qquad \Pr[j\text{ bad}]\cdot 72(c_j^*+c_j) \\
\leq\; & \sum_{j\in\mathcal{X}}(1+118\epsilon)c_j^*-(1-119\epsilon)c_j.
\end{aligned}
$$

Rearranging shows $\mathrm{cost}(\mathcal{S})\leq\frac{1+118\epsilon}{1-119\epsilon}\cdot\mathrm{cost}(\mathcal{O})$. ∎

## IV. THE PARTITIONING SCHEME: PROOF OF THEOREM 4

Here is the overview of our partitioning scheme. First we bucket the real values by defining bucket $a$ to have real values $r$ where $\epsilon^{-a}\leq r<\epsilon^{-(a+1)}$. Fix $b=\Theta(1/\epsilon)$ and then choose a random off-set $a'$ and consider $b$ consecutive buckets $a'+\ell\cdot b,\ldots,a'+(\ell+1)\cdot b-1$ to form "bands". We get a random partition of points $i\in\overline{\mathcal{O}}\cup\overline{\mathcal{S}}$ into bands $B_1,B_2,\ldots$ such that all centres $i$ in the same band $B_\ell$ have the property that their $D_i$ values are in buckets $a'+\ell\cdot b,\ldots,a'+(\ell+1)\cdot b-1$ (so $\epsilon^{-(a'+\ell\cdot b)}\leq D_i<\epsilon^{-(a'+(\ell+1)\cdot b-1)}$). Given that for each pair $(i,i^*)\in\mathcal{N}$ the $D_i,D_{i^*}$ values are within a factor $1/\epsilon$ of each other, the probability that they fall into different bands is small (like $\epsilon/4$). Next we impose a random grid partition on each band $B_\ell$ to cut it into cells with cell width roughly $\Theta(d/\epsilon^{\ell\cdot b})$; again the randomness comes from choosing a random off-set for our grid. The randomness helps to bound the probability of any pair $(i,i^*)\in\mathcal{N}$ being cut into two different cells to be small again. This will ensure the 3rd property of the theorem holds. Furthermore, since each $i\in B_\ell$ has $D_i\geq\epsilon^{-(a'+\ell\cdot b)}$ and by the filtering we did (Lemma 1), balls of radius $\frac{\epsilon}{2}\cdot\epsilon^{-(a'+\ell\cdot b)}$ around them must be disjoint; hence using a simple volume/packing argument we can bound the number of centres from each $B_\ell$ that fall into the same grid cell by a function of $\epsilon$ and $d$. This helps of establish the first property. We have to do some clean-up to ensure that for each pair $(i,\mathrm{cent}(i))\in\mathcal{T}$ they belong to the same part (this will imply property 2 of the Theorem) and that at the end for each part $P$, $|P\cap\mathcal{S}|=|P\cap\mathcal{O}|$. These are a bit technical and are described in details below.

We start by geometrically grouping the centres in $\overline{\mathcal{O}}\cup\overline{\mathcal{S}}$. For $a\in\mathbb{Z}$, let $G_a:=\left\{i\in\overline{\mathcal{O}}\cup\overline{\mathcal{S}}:\frac{1}{\epsilon^a}\leq D_i<\frac{1}{\epsilon^{a+1}}\right\}$. Finally, let $G_{-\infty}:=\{i\in\overline{\mathcal{O}}\cup\overline{\mathcal{S}}:D_i=0\}$. Note that each $i\in\overline{\mathcal{O}}\cup\overline{\mathcal{S}}$ appears in exactly one set among $\{G_a:a\in\mathbb{Z}\}\cup\{G_{-\infty}\}$.

We treat $G_{-\infty}$ differently in our partitioning algorithm. It is important to note that no pair in $\mathcal{T}$ or $\mathcal{N}$ has precisely one point in $G_{-\infty}$, as the following shows.

*Lemma 6:* For each pair of centres $(i^*,i)\in\mathcal{T}\cup\mathcal{N}$, $|\{i,i^*\}\cap G_{-\infty}|\neq 1$.

*Proof:* Consider some colocated pair $(i^*,i)\in\mathcal{S}\times\mathcal{O}$. As $D_i=D_{i^*}=0$ and because no other centre in $\mathcal{S}\cup\mathcal{O}$ is colocated with $i$ and $i^*$, then $i\in\overline{\mathcal{S}}$ and $i^*\in\overline{\mathcal{O}}$. Thus, $\overline{\phi}(i^*)=i$ and it is the unique closest facility in $\overline{\mathcal{O}}$ to $i$ so $i^*=\mathrm{cent}(i)$. This shows every pair $(i^*,i)\in\mathcal{T}$ with either $D_i=0$ or $D_{i^*}=0$ must have both $D_i=D_{i^*}=0$ so $i^*,i\in G_{-\infty}$.

Next consider some $(i^*,i)\in\mathcal{N}$. We know $\epsilon\cdot D_i\leq D_{i^*}$ by definition of $\mathcal{N}$. Thus, if $i^*\in G_{-\infty}$ then $i\in G_{-\infty}$ as well. Conversely, suppose $i\in G_{-\infty}$. Since $(i^*,i)\in\mathcal{N}$ then $\delta(i,i^*)\leq\epsilon^{-1}\cdot D_i=0$. So, $D_{i^*}=0$ meaning $i^*\in G_{-\infty}$ as well. ∎

### A. Partitioning $\{G_a:a\in\mathbb{Z}\}$

**Step 1: Forming Bands**

Fix $b$ to be the smallest integer that is at least $4/\epsilon$. We first partition $\{G_a:a\in\mathbb{Z}\}$ into *bands*. Sample an integer *shift* $a'$ uniformly at random in the range $\{0,1,\ldots,b-1\}$. For each $\ell\in\mathbb{Z}$, form the band $B_\ell$ as follows:

$$
B_\ell:=\bigcup_{0\leq j\leq b-1}G_{a'+j+\ell\cdot b}.
$$

**Step 2: Cutting Out Cells**

Focus on a band $B_\ell$. Let $W_\ell:=4d\cdot\epsilon^{-(\ell+2)\cdot b-2}$ (recall that $d$ is the dimension of the Euclidean metric). This will be the "width" of the cells we create. It is worth mentioning that this is the first time we are going to use the properties of Euclidean metrics as all our arguments so far were treating $\delta(\cdot,\cdot)$ as a general metric.

We consider a random grid with cells having width $W_\ell$ in each dimension. Choose an *offset* $\beta\in\mathbb{R}^d$ uniformly at random from the cube $[0,W_\ell]^d$. For emphasis, we let $\mathbf{p}_j^i$ refer to the $j$'th component of $i$'s point in Euclidean space (this is the only time we will refer specifically to the coordinates for a point). So centre $i$ has Euclidean coordinates $(\mathbf{p}_1^i,\mathbf{p}_2^i,\ldots,\mathbf{p}_d^i)$.

For any tuple of integers $\mathbf{a}\in\mathbb{Z}^d$ define the *cell* $C_\mathbf{a}^\ell$ as follows.

$$
C_\mathbf{a}^\ell=\{i\in B_\ell:\beta_j+W_\ell\cdot\mathbf{a}_j\leq\mathbf{p}_j^i<\beta_j+W_\ell\cdot(\mathbf{a}_j+1)\text{ for all }1\leq j\leq d\}.
$$

These are all points in the band $B_\ell$ that lie in the half-open cube with side length $W_\ell$ and lowest corner $\beta+W_\ell\cdot\mathbf{a}$. Each $i\in B_\ell$ lies in precisely one cell as these half-open cubes tile $\mathbb{R}^d$.

**Step 3: Fixing $\mathcal{T}$**

Let $\mathcal{I}\subseteq\overline{\mathcal{O}}$ be the centres of the form $\mathrm{cent}(i)$ for some $(\mathrm{cent}(i),i)\in\mathcal{T}$ where $i,\mathrm{cent}(i)\notin G_{-\infty}$ and $i$ and $\mathrm{cent}(i)$ lie in different cells after step 2. We will simply move each $i\in\mathcal{I}$ to the cell containing $\overline{\phi}(i)$. More precisely, for $\ell\in\mathbb{Z}$ and $\mathbf{a}\in\mathbb{Z}$ we define the *part* $P_\mathbf{a}^\ell$ as

$$
P_\mathbf{a}^\ell=(C_\mathbf{a}^\ell-\mathcal{I})\cup\{i\in\mathcal{I}:\overline{\phi}(i)\in C_\mathbf{a}^\ell\}.
$$

We will show that the parts $P_{\mathbf{a}}^\ell$ essentially satisfy most of the desired properties stated about the random partitioning scheme promised in Theorem 4. It is easy to see that property 2 holds for each part $P_{\mathbf{a}}^\ell$ since for each pair $(\mathrm{cent}(i), i) \in \mathcal{T}$ both of them belong to the same part. All that remains is to ensure they are *balanced* (i.e. have include the same number of centres in $\mathcal{O} \cup \mathcal{S}$) and to incorporate $G_{-\infty}$ and the centres in $\mathcal{O} \cup \mathcal{S}$ that were filtered out. These are relatively easy cleanup steps.

### B. Properties of the Partitioning Scheme

We will show the parts formed in the partitioning scheme so far have constant size (depending only on $\epsilon$ and $d$) and also that each pair in $\mathcal{N}$ is cut with low probability.

*Lemma 7:* For each $\ell \in \mathbb{Z}$ and $\mathbf{a} \in \mathbb{Z}^d$ we have $|P_{\mathbf{a}}^\ell| \le 2 \cdot (2d)^{2d} \cdot \epsilon^{-9d/\epsilon}$.

*Proof:* Note that each centre $i \in C_{\mathbf{a}}^\ell$ witnesses the inclusion of at most one additional centre of $\mathcal{I}$ into $P_{\mathbf{a}}^\ell$ (in particular, $\mathrm{cent}(i)$). So, $|P_{\mathbf{a}}^\ell| \le 2 \cdot |C_{\mathbf{a}}^\ell|$ meaning it suffices to show $|C_{\mathbf{a}}^\ell| \le (2d)^{2d} \cdot \epsilon^{-9d/\epsilon}$. For each $i \in C_{\mathbf{a}}^\ell$ we have $D_i \ge \epsilon^{-\ell \cdot b}$. By Lemma 1, the balls of radius $\frac{\epsilon}{2} \cdot \epsilon^{-\ell \cdot b}$ around the centres in $C_{\mathbf{a}}^\ell$ must be disjoint subsets of $\mathbb{R}^d$.

The volume of a ball of radius $R$ in $\mathbb{R}^d$ can be lower bounded by $\left(\frac{\sqrt{2\pi}}{d}\right)^d \cdot R^d$, so the total volume of all of these balls of radius $\frac{\epsilon}{2} \cdot \epsilon^{-\ell \cdot b}$ is at least $|C_{\mathbf{a}}^\ell| \cdot \left(\frac{\sqrt{\pi} \cdot \epsilon}{\sqrt{2} \cdot d}\right)^d \cdot \epsilon^{-d \cdot \ell \cdot b}$.

On the other hand, these balls are contained in a cube with side length $W_\ell + \epsilon^{-\ell \cdot b + 1}$. Thus, the total volume of all balls is at most $(W_\ell + \epsilon^{-\ell \cdot b + 1})^d \le (5d\epsilon^{-(2b+2)})^d \cdot \epsilon^{-d \cdot \ell \cdot b}$. Combining this with the lower bound on the total volume shows $|C_{\mathbf{a}}^\ell| \le (2d)^{2d} \cdot \epsilon^{-9d/\epsilon}$. ∎

*Lemma 8:* For each $(i^*, i) \in \mathcal{N}$ with $i, i^* \notin G_{-\infty}$ we have $\Pr[i, i^*$ lie in different parts$] \le \epsilon$.

*Proof:* We first bound the probability they lie in different bands. Note that $D_{i^*} \le \delta(i, i^*) \le D_i/\epsilon$ and also $\epsilon D_i \le D_{i^*}$ because $(i^*, i) \in \mathcal{N}$. Thus, if we say $i \in G_a$ and $i^* \in G_{a^*}$ then $|a - a^*| \le 1$. The probability that $G_a$ and $G_{a^*}$ are separated when forming the bands $B_\ell$ is at most $1/b \le \epsilon/4$.

Next, conditioned on the event that $i, i^*$ lie in the same band $B_\ell$, we bound the probability they lie in different cells. Say $i, i^* \in B_\ell$ and note

$$\delta(i, i^*) \le \epsilon^{-1} \cdot D_i \le \epsilon^{-(\ell+2) \cdot b - 1} \le \frac{\epsilon}{4d} \cdot W_\ell.$$

The probability that $i$ and $i^*$ are cut by the random offset along one of the $d$-dimensions is then at most $\frac{\epsilon}{4d}$. Taking the union bound over all dimensions, the probability that $i$ and $i^*$ lie in different cells is at most $\frac{\epsilon}{4}$.

Finally, we bound the probability that $i^*$ will be moved to a different part when fixing $\mathcal{T}$. If $(i^*, \overline{\phi}(i^*)) \notin \mathcal{T}$ or if $i^* \notin \mathcal{I}$ then this will not happen. So, we will bound $\Pr[i^* \in \mathcal{I}]$ if $(i^*, \overline{\phi}(i^*)) \in \mathcal{T}$. That is, we bound the probability that $i', i^*$ lie in different parts where $i'$ is such that $i^* = \mathrm{cent}(i')$ and $(i^*, i') \in \mathcal{T}$. Note that $D_{i'} \le \delta(i', i^*) \le (1 + \epsilon) \cdot D_{i^*} \le$

$\epsilon^{-1} \cdot D_{i^*}$ by Lemma 2 and $D_{i'} \ge \epsilon \cdot \delta(i', i^*) = \epsilon \cdot D_{i^*}$ by definition of $\mathcal{T}$. So $i'$ and $i^*$ lie in different bands with probability at most $\frac{\epsilon}{4}$ by the same argument as with $(i, i^*)$. Similarly, conditioned on $i', i^*$ lying in the same band, the probability they lie in different cells is at most $\frac{\epsilon}{4}$ again by the same arguments as with $(i, i^*)$.

Note that if $i, i^*$ lie in different parts then they were cut by the random band or by the random tiling, or else $i', i^*$ were cut by the random band or the random tiling (the latter may not apply if $i^*$ is not involved in a pair in $\mathcal{T}$). By the union bound, $\Pr[i, i^*$ are in different parts$] \le \epsilon$. ∎

### C. Balancing the Parts

We have partitioned $\{G_a : a \in \mathbb{Z}\}$ into parts $P_{\mathbf{a}}^\ell$ for various $\ell \in \mathbb{Z}$ and $\mathbf{a} \in \mathbb{Z}^d$. We extend this to a partition of all of $\mathcal{O} \cup \mathcal{S}$ into constant-size parts that are balanced between $\mathcal{O}$ and $\mathcal{S}$ to complete the proof of Theorem 4. Let $\mathcal{P} = \{P_{\mathbf{a}}^\ell : \ell \in \mathbb{Z}, \mathbf{a} \in \mathbb{Z}^d \text{ and } P_{\mathbf{a}}^\ell \ne \emptyset\}$ be the collection of nonempty parts formed so far.

The proof of Lemma 6 shows that $G_{-\infty}$ partitions naturally into colocated centres. So let $\mathcal{P}_{-\infty}$ denote the partition of $G_{-\infty}$ into these pairs. Finally let $\mathcal{P}'$ denote the partition of $(\mathcal{O} - \overline{\mathcal{O}}) \cup (\mathcal{S} - \overline{\mathcal{S}})$ into singleton sets.

We summarize important properties of $\mathcal{P} \cup \mathcal{P}_{-\infty} \cup \mathcal{P}'$.

- $\mathcal{P} \cup \mathcal{P}_{-\infty} \cup \mathcal{P}'$ is itself a partitioning of $\mathcal{O} \cup \mathcal{S}$ into parts with size at most $2 \cdot (2d)^{2d} \cdot \epsilon^{-9d/\epsilon}$ (Lemma 7).
- Every $(\mathrm{cent}(i), i) \in \mathcal{T}$ pair has both endpoints in the same part.
- Over the randomized formation of $\mathcal{P}$, each $(i^*, i) \in \mathcal{N}$ has endpoints in different parts with probability $\le \epsilon$.

From now on, we simply let $\overline{\mathcal{P}} = \mathcal{P} \cup \mathcal{P}_{-\infty} \cup \mathcal{P}'$. We finish by combining parts of $\overline{\mathcal{P}}$ into constant-size parts that are also balanced between $\mathcal{O}$ and $\mathcal{S}$. Since merging parts does not destroy the property of two centres lying together, we will still have that pairs of centres in $\mathcal{T}$ appear together and any pair of centres in $\mathcal{N}$ lie in different parts with probability at most $\epsilon$. For any subset $A \subseteq \mathcal{O} \cup \mathcal{S}$, let $\mu(A) = |A \cap \mathcal{O}| - |A \cap \mathcal{S}|$ denote the *imbalance* of $A$.

*Lemma 9:* Let $Y \ge 1$ be an integer and $\mathcal{A}$ a collection of disjoint, nonempty subsets of $\mathcal{O} \cup \mathcal{S}$ such that $\sum_{A \in \mathcal{A}} \mu(A) = 0$. If $|A| \le Y$ for each $A$ then there is some nonempty $\mathcal{B} \subseteq \mathcal{A}$ where $|\mathcal{B}| \le 2Y^3$ such that $\sum_{B \in \mathcal{B}} \mu(B) = 0$.

The proof uses simple counting and the pigeonhole principle and is deferred to the full version.

## V. Conclusion

We have presented a PTAS for $k$-MEANS in constant-dimensional Euclidean metrics that works, more generally, in doubling metrics and when the objective is to minimizing the $\ell_q^q$-norm of distances between points and their nearest centres for any constant $q \ge 1$. This is also the first approximation for $k$-MEDIAN in doubling metrics and the first demonstration that local search yields a true PTAS for $k$-MEDIAN even in the Euclidean plane.

The analysis extends to prove that a simple local-search algorithm for UNCAPACITATED FACILITY LOCATION yields a PTAS in doubling metrics even if the opening costs are nonuniform. Briefly, this is because the test swaps in the analysis is a partitioning of $\mathcal{S} \cup \mathcal{O}$. Note we do not need to "balance" the parts to swap the same number of facilities (i.e. Lemma 9 is not necessary for analyzing the UNCAPACITATED FACILITY LOCATION local search). The full details appear in [25].

The running time of a single step of the local search algorithm is $O(k^\rho)$ where $\rho = d^{O(d)} \cdot \epsilon^{-O(d/\epsilon)}$ for the case of $k$-MEANS when the metric has doubling dimension $d$. We have not tried to optimize the constants in the $O(\cdot)$ notations in $\rho$. The dependence on $d$ cannot be improved much under the Exponential Time Hypothesis (**ETH**): the hardness results for Euclidean $k$-MEANS in [10] shows that the running time of any PTAS must essentially depend *doubly exponentially* in $d$ under **ETH**.

It may still be possible to obtain an EPTAS for any constant dimension $d$. For example, there may be a PTAS with running time of the form $g(\epsilon, d) \cdot n^{O(1)}$ for some function. Finally, what is the fastest PTAS we can obtain in the special case of the Euclidean plane (i.e. $d = 2$)? It would be interesting to see if there is an EPTAS whose running time is linear or near linear in $n$ for any fixed constant $\epsilon$.

## ACKNOWLEDGMENT

## REFERENCES

[1] Daniel Aloise, Amit Deshpande, Pierre Hansen, and Preyas Popat. NP-hardness of Euclidean sum-of-squares clustering. *Mach. Learn.*, 75(2):245–248, 2009.

[2] Sanjeev Arora. Polynomial time approximation schemes for Euclidean traveling salesman and other geometric problems. *J. ACM*, 45(5):753–782, 1998.

[3] Sanjeev Arora, Prabhakar Raghavan, and Satish Rao. Approximation schemes for Euclidean K-medians and related problems. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing (STOC '98)*, pages 106–113. ACM, 1998.

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

[5] David Arthur and Sergei Vassilvitskii. How slow is the K-means method? In *Proceedings of the Twenty-second Annual Symposium on Computational Geometry (SoCG '06)*, pages 144–153. ACM, 2006.

[6] David Arthur and Sergei Vassilvitskii. K-means++: The advantages of careful seeding. In *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA '07)*, pages 1027–1035. SIAM, 2007.

[7] 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 the Thirty-third Annual ACM Symposium on Theory of Computing (STOC '01)*, pages 21–29. ACM, 2001.

[8] Vijay Arya, Naveen Garg, Rohit Khandekar, Adam Meyerson, Kamesh Munagala, and Vinayaka Pandit. Local search heuristics for K-median and facility location problems. *SIAM J. Comput.*, 33(3):544–562, 2004.

[9] Pranjal Awasthi, Avrim Blum, and Or Sheffet. Stability yields a PTAS for K-median and K-means clustering. In *Proceedings of the 2010 IEEE 51st Annual Symposium on Foundations of Computer Science (FOCS '10)*, pages 309–318. IEEE Computer Society, 2010.

[10] Pranjal Awasthi, Moses Charikar, Ravishankar Krishnaswamy, and Ali Kemal Sinop. The Hardness of Approximation of Euclidean K-Means. In *Proceedings of 31st International Symposium on Computational Geometry (SoCG '15)*, Leibniz International Proceedings in Informatics (LIPIcs), pages 754–767, 2015.

[11] Mihai Bādoiu, Sariel Har-Peled, and Piotr Indyk. Approximate clustering via Coresets. In *Proceedings of the Thiry-fourth Annual ACM Symposium on Theory of Computing (STOC '02)*, pages 250–257. ACM, 2002.

[12] Sayan Bandyapadhyay and Kasturi Varadarajan. On variants of K-means clustering. In *Proceedings of the 32nd International Symposium on Computational Geometry (SoCG '16)*, Leibniz International Proceedings in Informatics (LIPIcs), 2016.

[13] Johannes Blömer, Christiane Lammersen, Melanie Schmidt, and Christian Sohler. Theoretical analysis of the K-means algorithm - A survey. *CoRR*, abs/1602.08254, 2016.

[14] Jaroslaw Byrka, Thomas Pensyl, Bartosz Rybicki, Aravind Srinivasan, and Khoa Trinh. An improved approximation for K-median, and positive correlation in budgeted optimization. In *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA '15)*, pages 737–756. SIAM, 2015.

[15] Ke Chen. On Coresets for K-median and K-means clustering in metric and Euclidean spaces and their applications. *SIAM J. Comput.*, 39:923–947, 2009.

[16] Vincent Cohen-Addad, Philip N Klein, and Claire Mathieu. Local search yields approximation schemes for k-means and k-median in euclidean and minor-free metrics. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science (FOCS '16)*. IEEE Computer Society, 2016.

[17] Vincent Cohen-Addad, Philip N Klein, and Claire Mathieu. Local search yields approximation schemes for k-means and k-median in euclidean and minor-free metrics. *arXiv preprint arXiv:1603.09535*, 2016.

[18] Vincent Cohen-Addad, and Claire Mathieu. The Unreasonable Success of Local Search: Geometric Optimization *arXiv preprint arXiv:1410.0553*, 2014.

[19] Sanjoy Dasgupta. How fast is K-means? In *Proceedings of the 16th Annual Conference on Learning Theory and 7th Kernel Workshop (COLT/Kernel '03)*, pages 735–735, 2003.

[20] W. Fernandez de la Vega, Marek Karpinski, Claire Kenyon, and Yuval Rabani. Approximation schemes for clustering problems. In *Proceedings of the Thirty-fifth Annual ACM Symposium on Theory of Computing (STOC '03)*, pages 50–58. ACM, 2003.

[21] P. Drineas, A. Frieze, R. Kannan, S. Vempala, and V. Vinay. Clustering large graphs via the singular value decomposition. *Mach. Learn.*, 56(1-3):9–33, 2004.

[22] Jittat Fakcharoenphol, Satish Rao, and Kunal Talwar. A tight bound on approximating arbitrary metrics by tree metrics. In *Proceedings of the Thirty-fifth Annual ACM Symposium on Theory of Computing (STOC '03)*, pages 448–455. ACM, 2003.

[23] Dan Feldman, Morteza Monemizadeh, and Christian Sohler. A PTAS for K-means clustering based on weak Coresets. In *Proceedings of the Twenty-third Annual Symposium on Computational Geometry (SoCG '07)*, SoCG '07, pages 11–18. ACM, 2007.

[24] Dan Feldman, Melanie Schmidt, and Christian Sohler. Turning big data into tiny data: Constant-size Coresets for K-means, PCA and projective clustering. In *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA '13)*, pages 1434–1453. SIAM, 2013.

[25] Zachary Friggstad, Mohsen Rezapour, and Mohammad R. Salavatipour. Local search yields a PTAS for $k$-means in doubling metrics. *arXiv preprint arXiv:1603.08976*, 2016.

[26] A. Gupta and T. Tangwongsan. Simpler analyses of local search algorithms for facility location. *CoRR, abs/0809.2554*, 2008.

[27] Sariel Har-Peled and Akash Kushal. Smaller Coresets for K-median and K-means clustering. In *Proceedings of the Twenty-first Annual Symposium on Computational Geometry (SoCG '05)*, pages 126–134. ACM, 2005.

[28] Sariel Har-Peled and Soham Mazumdar. On Coresets for K-means and K-median clustering. In *Proceedings of the Thirty-sixth Annual ACM Symposium on Theory of Computing (STOC '04)*, pages 291–300. ACM, 2004.

[29] Sariel Har-Peled and Bardia Sadri. How fast is the K-means method? *Algorithmica*, 41(3):185–202, 2005.

[30] Mary Inaba, Naoki Katoh, and Hiroshi Imai. Applications of weighted Voronoi diagrams and randomization to variance-based K-clustering. In *Proceedings of the tenth annual Symposium on Computational Geometry (SoCG '94)*, pages 332–339. ACM, 1994.

[31] Anil K. Jain. Data clustering: 50 years beyond K-means. *Pattern Recogn. Lett.*, 31(8):651–666, 2010.

[32] Tapas Kanungo, David M. Mount, Nathan S. Netanyahu, Christine D. Piatko, Ruth Silverman, and Angela Y. Wu. A local search approximation algorithm for K-means clustering. *Comput. Geom. Theory Appl.*, 28(2-3):89–112, 2004.

[33] Stavros G. Kolliopoulos and Satish Rao. A nearly linear-time approximation scheme for the Euclidean Kappa-median problem. In *Proceedings of the 7th Annual European Symposium on Algorithms (ESA '99)*, pages 378–389. Springer-Verlag, 1999.

[34] Amit Kumar and Ravindran Kannan. Clustering with spectral norm and the K-means algorithm. In *Proceedings of the 2010 IEEE 51st Annual Symposium on Foundations of Computer Science (FOCS '10)*, pages 299–308. IEEE Computer Society, 2010.

[35] Amit Kumar, Yogish Sabharwal, and Sandeep Sen. A simple linear time $(1+\epsilon)$-approximation algorithm for K-means clustering in any dimensions. In *Proceedings of the 45th Annual IEEE Symposium on Foundations of Computer Science (FOCS '04)*, pages 454–462. IEEE Computer Society, 2004.

[36] Amit Kumar, Yogish Sabharwal, and Sandeep Sen. Linear-time approximation schemes for clustering problems in any dimensions. *J. ACM*, 57(2):5:1–5:32, 2010.

[37] Michael Langberg and Leonard J. Schulman. Universal $\epsilon$-approximators for integrals. In *Proceedings of the Twenty-first Annual ACM-SIAM Symposium on Discrete Algorithms (SODA '10)*, pages 598–607. SIAM, 2010.

[38] Shi Li and Ola Svensson. Approximating K-median via pseudo-approximation. In *Proceedings of the Forty-fifth Annual ACM Symposium on Theory of Computing (STOC '13)*, pages 901–910. ACM, 2013.

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

[40] Meena Mahajan, Prajakta Nimbhorkar, and Kasturi Varadarajan. The planar K-means problem is NP-hard. In *Proceedings of the 3rd International Workshop on Algorithms and Computation (WALCOM '09)*, pages 274–285. Springer-Verlag, 2009.

[41] Jırı Matoušek. On approximate geometric k-clustering. *Discrete & Computational Geometry*, 24(1):61–84, 2000.

[42] Rafail Ostrovsky and Yuval Rabani. Polynomial-Time Approximation Schemes for geometric min-sum median clustering. *J. ACM*, 49(2):139–156, 2002.

[43] 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, 2013.

[44] Kunal Talwar. Bypassing the embedding: Algorithms for low dimensional metrics. In *Proceedings of the Thirty-sixth Annual ACM Symposium on Theory of Computing (STOC '04)*, pages 281–290. ACM, 2004.

[45] Andrea Vattani. The hardness of K-means clustering in the plane. *Manuscript*, 2009.

[46] Andrea Vattani. K-means requires exponentially many iterations even in the plane. *Discrete Comput. Geom.*, 45(4):596–616, 2011.