

Approximate Clustering without the Approximation

Maria-Florina Balcan*

Avrim Blum†

Anupam Gupta‡

Abstract

Approximation algorithms for clustering points in metric spaces is a flourishing area of research, with much research effort spent on getting a better understanding of the approximation guarantees possible for many objective functions such as k -median, k -means, and min-sum clustering.

This quest for better approximation algorithms is further fueled by the implicit hope that these better approximations also yield more accurate clusterings. E.g., for many problems such as clustering proteins by function, or clustering images by subject, there is some unknown correct “target” clustering and the implicit hope is that approximately optimizing these objective functions will in fact produce a clustering that is close pointwise to the truth.

In this paper, we show that if we make this implicit assumption explicit—that is, if we assume that any c -approximation to the given clustering objective Φ is ϵ -close to the target—then we can produce clusterings that are $O(\epsilon)$ -close to the target, *even for values c for which obtaining a c -approximation is NP-hard*. In particular, for k -median and k -means objectives, we show that we can achieve this guarantee for any constant $c > 1$, and for the min-sum objective we can do this for any constant $c > 2$.

Our results also highlight a surprising conceptual difference between assuming that the *optimal* solution to, say, the k -median objective is ϵ -close to the target, and assuming that any *approximately optimal* solution is ϵ -close to the target, even for approximation factor say $c = 1.01$. In the former case, the problem of finding a solution that is $O(\epsilon)$ -close to the target remains computationally hard, and yet for the latter we have an efficient algorithm.

1 Introduction

The field of approximation algorithms for clustering points in metric spaces is a very active one, with a large number of algorithms having been developed for clustering objectives like k -median, k -means, and min-sum clustering. The k -median problem has a $3 + \epsilon$ -approximation [AGK⁺04], and it is NP-hard to approximate to better than $1 + 2/e$ [JMS02]. The k -means problem for general metric spaces has a

constant-factor approximation, and admits a PTAS in Euclidean spaces for constant number of clusters k [KSS04]. The min-sum clustering problem admits an $O(\log^{1+\delta} n)$ -approximation for general metric spaces, and admits a PTAS when k is a constant [dIVKKR03]. For most of these problems, the approximation guarantees do not match the known hardness results, and much effort is spent on obtaining tighter approximation guarantees.

However, this search for better approximation algorithms is motivated not just by the desire to pin down the tractability threshold for these objectives: there is the underlying hope that better approximations will give more meaningful clusterings of the underlying data. Indeed, for many clustering problems, such as clustering proteins by function, or clustering images by subject, the real goal is to classify the points correctly, and these objectives are only a proxy. That is, there is some unknown correct “target” clustering—such as grouping the proteins by their actual functions, or grouping the images by who is actually in them—and the implicit hope is that approximately optimizing these objectives will in fact produce a clustering that is close in symmetric difference to the truth. In other words, implicit in taking the approximation-algorithms approach is the hope that any c -approximation to our given objective will be pointwise close to the true answer, and our motivation for improving a c_2 -approximation to a c_1 -approximation (for $c_1 < c_2$) is that perhaps this closeness property holds for c_1 but not c_2 .

In this paper, we show that if we make this implicit assumption explicit, and assume that any c -approximation to the given objective Φ is ϵ -close pointwise to the target clustering, then we can in fact produce a clustering that is $O(\epsilon)$ -close to the target, *even for values c for which obtaining a c -approximation is provably NP-hard*. In particular, for k -median and k -means objectives, we achieve this guarantee for any constant $c > 1$, and for min-sum we do this for any constant $c > 2$ when the target clusters are “large”. Moreover, if clusters are sufficiently large compared to $\frac{\epsilon n}{c-1}$ then for k -median we can actually get ϵ -close (rather than $O(\epsilon)$ -close) to the target.

Thus, we show that we do not need to find a better approximation algorithm in order to get the properties that such algorithms would imply: we can approximate the target without approximating the objective (up to a constant factor loss in the error ϵ in some cases). Moreover, the problem of finding a c -approximation to these objectives even with this assumption is as hard as finding a c -approximation to them without it (see Theorem A.2 in Appendix A) so we *must* bypass the objective to do so.

*School of Computer Science, Carnegie Mellon University. Supported in part by NSF grant CCF-0514922, by an IBM Graduate Fellowship, and by a Google Research Grant. ninamf@cs.cmu.edu

†School of Computer Science, Carnegie Mellon University. Supported in part by NSF grant CCF-0514922, and by a Google Research Grant. avrim@cs.cmu.edu

‡Computer Science Department, Carnegie Mellon University. Supported in part by NSF awards CCF-0448095 and CCF-0729022, and an Alfred P. Sloan Fellowship.

Our results also show that there is a perhaps unexpected conceptual difference between assuming that the *optimal* solution to, say, the k -median objective is ϵ -close to the target, and assuming that any *approximately optimal* solution is ϵ -close to the target, even for approximation factor $c = 1.01$ (say). In the former case, the problem of finding a solution that is $O(\epsilon)$ -close to the target remains computationally hard (see Section 2.1 and Appendix A), and yet for the latter case we give efficient algorithms.

1.1 Related Work

Work on approximation algorithms: For k -median, $O(1)$ -approximations were given by [CGTS99, JV01, CG99] and the best approximation guarantee known is $(3 + \epsilon)$ due to [AGK⁺04]. A reduction from max- k -coverage shows an easy $(1 + 2/e)$ -hardness of approximation [GK99, JMS02]. The k -median problem on constant-dimensional Euclidean spaces admits a PTAS [ARR99].

For k -means on general metric spaces, one can derive a constant approximation using ideas from k -median—the squared distances do not form a metric, but are close enough for the proofs to go through; an approximation-hardness of $1 + 8/e$ follows from the ideas of [GK99, JMS02]. This problem is very often studied in Euclidean space, where a near-linear time $(1 + \epsilon)$ -approximation algorithm is known for the case of constant k and ϵ [KSS04]. Lloyd’s local search algorithm [Llo82] is often used in practice, despite having poor worst-case performance [AV06]. Ostrovsky et al. [ORSS06] study ways of seeding Lloyd’s local search algorithm: they show that on instances satisfying an ϵ -separation property, this seeding results in solutions with provable approximation guarantees. We show in Section 2.2 that their assumption can be quite a bit stronger than ours—though their goal is different (they want to approximate the objective whereas we want to approximate the target).

Min-sum k -clustering on general metric spaces admits a PTAS for the case of constant k by Fernandez de la Vega et al. [dIVKKR03] (see also [Ind99]). For the case of arbitrary k there is an $O(\delta^{-1} \log^{1+\delta} n)$ -approximation algorithm in time $n^{O(1/\delta)}$ due to Bartal et al. [BCR01]. The problem has also been studied in geometric spaces for constant k by Schulman [Sch00] who gave an algorithm for (R^d, ℓ_2^2) that either output a $(1 + \epsilon)$ -approximation, or a solution that agreed with the *optimum* clustering on $(1 - \epsilon)$ -fraction of the points (but could have much larger cost than optimum); the runtime is $O(n^{\log \log n})$ in the worst case and linear for sublogarithmic dimension d .

Related work on error to a target clustering: There has been significant work in machine learning and theoretical computer science on clustering or learning with mixture models [AM05, AK05, DHS01, DGL96, KSV05, VW04, Das99]. That work, like ours, has an explicit notion of a correct ground-truth clustering of the data points; however, it makes strong probabilistic assumptions about the data.

In recent work, Balcan et al. [BBV08] investigated the goal of approximating a desired target clustering without probabilistic assumptions. They analyzed what properties of

a pairwise similarity function are sufficient to produce a tree such that some unknown pruning is close to the target, or a small list of clusterings such that the target is close to one of them. In relation to implicit assumptions about approximation algorithms, [BBV08] made the observation that for k -median, the assumption that any 2-approximation is ϵ -close to the target implies that most of the data satisfies a certain separation property, which they then use to construct a hierarchical clustering such that the target clustering is close to some pruning of the hierarchy. Inspired by their approach, in this paper we initiate a systematic investigation of the consequences of such assumptions about approximation algorithms. Moreover, the goals in this paper will be stronger—we want to output a *single* approximately correct clustering (as opposed to a list of clusterings or a hierarchy), and we want to succeed for any $c > 1$.

2 Definitions and Preliminaries

The clustering problems in this paper fall into the following general framework: we are given a metric space $\mathcal{M} = (X, d)$ with point set X and a distance function $d : \binom{X}{2} \rightarrow R_{\geq 0}$ satisfying the triangle inequality—this is the ambient space. We are also given the actual point set $S \subseteq X$ we want to cluster; we use n to denote the cardinality of S . A k -clustering \mathcal{C} is a partition of S into k sets C_1, C_2, \dots, C_k . In this paper, we always assume that there is a *true* or *target* k -clustering \mathcal{C}_T for the point set S .

Commonly used clustering algorithms seek to minimize some objective function or “score”; e.g., the k -median clustering objective assigns to each cluster C_i a “median” $c_i \in X$ and seeks to minimize $\Phi_1(\mathcal{C}) = \sum_{i=1}^k \sum_{x \in C_i} d(x, c_i)$, k -means clustering minimizes $\Phi_2(\mathcal{C}) = \sum_{i=1}^k \sum_{x \in C_i} d(x, c_i)^2$, and *min-sum clustering* minimizes $\Phi_\Sigma = \sum_{i=1}^k \sum_{x, y \in C_i} d(x, y)$. Given a function Φ and instance (\mathcal{M}, S) , let $\text{OPT}_\Phi = \min_{\mathcal{C}} \Phi(\mathcal{C})$, where the minimum is over all k -clusterings of (\mathcal{M}, S) .

We define the distance $\text{dist}(\mathcal{C}, \mathcal{C}')$ between two k -clusterings $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ and $\mathcal{C}' = \{C'_1, C'_2, \dots, C'_k\}$ as the fraction of points on which they disagree under the optimal matching of clusters in \mathcal{C} to clusters in \mathcal{C}' ; i.e., $\text{dist}(\mathcal{C}, \mathcal{C}') = \min_{\sigma \in S_k} \frac{1}{n} \sum_{i=1}^k |C_i - C'_{\sigma(i)}|$, where S_k is the set of bijections $\sigma : [k] \rightarrow [k]$. We say that two clusterings \mathcal{C} and \mathcal{C}' are ϵ -close if $\text{dist}(\mathcal{C}, \mathcal{C}') < \epsilon$. Note that if \mathcal{C} and \mathcal{C}' are ϵ -close and all clusters C_i have size at least $2\epsilon n$, then the bijection σ minimizing $\frac{1}{n} \sum_{i=1}^k |C_i - C'_{\sigma(i)}|$ is unique; in this case we call this the *optimal bijection* σ and we say that \mathcal{C} and \mathcal{C}' agree on x if $x \in C_i \cap C'_{\sigma(i)}$ for some i .

The following definition is central to our discussion:

DEFINITION 1. (THE (c, ϵ) -PROPERTY) *Given an objective function Φ (such as k -median, k -means, or min-sum), we say that instance (\mathcal{M}, S) satisfies the (c, ϵ) -property for Φ if all clusterings \mathcal{C} with $\Phi(\mathcal{C}) \leq c \cdot \text{OPT}_\Phi$ are ϵ -close to the target clustering \mathcal{C}_T for (\mathcal{M}, S) .*

The above assumption is often implicitly made when propos-

ing to use a c -approximation for objective Φ to solve a clustering problem in which the true goal is to classify data points correctly; similarly, the motivation for improving a c_2 approximation to a $c_1 < c_2$ approximation is that perhaps the data satisfies the (c_1, ϵ) property for Φ but not the (c_2, ϵ) property.

Note that for any $c > 1$, the (c, ϵ) -property does not require that the target clustering \mathcal{C}_T exactly coincide with the optimal clustering \mathcal{C}^* under objective Φ . However, it does imply the following simple facts:

FACT 2.1. *If (\mathcal{M}, S) satisfies the (c, ϵ) -property for Φ , then:*

- (a) *The target clustering \mathcal{C}_T , and the optimal clustering \mathcal{C}^* are ϵ -close.*
- (b) *The distance between k -clusterings is a metric, and hence a (c, ϵ) property with respect to the target clustering \mathcal{C}_T implies a $(c, 2\epsilon)$ property with respect to the optimal clustering \mathcal{C}^* .*

Thus, we can act as if the optimal clustering is indeed the target up to a constant factor loss in the error rate.

2.1 Two Strawman Solutions, and Why They Fail

Before proceeding to our results, we first consider two “strawman” approaches to achieving our goals, and indicate why they do not work.

- First, suppose that the (c, ϵ) -property for some objective Φ implied, say, the $(2c, 2\epsilon)$ property. Then it would be sufficient to simply apply an $O(c)$ approximation in order to have error $O(\epsilon)$ with respect to the target. However, for any $c_1 < c_2$ and any $\epsilon, \alpha > 0$, for each of the three objectives we consider (k -median, k -means, and min-sum), there exists a family of metric spaces and target clusterings satisfying the (c_1, ϵ) property for that objective, and yet that do not satisfy even the $(c_2, 1/2 - \alpha)$ property (See Appendix, Theorem A.1). Thus, the result of a direct application of a c_2 -approximation is nearly as poor as possible.
- Second, perhaps the (c, ϵ) assumption implies that finding a c -approximation is somehow trivial. However, this is not the case either: for any $c > 1$, the problem of finding a c -approximation to any of the three objectives we consider under the (c, ϵ) assumption is as hard as finding a c -approximation in general (Theorem A.2).

It is also interesting to note that results of the form we are aiming for are *not possible given only the $(1, \epsilon)$ property*. Indeed, because the standard hardness-of-approximation reduction for k -median produces a metric in which all pairwise distances lie in a bounded range, the reduction also implies that it is NP-hard, given a data set satisfying the $(1, \epsilon)$ property, to find a clustering of error $O(\epsilon)$; see Theorem A.3.

2.2 Relationship to Similar Concepts

Ostrovsky et al. [ORSS06] study k -means in Euclidean space; they call a k -means instance ϵ -separated if the optimal k -means cost is at most ϵ^2 times the cost of optimally opening $k - 1$ means; under this assumption on the

input, they show how to seed Lloyd’s method to obtain a $1 + O(\epsilon^2)$ approximation in d -dimensional Euclidean space in time $O(nkd + k^3d)$, and a $(1 + \delta)$ -PTAS with run-time $nd2^{k(1+\epsilon^2)/\delta}$. In Theorem 5.1 of their paper, they show that their ϵ -separatedness assumption implies that any near-optimal solution to k -means is $O(\epsilon^2)$ -close to the optimal clustering. *However, the converse is not true:* an instance could satisfy our property without being ϵ -separated.¹ For example, consider $k = 2$ where target cluster C_1 has $(1 - \alpha)n$ points and target cluster C_2 has αn points. Any two points inside the same cluster have distance 1 and any two points inside different clusters have distance $1 + 1/\epsilon$. For any $\alpha \in (\epsilon, 1 - \epsilon)$, this satisfies the $(2, \epsilon)$ property for k -median (and the $(2, \epsilon^2)$ property for k -means for any $\alpha \in (\epsilon^2, 1 - \epsilon^2)$). However, it need not satisfy the ϵ -separation property: for $\alpha = 2\epsilon$, the optimal 2-median solution has cost $n - 2$, but the optimal 1-median has cost $< 3n$. Likewise for $\alpha = 2\epsilon^2$, the optimal 2-means solution has cost $n - 2$, but the optimal 1-means has cost $< (3 + 4\epsilon)n$. Thus, the ratio of costs for $k = 1$ and $k = 2$ is not so large.

3 The k -Median Problem

We first study k -median clustering under the (c, ϵ) -property. Our main results are that for any constant $c > 1$, (1) if all clusters are “large”, then this property allows us to efficiently find a clustering that is ϵ -close to the target clustering, and (2) for *any* cluster sizes, we can efficiently find a clustering that is $O(\epsilon)$ -close to the target. To prove these results, we first investigate the implications of the (c, ϵ) -property in Section 3.1. We then give our algorithm for the case that all clusters are large in Section 3.2, and our algorithm for arbitrary cluster sizes in Section 3.3.

3.1 Implications of the (c, ϵ) -Property

Given an instance of k -median specified by a metric space $\mathcal{M} = (X, d)$ and a set of points $S \subseteq X$, fix an optimal k -median clustering $\mathcal{C}^* = \{C_1^*, \dots, C_k^*\}$, and let c_i^* be the center point for C_i^* . Let $w(x) = \min_i d(x, c_i^*)$ be the contribution of x to the k -median objective in \mathcal{C}^* (i.e., x ’s “weight”), and let $w_2(x)$ be x ’s distance to the second-closest center point among $\{c_1^*, c_2^*, \dots, c_k^*\}$. Also, let $w = \frac{1}{n} \sum_{i=1}^n w(x) = \frac{\text{OPT}}{n}$ be the average weight of the points. Finally, let $\epsilon^* = \text{dist}(\mathcal{C}_T, \mathcal{C}^*)$; so, by our assumption we have $\epsilon^* < \epsilon$.

LEMMA 3.1. *If the k -median instance (\mathcal{M}, S) satisfies the $(1 + \alpha, \epsilon)$ -property with respect to \mathcal{C}_T , and each cluster in \mathcal{C}_T has size at least $2\epsilon n$, then*

- (a) *less than $(\epsilon - \epsilon^*)n$ points $x \in S$ on which \mathcal{C}_T and \mathcal{C}^* agree have $w_2(x) - w(x) < \frac{\alpha w}{\epsilon}$, and*
 - (b) *at most $5\epsilon n/\alpha$ points $x \in S$ have $w(x) \geq \frac{\alpha w}{5\epsilon}$.*
- For the case of general cluster sizes in \mathcal{C}_T we replace (a) and (b) with:*
- (a’) *less than $6\epsilon n$ points $x \in S$ have $w_2(x) - w(x) < \frac{\alpha w}{2\epsilon}$.*
 - (b’) *at most $10\epsilon n/\alpha$ points $x \in S$ have $w(x) \geq \frac{\alpha w}{10\epsilon}$.*

¹[ORSS06] shows an implication in this direction (Theorem 5.2); however, the notion of closeness used there is much stronger.

Proof: To prove Property (a), assume to the contrary. Then one could take C^* and move $(\epsilon - \epsilon^*)n$ points x on which \mathcal{C}_T and C^* agree to their second-closest clusters, increasing the objective by at most αOPT . Moreover, this new clustering $C' = \{C'_1, \dots, C'_k\}$ has distance at least ϵ from \mathcal{C}_T , because we begin at distance ϵ^* from \mathcal{C}_T and each move increases this distance by $\frac{1}{n}$ (here we use the fact that because each cluster in \mathcal{C}_T has size at least $2\epsilon n$, the optimal bijection between \mathcal{C}_T and C' remains the same as the optimal bijection between \mathcal{C}_T and C^*). Hence we have a clustering that is not ϵ -close to \mathcal{C}_T with cost only $(1 + \alpha)\text{OPT}$, a contradiction. Property (b) follows from the definition of the average weight w , and Markov's inequality. For Property (a'), we use Lemma A.1 in the Appendix which addresses the case of small clusters. Specifically, assuming for contradiction that $6\epsilon n$ points satisfy (a'), Lemma A.1 states that we can find a subset of $2\epsilon n$ of them such that starting from C^* , for each one that we move to its second-closest cluster, the distance from C^* increases by $\frac{1}{n}$. Therefore, by increasing the objective by at most αOPT we can create a clustering C' that is distance at least 2ϵ from C^* , and so is not ϵ -close to \mathcal{C}_T . Property (b') again follows from Markov's inequality. ■

For the case that each cluster in \mathcal{C}_T has size at least $2\epsilon n$, define the *critical distance* $d_{crit} = \frac{\alpha w}{5\epsilon}$, else define $d_{crit} = \frac{\alpha w}{10\epsilon}$; i.e., these are the values in properties (b) and (b') respectively of Lemma 3.1. We call point x *good* if both $w(x) < d_{crit}$ and $w_2(x) - w(x) \geq 5d_{crit}$, else x is called *bad*; by Lemma 3.1 and the definition of ϵ^* , if all clusters in the target have size greater than $2\epsilon n$ then at most a $(1 + 5/\alpha)\epsilon$ fraction of points are bad, and in general at most a $(6 + 10/\alpha)\epsilon$ fraction of points are bad. Let X_i be the *good* points in the optimal cluster C_i^* , and let $B = S \setminus \cup X_i$ be the bad points. Let $b = |B|$.

LEMMA 3.2. (THRESHOLD GRAPH) Define the τ -threshold graph $G_\tau = (S, E_\tau)$ by connecting all pairs $\{x, y\} \in \binom{S}{2}$ with $d(x, y) \leq \tau$. For an instance satisfying the $(1 + \alpha, \epsilon)$ -property and $\tau = 2d_{crit}$, the threshold graph G_τ has the following properties:

- (i) For all x, y in the same X_i , the edge $\{x, y\} \in E(G_\tau)$.
 - (ii) For $x \in X_i$ and $y \in X_{j \neq i}$, $\{x, y\} \notin E(G_\tau)$.
- Moreover, such points x, y do not share any neighbors in G_τ .

Proof: For part (i), since x, y are both good, they are at distance less than d_{crit} to their cluster center, by Lemma 3.1 (b or b'). By the triangle inequality, the distance $d(x, y) \leq d(x, c_i^*) + d(c_i^*, y) \leq 2 \times d_{crit} = \tau$. For part (ii), the distance from x to y 's cluster center c_j^* is at least $5d_{crit}$, by Lemma 3.1 (a or a'). Again by the triangle inequality, $d(x, y) \geq d(x, c_j^*) - d(y, c_j^*) > 5d_{crit} - d_{crit} = 2\tau$. Since each edge in G_τ is between points at distance at most τ , the points x, y cannot share any common neighbors. ■

Hence, the graph G_τ for the above value of τ is fairly simple to describe: each X_i forms a clique, and the neighborhood $N_{G_\tau}(X_i)$ of X_i lies entirely in the bad bucket B with no

edges going between X_i and $X_{j \neq i}$, or between X_i and $N_{G_\tau}(X_{j \neq i})$. We now show how we can use this to find a clustering of error at most ϵ if the size of each X_i is large (Section 3.2) and how we can get error $O(\epsilon)$ for general cluster sizes (Section 3.3).

3.2 An algorithm for large clusters We begin with the following lemma. Recall that $b = |B|$.

LEMMA 3.3. Given a graph $G = (S, E)$ satisfying properties (i), (ii) of Lemma 3.2 and where each $|X_i| \geq b + 2$, there is an efficient algorithm that outputs a k -clustering with each X_i contained in a distinct cluster.

Proof: Construct a graph $H = (S, E')$ where we place an edge $\{x, y\} \in E'$ if x and y have at least b common neighbors in G . By property (i) each X_i is a clique of size $\geq b + 2$ in G , so each pair $x, y \in X_i$ has at least b common neighbors in G and hence $\{x, y\} \in E'$. Now consider $x \in X_i \cup N_G(X_i)$, and $y \notin X_i \cup N_G(X_i)$: we claim $\{x, y\} \notin E'$. Indeed, by property (ii), x and y cannot share neighbors that lie in X_i (since $y \notin X_i \cup N_G(X_i)$), nor in some $X_{j \neq i}$ (since $x \notin X_j \cup N_G(X_j)$). Hence the common neighbors of x, y all lie in B , which has size b . Moreover, at least one of x and y must itself belong to B , else they would have no common neighbors by property (ii); hence, the number of distinct common neighbors is at most $b - 1$, which implies that $\{x, y\} \notin E'$.

Thus each X_i is contained within a distinct component of the graph H ; the remaining components of H contain vertices from the "bad bucket" B . Since the X_i 's are larger than B , we can obtain the claimed clustering by taking the largest k components in H , adding the vertices of all other smaller components to any of these, and using this as the k -clustering. ■

We now show how we can use Lemma 3.3 to find a clustering that is ϵ -close to \mathcal{C}_T . For simplicity, we begin by assuming that we are given the value of $w = \frac{\text{OPT}}{n}$, and then we show how this assumption can be removed.

THEOREM 3.1. (THE "KNOWN w " CASE) If the k -median instance satisfies the $(1 + \alpha, \epsilon)$ -property and each cluster in \mathcal{C}_T has size at least $(3 + 10/\alpha)\epsilon n + 2$, then given w we can efficiently find a clustering that is ϵ -close to \mathcal{C}_T .

Proof: Since each cluster in the target clustering has at least $(3 + 10/\alpha)\epsilon n + 2$ points, and the optimal k -median clustering C^* differs from the target clustering by $\epsilon^* n \leq \epsilon n$ points, each cluster in C^* must have at least $(2 + 10/\alpha)\epsilon n + 2$ points. Moreover, by Lemma 3.1, the bad points B constitute at most $(1 + 5/\alpha)\epsilon n$ points, and hence each $|X_i| = |C_i^* \setminus B| \geq (1 + 5/\alpha)\epsilon n + 2 = b + 2$.

Now, given w , we can construct the graph G_τ with $\tau = 2d_{crit}$ (which we can compute from the given value of w), and apply Lemma 3.3 to find a k -clustering C' where each X_i is contained within a distinct cluster. Note that this clustering C' differs from the optimal clustering C^* only in the bad points which constitute an $O(\epsilon/\alpha)$ fraction of the

total. Hence, it is at distance $O(\epsilon/\alpha + \epsilon)$ from the target. However, our goal is to get ϵ -close to the target, which we do as follows.

Call a point x “red” if it satisfies condition (a) in Lemma 3.1 (i.e., $w_2(x) - w(x) < 5d_{crit}$), “yellow” if it is not red but satisfies condition (b) in Lemma 3.1 (i.e., $w(x) \geq d_{crit}$), and “green” otherwise. So, the green points are those in the sets X_i , and we have partitioned the bad set B into red points and yellow points. Let $C' = \{C'_1, \dots, C'_k\}$ and recall that C' agrees with C^* on the green points, so without loss of generality we may assume $X_i \subseteq C'_i$. We now construct a new clustering C'' that agrees with C^* on both the green and yellow points. Specifically, for each point x and each cluster C'_j , compute the median distance $d_{med}(x, j)$ between x and all points in C'_j ; then insert x into the cluster C''_i for $i = \operatorname{argmin}_j d_{med}(x, j)$. Since each non-red point x satisfies $w_2(x) - w(x) \geq 5d_{crit}$, and all green points g satisfy $w(g) < d_{crit}$, this means that any non-red point x must satisfy the following two conditions: (1) for a green point g_1 in the same cluster as x in C^* we have $d(x, g_1) \leq w(x) + d_{crit}$, and (2) for a green point g_2 in a different cluster than x in C^* we have $d(x, g_2) \geq w_2(x) - d_{crit} \geq w(x) + 4d_{crit}$. Therefore, $d(x, g_1) < d(x, g_2)$. Since each cluster in C' has a strict majority of green points (even with point x removed) all of which are clustered as in C^* , this means that for a non-red point x , the median distance to points in its correct cluster with respect to C^* is less than the median distance to points in any incorrect cluster. Thus, C'' agrees with C^* on all non-red points. Finally, since there are at most $(\epsilon - \epsilon^*)n$ red points on which C_T and C^* agree by Lemma 3.1—and C'' and C_T might disagree on all these points—this implies $\operatorname{dist}(C'', C_T) \leq (\epsilon - \epsilon^*) + \epsilon^* = \epsilon$ as desired. ■

We now extend the above argument to the case where we are not given the value of w .

THEOREM 3.2. (THE “UNKNOWN w ” CASE) *If the k -median instance satisfies the $(1 + \alpha, \epsilon)$ -property and each cluster in C_T has size at least $(4 + 15/\alpha)\epsilon n + 2$, then we can efficiently find a clustering that is ϵ -close to C_T .*

Proof: If we are not given the value w , we instead run the algorithm of Lemma 3.3 repeatedly for different values of w , starting with $w = 0$ (so the graph G_τ is empty) and at each step increasing w to the next value such that G_τ contains at least one new edge (so we have at most n^2 different guesses to try). If some guess for w causes the k largest components of H to miss more than $b = (2 + 10/\alpha)\epsilon n$ points, or if any of these components have size $\leq b$, then we reject, and increase w . Otherwise, we define C' to be the k largest components in H (so up to b points may be unclustered) and continue to the second phase of the algorithm for the known- w case constructing clustering C'' .

Note that we still might have too small a guess for w , but this just means that the resulting graphs G_τ and H can only have fewer edges than the corresponding graphs for the correct w . Hence, some of the X_i 's might not have fully formed into connected components in H . However,

if the k largest components have size greater than b , then we never misclassify the good points lying in these largest components. We might misclassify all the bad points (at most b of these), and might fail to cluster at most b of the points in the actual X_i 's (i.e., those not lying in the largest k components), but this nonetheless guarantees that each cluster C'_i contains at least $|X_i| - b \geq b + 2$ correctly clustered green points (with respect to C^*) and at most b misclassified points. Therefore, as shown in the proof of Theorem 3.1, the resulting clustering C'' will correctly cluster all non-red points as in C^* and so is at distance at most $(\epsilon - \epsilon^*) + \epsilon^* = \epsilon$ from C_T . ■

3.3 An Algorithm for the General Case The algorithm in the previous section required the minimum cluster size in the target to be large (of size $\Omega(\epsilon n)$). In this section, we show how this requirement can be removed using a different algorithm that finds a clustering that is $O(\epsilon/\alpha)$ -close to the target; while the algorithm is just as simple, we need to be a bit more careful in the analysis. (Again, we will assume we know $w = \frac{\text{OPT}}{n}$, and discharge this assumption later.)

Algorithm 1 k -median Algorithm: General Case

Input: $w, \epsilon \leq 1, \alpha > 0, k$.

Step 1: Construct the τ -threshold graph G_τ with $\tau = 2d_{crit} = \frac{1}{5} \frac{\alpha w}{\epsilon}$.

Step 2: For $j = 1$ to k do:

Pick the vertex v_j of highest degree in G_τ .

Remove v_j and its neighborhood from G_τ and call this cluster $C(v_j)$.

Step 3: Output the k clusters $C(v_1), \dots, C(v_{k-1}), S - \cup_{i=1}^{k-1} C(v_i)$.

THEOREM 3.3. (k -MEDIAN: GENERAL CASE) *If the k -median instance satisfies the $(1 + \alpha, \epsilon)$ -property and we are given the value of w , the above algorithm produces a clustering which is $O(\epsilon/\alpha)$ -close to the target.*

Proof: Recall the notation from Section 3.1: the graph G_τ satisfies properties (i),(ii) of Lemma 3.2. We show that the greedy method of Step 2 above correctly captures most of the cliques X_1, X_2, \dots, X_k in G_τ —in particular, we show there is a bijection $\sigma : [k] \rightarrow [k]$ such that $\sum_i |X_{\sigma(i)} \setminus C(v_i)| = O(b)$. Since the b bad points (i.e., those in $B = S \setminus \cup_{i=1}^k X_i$) may potentially all be misclassified, this gives an additional error of b .

Let us think of each clique X_i as initially “unmarked”, and then “marking” it the first time we choose a cluster $C(v_j)$ that intersects it. We now consider two cases. If the j^{th} cluster $C(v_j)$ intersects some unmarked clique X_i , we will assign $\sigma(j) = i$. (Note that it is not possible for $C(v_j)$ to intersect two cliques X_i and $X_{j \neq i}$, since by Lemma 3.2(ii) these cliques have no common neighbors.) If $C(v_j)$ misses r_i points from X_i , then since the vertex v_j defining this cluster had maximum degree and X_i is a clique, we must have picked at least r_i elements from B in $C(v_j)$. Therefore the total sum of these r_i can be at most $b = |B|$, and hence

$\sum_j |X_{\sigma(j)} \setminus C(v_j)| \leq b$, where the sum is over j 's that correspond to the first case.

The other case is if $C(j)$ intersects a previously marked clique X_i . In this case we assign $\sigma(j)$ to any arbitrary clique $X_{i'}$ that is not marked by the end of the process. Note that the total number of points in such $C(j)$'s must be at most the number of points remaining in the marked cliques (i.e., $\sum_j r_j$), and possibly the bad points (at most b of them). Since the cliques $X_{i'}$ were unmarked at the end, their sizes must be bounded by the size of the $C(j)$'s, and hence by $|B| + \sum_i r_i \leq 2b$. This shows that the sum over such j 's, $\sum_j |X_{\sigma(j)} \setminus C(v_j)| \leq 2b$. Therefore, overall, the total error over all $C(v_j)$ with respect to the k -median optimal is the two sums above, plus potentially the bad points, which gives us at most $4b$ points. Adding in the extra ϵn to account for the distance between the k -median optimum and the target clustering yields the claimed $4b + \epsilon n = O(\epsilon/\alpha)n$ result for the case that we are given the value of w . ■

Not Knowing the Value of w . If we do not know the value of w (and hence of τ), unfortunately the method used in the proof of Theorem 3.2 may not work, because we might split some large cluster causing substantial error, and not be able to recognize our mistake (because we only miss small clusters which do not result in very many points being left over). However, we can instead run an off-the-shelf k -median approximation algorithm to produce an estimate for w that is off by only a constant factor, and use this estimate instead. In particular, if we have a β -approximation \tilde{w} (i.e., say $w \leq \tilde{w} \leq \beta\tilde{w}$, an analog of Lemma 3.2 holds for the threshold graph $G_{\tau'}$ with the altered threshold $\tau' = \frac{1}{5} \frac{\alpha\tilde{w}}{\epsilon\beta}$, with the number of bad points now bounded by $b' = (6 + 10\beta/\alpha)\epsilon$. The rest of the proof follows unchanged with all b s replaced by b' s, to give us a final bound of $O(\beta\epsilon/\alpha)$ on the number of misclassified points.

4 The k -Means Problem

The algorithm in Section 3.3 for the k -median problem can be easily altered to work for the k -means problem as well. Indeed, if we can prove the existence of a structure like that promised by Lemma 3.1 and Lemma 3.2 (albeit with different parameters), the same algorithm and proof would give a good clustering for any objective function. Given some optimal solution for k -means define $w(x) = \min_i d(x, c_i)$ to be the distance of x to its center, which is the square root of x 's contribution to the k -means objective function; hence $\text{OPT} = \sum_x w(x)^2$. Again, let $w_2(x) = \min_{j \neq i} d(x, c_j)$ be the distance to the second-closest center, and let $\epsilon^* = \text{dist}(C_T, C^*)$.

LEMMA 4.1. *If the k -means instance (\mathcal{M}, S) satisfies the $(1 + \alpha, \epsilon)$ -property and each cluster in the target has size at least $2\epsilon n$, then*

- (a) *less than $(\epsilon - \epsilon^*)n$ points $x \in S$ on which C_T and C^* agree have $w_2(x) < (\frac{\alpha\text{OPT}}{\epsilon n})^{1/2}$, and*
- (b) *at most $25\epsilon n/\alpha$ points $x \in S$ have $w(x) > \frac{1}{5}(\frac{\alpha\text{OPT}}{\epsilon n})^{1/2}$.*

For the case of general cluster sizes we replace (a) and (b)

with:

- (a') *less than $6\epsilon n$ points $x \in S$ have $w_2(x) < (\frac{\alpha\text{OPT}}{2\epsilon n})^{1/2}$.*
- (b') *at most $50\epsilon n/\alpha$ points $x \in S$ have $w(x) > \frac{1}{5}(\frac{\alpha\text{OPT}}{2\epsilon n})^{1/2}$.*

The proof is similar to the proof for Lemma 3.1, and is omitted here. Note that the threshold for $w_2(x)$ in part (a) above is again 5 times the threshold for $w(x)$ in part (b), and similarly for (a') and (b'). We can thus define the critical distance d_{crit} as the value in (b) or (b') respectively, and define the $b = (1 + 25/\alpha)\epsilon n$ points that satisfy either (a) or (b) above (in the large-cluster case) or the $b = (6 + 50/\alpha)\epsilon n$ points satisfying (a') or (b') (in the general case) as *bad*. The rest of the proof for achieving an $O(\epsilon/\alpha)$ -close clustering for k -median now goes through unchanged in the k -means case as well. Note that k -means also has a constant-factor approximation, so the results for the case of unknown w go through similarly, with different constants. Unfortunately, the argument for exact ϵ -closeness breaks down because property (a) in Lemma 4.1 is weaker than property (a) in Lemma 3.1. We therefore have the following theorem.

THEOREM 4.1. *If the instance satisfies the $(1 + \alpha, \epsilon)$ -property for the k -means objective, we can efficiently produce a clustering which is $O(\epsilon/\alpha)$ -close to the target.*

5 The Min-sum Clustering Problem

Recall that the min-sum k -clustering problem asks to find a k -clustering $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ to minimize the objective function

$$\Phi(\mathcal{C}) = \sum_{i=1}^k \sum_{x, y \in C_i} d(x, y).$$

In this section, we show that if our data satisfies the $(2 + \alpha, \epsilon)$ -property for the min-sum objective, and if all the clusters in the target are “large”, then we can find a clustering that is $O(\epsilon)$ -close to the target \mathcal{C}_T . The general idea is reduce to a problem known as “balanced k -median” (which is within a factor of 2 of the min-sum objective) and extend the techniques from the previous sections to this problem.

5.1 Properties of Min-Sum Clustering

The *balanced k -median* clustering objective assigns to each cluster C_i a “median” $c_i \in X$ and seeks to minimize $\Psi(\mathcal{C}) = \sum_{i=1}^k |C_i| \sum_{x \in C_i} d(x, c_i)$. We begin with a useful lemma, which shows that the two objective functions Φ (for min-sum clustering) and Ψ (for balanced- k -median) are related to within a factor of 2.

LEMMA 5.1. ([BCR01]) *Let Ψ be the balanced k -median objective and let Φ be the min-sum objective. For any k -clustering \mathcal{C} of S we have: $\Psi(\mathcal{C})/2 \leq \Phi(\mathcal{C}) \leq \Psi(\mathcal{C})$.*

LEMMA 5.2. *If the instance (\mathcal{M}, S) satisfies the $(2(1 + \alpha), \epsilon)$ -property for the min-sum objective, then (\mathcal{M}, S) satisfies the $(1 + \alpha, \epsilon)$ -property for balanced k -median.*

Henceforth, we will work with the balanced k -median objective function. Let the balanced k -median optimal

clustering be $\mathcal{C}^* = \{C_1^*, \dots, C_k^*\}$ with objective function value $\text{OPT} = \Psi(\mathcal{C}^*)$. For each cluster C_i^* , let c_i^* be the median point in the cluster. For $x \in C_i^*$, define $w(x) = |C_i^*|d(x, c_i^*)$ and let $w = \text{avg}_x w(x) = \frac{\text{OPT}}{n}$. Define $w_2(x) = \min_{j \neq i} d(x, c_j^*)|C_j^*|$. Let $d_{C_i^*} = \sum_{x \in C_i^*} d(x, c_i^*)$, and hence $\text{OPT} = \sum_i |C_i^*|d_{C_i^*}$.

LEMMA 5.3. *If the balanced k -median instance (\mathcal{M}, S) satisfies the $(1 + \alpha, \epsilon)$ -property with respect to the target clustering, then as long as the minimum cluster size is at least $\max(6, 6/\alpha) \cdot \epsilon n$ we have:*

- (a) *at most 2ϵ -fraction of points $x \in S$ have $w_2(x) < \frac{\alpha w}{4\epsilon}$,*
- (b) *at most $60\epsilon/\alpha$ -fraction of $x \in S$ have $w(x) > \frac{\alpha w}{60\epsilon}$.*

Proof: To prove Property (a), assume to the contrary. Then one could move a 2ϵ fraction of points from their clusters in the optimal clustering \mathcal{C}^* to the clusters that define their w_2 value. This may increase the sizes of the clusters; let the new clustering be $\mathcal{C}' = (C'_1, \dots, C'_k)$, where $|C'_i \setminus C_i^*| = \delta_i n$, so that $\sum_i \delta_i = 2\epsilon$. If a point x moves to cluster C'_i from some other cluster, then it now contributes $w_2(x) \cdot \frac{|C'_i|}{|C_i^*|}$. Summing over all the points, we get that the cost $\Psi(\mathcal{C}')$ is at most

$$\Psi(\mathcal{C}') \leq \sum_{i=1}^k \left((|C_i^*| + \delta_i n) d_{C_i^*} + \delta_i n \cdot \frac{\alpha w}{4\epsilon} \cdot \frac{|C_i^*| + \delta_i n}{|C_i^*|} \right)$$

However, $\delta_i n \leq \sum_i \delta_i n \leq 2\epsilon n \leq \frac{\min(1, \alpha)}{3} |C_i^*|$ (since each cluster size is at least $\max(6, 6/\alpha) \cdot \epsilon n$). Hence, we have

$$\begin{aligned} \Psi(\mathcal{C}') &\leq \sum_{i=1}^k \left(1 + \frac{\alpha}{3} \right) |C_i^*| d_{C_i^*} + \frac{4}{3} \sum_{i=1}^k \frac{\delta_i \alpha \text{OPT}}{4\epsilon} \\ &\leq (1 + \alpha) \text{OPT}. \end{aligned}$$

This would give a clustering with cost at most $(1 + \alpha)\text{OPT}$ that is not 2ϵ -close to the optimal clustering \mathcal{C}^* , which is impossible by Fact 2.1(b). Property (b) above follows from Markov's inequality. ■

We call point x *good* if it both $w(x) \leq \frac{\alpha w}{60\epsilon}$ and $w_2(x) \geq \frac{\alpha w}{4\epsilon}$, else x is called *bad*; let X_i be the good points in the optimal cluster C_i^* , and let $B = S \setminus \cup X_i$ be the bad points.

LEMMA 5.4. *If the balanced- k -median instance (\mathcal{M}, S) satisfies the $(1 + \alpha, \epsilon)$ -property, then as long as the minimum cluster size is at least $\max(6, 6/\alpha) \cdot \epsilon n$ we have:*

- (i) *For all x, y in the same X_i , we have $d(x, y) < \frac{\alpha}{30} \frac{1}{\epsilon} \frac{w}{|C_i^*|}$*
- (ii) *For $x \in X_i$ and $y \in X_{j \neq i}$, $d(x, y) > \frac{\alpha}{5} \frac{1}{\epsilon} \frac{w}{\min(|C_i^*|, |C_j^*|)}$*
- (iii) *The number of bad points $|B| = |S \setminus \cup X_i|$ is at most $b := (2 + 60/\alpha)\epsilon n$.*

Proof: For part (i), since $x, y \in X_i \subseteq C_i^*$ are both good, they are at distance less than $\frac{\alpha}{60} \frac{1}{\epsilon} \frac{w}{|C_i^*|}$ to their cluster center (Lemma 5.3(a)), and hence at distance at most $\frac{\alpha}{30} \frac{1}{\epsilon} \frac{w}{|C_i^*|}$ to each other. For part (ii) assume without loss of generality that $|C_i^*| \geq |C_j^*|$; using both parts of Lemma 5.3 and the fact that both $x \in C_i^*, y \in C_j^*$ are good, we have $d(x, c_j) \leq \frac{\alpha}{60} \frac{1}{\epsilon} \frac{w}{|C_j^*|}$, and $d(x, c_j) \geq w_2(x) > \frac{\alpha}{4} \frac{1}{\epsilon} \frac{w}{|C_j^*|}$, so

$$(5.1) \quad d(x, y) \geq \alpha w \cdot \left(\frac{1}{4} - \frac{1}{60} \right) \frac{1}{\epsilon} \frac{w}{|C_j^*|} > \frac{\alpha}{5} \frac{1}{\epsilon} \frac{w}{\min(|C_i^*|, |C_j^*|)}$$

where we use that $|C_j^*| = \min(|C_i^*|, |C_j^*|)$. Part (iii) follows from Lemma 5.3 and the trivial union bound. ■

While Lemma 5.4 is similar in spirit to Lemma 3.2, there is a crucial difference: the distance between the good points in X_i and X_j is no longer lower bounded by some absolute value τ , but rather the bound depends on the sizes of X_i and X_j . However, a redeeming feature is that the separation between them is large compared to the diameters of both X_i and X_j ; we will use this feature crucially in our algorithm.

5.2 Algorithm for Min-Sum/Balanced- k -Median Clustering

For the algorithm below, define *critical thresholds* $\tau_0, \tau_1, \tau_2, \dots$ as: $\tau_0 = 0$ and τ_i is the i th smallest distinct distance $d(x, y)$ for $x, y \in S$. Thus, $G_{\tau_0}, G_{\tau_1}, \dots$ are the only distinct threshold graphs possible.

THEOREM 5.1. *If the balanced k -median instance satisfies the $(1 + \alpha, \epsilon)$ -property and we are given the value of w , then so long as the smallest correct cluster has size greater than $(6 + 120/\alpha)\epsilon n$, Algorithm 2 produces a clustering that is $O(\epsilon/\alpha)$ -close to the target. If we are not given w , then we can use Algorithm 2 as a subroutine to produce a clustering that is $O(\epsilon/\alpha)$ -close to the target.*

Algorithm 2 Balanced k -median Algorithm

Input: (\mathcal{M}, S) , w , $\epsilon \leq 1$, $\alpha > 0$, $k, b := (2 + 60/\alpha)\epsilon n$.

Let the initial threshold $\tau = \tau_0$.

Step 1: If $k = 0$ or $S = \emptyset$, stop.

Step 2: Construct the τ -threshold graph G_τ on the current set S of points.

Step 3: Create a new graph H by connecting two points by in S an edge if they share at least b neighbors in common in G_τ .

Step 4: Let C be largest connected component in H . **If** $|C| \geq \frac{1}{20} \frac{\alpha w}{\epsilon \tau}$, **then** output C as a cluster, set $k \leftarrow k - 1$, $S \leftarrow S \setminus C$, and go to Step 1. **Else** increase τ to the next critical threshold and go to Step 1.

Proof: Since each cluster in the target clustering has more than a $(6 + 120/\alpha)\epsilon$ fraction of the points by the assumption, the optimal balanced- k -median clustering \mathcal{C}^* must differ from the target clustering by fewer than ϵn points, and hence each cluster in \mathcal{C}^* must have at least $(5 + 120/\alpha)\epsilon n$ points. Moreover, by Lemma 5.3, the bad points B constitute at most $(2 + 60/\alpha)\epsilon$ fraction of points, and hence each $|X_i| = |C_i^* \setminus B| > (3 + 60/\alpha)\epsilon n \geq (2 + 60/\alpha)\epsilon n + 2 = b + 2$.

Assume we know w . Consider what happens in the execution of the algorithm: as we increase τ , the sizes of the H -components increase (since we are adding more edges in G_τ). This happens until the largest H -component is “large enough” (i.e., the condition in Step 4 gets satisfied), in which case we output it and then go back to raising τ .

We claim that every time we output a cluster in Step 4, this cluster completely contains some X_i and includes no

points in any $X_{j \neq i}$. More specifically, we show that as we increase τ , the condition in Step 4 will be satisfied *after* all the good points in the some cluster have been fully connected, but *before* any edges appear between good points in different clusters. It suffices to show that the first cluster output by the algorithm contains some X_i entirely; the claim for the subsequent output clusters is the same. Assume that $|C_1^*| \geq |C_2^*| \geq \dots \geq |C_k^*|$, and let $n_i = |C_i^*|$. Define $d_i = \frac{\alpha w}{30 \epsilon} \frac{1}{|C_i^*|}$ and recall that $\min_{x,y \in X_i} d(x,y) \leq d_i$.

We first claim that as long as $\tau \leq 3d_1$, no two points belonging to different X_i 's can lie in the same H -component. Since the distance between points in any X_i and $X_{j \neq i}$ is strictly greater than $\frac{\alpha}{5} \frac{1}{\epsilon} \frac{w}{\min(|C_i^*|, |C_j^*|)} \geq 2\tau$ for any $\tau \leq 3d_1$, every $x \in X_i$ and $y \in X_j$ share no common neighbors; hence, by an argument identical to that in Lemma 3.3, x and y belong to different components of H .

Next, we claim that for values of $\tau < \min\{d_i, 3d_1\}$, the H -component containing X_i cannot be output by Step 4. Indeed, since $\tau < 3d_1$, no X_i and X_j belong to the same H -component by the above claim, and hence any H -component containing points from X_i has size at most $|C_i^*| + |B| < \frac{3n_i}{2}$; however, the minimum size bound $\frac{1}{20} \frac{\alpha w}{\epsilon \tau} > \frac{3n_i}{2}$ for values of $\tau < d_i$, and hence the condition of Step 4 is not satisfied. Note that when $\tau \geq d_i$, all the points of X_i lie in the same H -component.

Finally, we show that the condition in Step 4 becomes true for some H -component fully containing some X_i for some value $\tau = [d_1, 3d_1]$. (By the argument in the previous paragraph, $\tau \geq d_i$, and hence the output component will fully contain X_i .) For the sake of contradiction, suppose not. But note at time $\tau = 3d_1$, at least the H -component containing X_1 has size at least $|C_1^*| - |B| > n_1/2$ and will satisfy the condition (which at time $\tau = 3d_1$ requires a cluster of size $\frac{1}{20} \frac{\alpha w}{\epsilon} \frac{30 \epsilon n_i}{3 \alpha w} = n_1/2$), giving the contradiction.

To recap, by time $3d_1$ none of the clusters have merged together, and the Step 4 condition was satisfied for at least the component containing X_1 (and hence for the largest component) at some time prior to that. Moreover, this largest component must fully contain some set X_i and no points in $X_{j \neq i}$. Finally, we can iterate the same argument on the set $S \setminus X_i$ to complete the proof for the case when we know w .

The case when we do not know w . In this case, we do not want to use a β -approximation algorithm for balanced k -median to obtain a clustering that is $O(\beta\epsilon/\alpha)$ -close to the target, because the balanced- k -median (and minsum clustering) problems only have a logarithmic approximation for arbitrary k , and hence our error would blow up by a logarithmic factor. Instead, we use the idea of trying increasing values of w : we then stop the first time when we output k clusters that cover at least $n - b = (1 - O(\epsilon/\alpha))n$ of the points in S . Clearly, if we reached the correct value of w we would succeed in covering all the good $n - b$ points using our k clusters; we now argue that we will never mistakenly output a high-error clustering.

The argument is as follows. Let us say we *mark* X_i the first time we output a cluster containing at least one point

from it. There are three possible sources of mistakes: (a) we may output a cluster prematurely, it may contain some but not all points from X_i , (b) we may output a cluster which contains points from one or more previously marked sets X_j (but no unmarked X_i), or (c) we may output a cluster with points from an unmarked X_i and one or more previously marked X_j . In case (a), if we end up with all but $O(\epsilon/\alpha)$ -fraction of the points, we did not miss too many points from the X_i 's, so our error is $O(\epsilon/\alpha)$. In case (b), we use up too many clusters and would end with missing some X_i completely, which would result in more than b unclustered points, and we would try a larger guess for w . The dangerous case is case (c), but we claim case (c) in fact cannot happen. Indeed, the value of τ at which we would form connected components containing points from both X_i and X_j is a constant times larger than the value $\tau_{<}$ at which all of X_i would be in a single H -component. Moreover, since our guess for w is too small, this H -component would certainly satisfy the condition of Step 4 and be output as a cluster instead. ■

6 Conclusions and Open Questions

A concrete open question is designing an efficient algorithm for the min-sum property which works in the presence of small target clusters. Another natural direction for investigation is designing faster algorithms for all the properties analyzed in this paper. The case of large clusters can be handled by using standard sampling ideas [MOP01, CS04, BD07], however these techniques do not seem to immediately apply in the case where the target clusters are small.

More broadly, it would be interesting to further explore and analyze other commonly used clustering objective functions in our framework.

References

- [AGK⁺04] V. Arya, N. Garg, R. Khandekar, A. Meyerson, K. Munagala, and V. Pandit. Local search heuristics for k -median and facility location problems. *SIAM J. Comput.*, 33(3):544–562, 2004.
- [AK05] S. Arora and R. Kannan. Learning mixtures of arbitrary gaussians. In *Proc. 37th STOC*, 2005.
- [AM05] D. Achlioptas and F. McSherry. On spectral learning of mixtures of distributions. In *COLT*, 2005.
- [ARR99] S. Arora, P. Raghavan, and S. Rao. Approximation schemes for Euclidean k -medians and related problems. In *STOC*, pages 106–113. 1999.
- [AV06] D. Arthur and S. Vassilvitskii. Worst-case and smoothed analyses of the icp algorithm, with an application to the k -means method. In *Proc. 47th FOCS*, 2006.
- [BBV08] M.-F. Balcan, A. Blum, and S. Vempala. A discriminative framework for clustering via similarity functions. In *Proc. 40th STOC*, 2008.
- [BCR01] Y. Bartal, M. Charikar, and D. Raz. Approximating minsum k -clustering in metric spaces. In *Proc. 33rd STOC*, 2001.
- [BD07] S. Ben-David. A framework for statistical clustering with constant time approximation for k -median and k -means clustering. *Mach. Learn.*, 66(2-3):243 – 257, 2007.

- [CG99] M. Charikar and S. Guha. Improved combinatorial algorithms for the facility location and k-median problems. In *Proc. 4th FOCS*, 1999.
- [CGTS99] M. Charikar, S. Guha, E. Tardos, and D. B. Shmoy. A constant-factor approximation algorithm for the k-median problem. In *STOC*, 1999.
- [CS04] A. Czumaj and C. Sohler. Sublinear-time approximation for clustering via random samples. In *Proc. 31st ICALP*, pages 396–407, 2004.
- [Das99] S. Dasgupta. Learning mixtures of gaussians. In *Proc. 40th FOCS*, 1999.
- [DGL96] L. Devroye, L. Györfi, and G. Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer-Verlag, 1996.
- [DHS01] R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern Classification*. Wiley, 2001.
- [dlVKKR03] W. Fernandez de la Vega, Marek Karpinski, Claire Kenyon, and Yuval Rabani. Approximation schemes for clustering problems. In *Proc. 35th STOC*, 2003.
- [GK99] S. Guha and S. Khuller. Greedy strikes back: Improved algorithms for facility location. *Journal of Algorithms*, 31(1):228–248, 1999.
- [Ind99] P. Indyk. Sublinear time algorithms for metric space problems. In *Proc. 31st STOC*, 1999.
- [JMS02] K. Jain, M. Mahdian, and A. Saberi. A new greedy approach for facility location problems. In *Proc. 34th STOC*, 2002.
- [JV01] K. Jain and V. V. Vazirani. Approximation algorithms for metric facility location and k-median problems using the primal-dual schema and lagrangian relaxation. *JACM*, 48(2):274–296, 2001.
- [KSS04] A. Kumar, Y. Sabharwal, and S. Sen. A simple linear time $(1 + \epsilon)$ -approximation algorithm for k-means clustering in any dimensions. In *Proc. 45th FOCS*, pages 454–462, Washington, DC, USA, 2004.
- [KSV05] R. Kannan, H. Salmasian, and S. Vempala. The spectral method for general mixture models. In *Proc. 18th COLT*, 2005.
- [Llo82] S.P. Lloyd. Least squares quantization in PCM. *IEEE Trans. Inform. Theory*, 28(2):129–137, 1982.
- [Mei06] M. Meila. The uniqueness of a good clustering for k-means. In *Proc. 23rd ICML*, 2006.
- [MOP01] N. Mishra, D. Oblinger, and L. Pitt. Sublinear time approximate clustering. In *SODA*, pages 439–447, 2001.
- [ORSS06] R. Ostrovsky, Y. Rabani, L. Schulman, and C. Swamy. The effectiveness of lloyd-type methods for the k-means problem. In *Proc. 47th FOCS*, 2006.
- [Sch00] L.J. Schulman. Clustering for edge-cost minimization. In *Proc. STOC*, pages 547–555, 2000.
- [VW04] S. Vempala and G. Wang. A spectral algorithm for learning mixture models. *JCSS*, 68(2):841–860, 2004.

A Appendix

THEOREM A.1. *For any $1 \leq c_1 < c_2$, any $\epsilon, \alpha > 0$, there exists a family of metric spaces G and target clusterings that satisfy the (c_1, ϵ) property for the k-median objective (likewise, k-means and min-sum) and yet do not satisfy even the $(c_2, 1/2 - \alpha)$ property for that objective.*

Proof: We focus first on the k-median objective. Consider a set of n points such that the target clustering consists of one cluster C_1 with $n(1 - 2\alpha)$ points all at the same location ($d(u, v) = 0$ for all $u, v, \in C_1$) and $k - 1$ clusters C_2, \dots, C_k

each consisting of $\frac{2\alpha n}{k-1}$ points, all at distance 1. The distance between points in any two distinct clusters C_i, C_j for $i, j \geq 2$ is D , where $D > 1$ will be defined below. Points in C_1 are at distance greater than $c_2 n$ from any of the other clusters.

In this construction, the target clustering is the optimal k-median solution, and has a total k-median cost of $2\alpha n - (k - 1)$. We now define D so that there (just barely) exists a c_2 approximation that splits cluster C_1 . In particular, consider the solution that merges C_2 and C_3 into a single cluster (C_4, \dots, C_k will each be their own cluster) and uses 2 clusters to evenly split C_1 . This clearly has error at least $1/2 - \alpha$, and furthermore this solution has a cost of $\frac{2\alpha n}{k-1} D + (2\alpha n - \frac{2\alpha n}{k-1} - (k - 2))$, and we define D to set this equal to $c_2(2\alpha n - (k - 1)) = c_2 \text{OPT}$.

Any c_1 approximation, however, must be ϵ -close to the target for $k > 1 + 2\alpha/\epsilon$. In particular, by definition of D , any c_1 -approximation cannot merge two C_i, C_j into a single cluster: it must have one median inside each C_i , and can have error on fewer than $\frac{2\alpha n}{k-1}$ points. This is less than ϵn by definition of k .

The same construction, with D defined appropriately, applies to k-means and min-sum objectives as well. ■

THEOREM A.2. *For k-median, k-means, and min-sum objectives, the problem of finding a c-approximation can be reduced to the problem of finding a c-approximation under the (c, ϵ) assumption. Therefore, the problem of finding a c-approximation under the (c, ϵ) assumption is as hard as the problem of finding a c-approximation in general.*

Proof: Given a metric G with n nodes and a value k (a generic instance of the clustering problem) we construct a new instance satisfying the (c, ϵ) assumption. In particular we create a new graph G' by adding an extra n/ϵ nodes that are all very far away from each other and from the nodes in G (call this distance D). We now let $k' = k + n/\epsilon$ and define the target clustering to be the optimal (k-median, k-means, or min-sum) solution on G , together with each of the points in $G' \setminus G$ in its own singleton cluster.

We first claim that G' satisfies the (c, ϵ) property. This is because for sufficiently large D , any solution that does not put each of the new nodes into its own singleton cluster will incur too high a cost. So a c-approximation can only differ from the target on G (which has less than an ϵ fraction of the nodes). Furthermore, a c-approximation in G' yields a c-approximation in G because the singleton clusters do not contribute to the overall cost. ■

The following shows that unlike the $(1.01, \epsilon)$ -property, obtaining an $O(\epsilon)$ -close clustering is NP-hard under the $(1, \epsilon)$ -property.

THEOREM A.3. *For any constant c' , for any $\epsilon < 1/(c'\epsilon)$, it is NP-hard to find a clustering of error at most $c'\epsilon$ for the k-median and k-means problem under the $(1, \epsilon)$ -property.*

Proof Sketch: We start from the hard instances of k-median arising from max-k-coverage (using edges of cost 1 and 3): the reduction implies that it is hard to distinguish

cases when there are k medians that cover all the points at distance 1 (the “yes” case), from instances where any set of k medians covers at least a $(1/e - \delta)$ -fraction of the points at distance 3 (the “no” case) for any constant $\delta > 0$. Let us add infinitesimal noise to make a unique optimal solution and call this the target; the uniqueness of the optimal solution ensures that we satisfy the $(1, \epsilon)$ assumption.

Now, in the “yes” case, any clustering with error $c'\epsilon$ will have cost at most $n[(1 - c'\epsilon) + 3c'\epsilon]$. This is less than the cost of the optimal solution in the “no” case (which is still at least $n[(1 - 1/e + \delta) + 3(1/e - \delta)]$) as long as $c'\epsilon \leq 1/e - \delta$, and would allow us to distinguish the “yes” and “no” instances. This completes the proof for the k -median case, and the proof can be altered slightly to work for the k -means problem as well. ■

LEMMA A.1. *Let $\mathcal{C} = \{C_1, \dots, C_k\}$ be a k -clustering in which each cluster is nonempty, and let $R = \{(x_1, j_1), (x_2, j_2), \dots, (x_t, j_t)\}$ be a set of t reassignments of points to clusters (assume that $x_i \notin C_{j_i}$ for all i). Then there must exist a set $R' \subseteq R$ of size at least $t/3$ such that the clustering \mathcal{C}' produced by reassigning points in R' has distance at least $\frac{1}{n}|R'|$ from \mathcal{C} .*

Before proving the lemma, note that we cannot necessarily just choose $R' = R$ because, for instance, it could be that R moves all points in C_1 to C_2 and all points in C_2 to C_1 : in this case, performing all reassignments in R produces the exact same clustering as we started with (just with different indices). Instead, we need to ensure that each reassignment in R' has an associated certificate ensuring that if implemented, it will increase the resulting distance from \mathcal{C} . Note also that if \mathcal{C} consists of 3 singleton clusters: $C_1 = \{x\}, C_2 = \{y\}, C_3 = \{z\}$, and if $R = \{(x, 2), (y, 3), (z, 1)\}$, then any subset of reassignments in R will produce a clustering that differs in at most one element from \mathcal{C} ; thus, the factor of 3 is tight.

Proof: The proof is based on the following lower-bounding technique. Given two clusterings \mathcal{C} and \mathcal{C}' , suppose we can produce a list L of disjoint subsets S_1, S_2, \dots , such that for each i , all points in S_i are in the same cluster in one of \mathcal{C} or \mathcal{C}' and they are all in different clusters in the other. Then \mathcal{C} and \mathcal{C}' must have distance at least $\frac{1}{n} \sum_i (|S_i| - 1)$. In particular, any bijection σ on the indices can have agreement between \mathcal{C} and \mathcal{C}' on at most one point from each S_i .

We construct R' and witness-list L as follows. While there exists a reassignment $(x, j) \in R$ such that x is in a cluster $C(x)$ with at least 3 points: choose an arbitrary point $y \in C(x)$ and add $\{x, y\}$ to L , add (x, j) to R' , and remove (x, j) from R as well as the reassignment involving y if one exists. In addition, remove x and y from the point set S . This process guarantees that all pairs added to L are disjoint, and we remove at most twice as many reassignments from R as we add to R' . (So, if R becomes empty, we will have achieved our desired result with $|R'| = t/2$). Moreover, because we only perform this step if $|C(x)| \geq 3$, this process does not produce any empty clusters.

We now have that for all reassignments $(x, j) \in R$, x is in a singleton or doubleton cluster. Let R_{single} be the set of reassignments $(x, j) \in R$ such that x is in a singleton cluster. Viewing these reassignments as directed edges, R_{single} forms a graph on the clusters C_i where each node has outdegree ≤ 1 . Therefore, each component of this graph must be an arborescence with possibly one additional edge from the root. We now proceed as follows. While R_{single} contains a source (a node of outdegree 1 and indegree 0), choose an edge (x, j) such that (a) x is a source and (b) for all other edges (y, j) , y is either a source or part of a cycle. We then consider two cases:

1. Node j is not a sink in R_{single} : that is, there exists an edge $(z, j_z) \in R_{single}$ for $z \in C_j$. In this case, we add to R' the edge (x, j) and all other edges (y, j) such that y is a source, and we remove from R (and from R_{single}) the edges (z, j_z) , (x, j) , and all edges (y, j) (including the at most one edge (y, j) such that y is part of a cycle). We then add to L the set $\{x\} \cup \{z\} \cup \{y : (y, j) \text{ was just added to } R'\}$ and remove these points from S . Note that the number of edges removed from R is at most the number of edges added to R' plus 2, giving a factor of 3 in the worst case. Note also that we maintain the invariant that no edges in R_{single} point to empty clusters, since we deleted all edges into C_j , and the points x and y added to L were sources in R_{single} .
2. Otherwise, node j is a sink in R_{single} . In this case, we add to R' the edge (x, j) along with all other edges $(y, j) \in R_{single}$ (removing those edges from R and R_{single}). We choose an arbitrary point $z \in C_j$ and add to L the set $\{x\} \cup \{z\} \cup \{y : (y, j) \text{ was just added to } R'\}$, removing those points from S . In addition, we remove from R all (at most two) edges exiting from C_j (we are forced to remove any edge exiting from z since z was added to L , and there might be up to one more edge if C_j is a doubleton). Again, the number of edges removed from R is at most the number of edges added to R' plus 2, giving a factor of 3 in the worst case.

At this point, if R_{single} is nonempty, its induced graph must be a collection of disjoint cycles. For each such cycle, we choose every other edge (half the edges in an even-length cycle, at least $1/3$ of the edges in an odd cycle), and for each edge (x, j) selected, we add (x, j) to R' , remove (x, j) and (z, j_z) for $z \in C_j$ from R and R_{single} , and add the pair $\{x, z\}$ to L .

Finally, R_{single} is empty and we finish off any remaining doubleton clusters using the same procedure as in the first part of the argument. Namely, while there exists a reassignment $(x, j) \in R$, choose an arbitrary point $y \in C(x)$ and add $\{x, y\}$ to L , add (x, j) to R' , and remove (x, j) from R as well as any reassignment involving y if one exists.

By construction, the set R' has size at least $|R|/3$, and the set L ensures that each reassignment in R' increases the resulting distance from \mathcal{C} as desired. ■