# Approximate Clustering without the Approximation

MARIA-FLORINA BALCAN
School of Computer Science, Georgia Institute of Technology
AVRIM BLUM
Computer Science Department, Carnegie Mellon University
ANUPAM GUPTA
Computer Science Department, Carnegie Mellon University

The design of approximation algorithms for clustering points in metric spaces is a flourishing area of research, with much 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 fueled in part by the implicit hope that better approximations will also yield more accurate clusterings. In particular, for problems such as clustering proteins by function, or clustering images by subject, the true aim is to match 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 to this target in terms of how the data points are clustered.

In this paper, we show that if we make this implicit assumption explicit—that is, if we assume our instance satisfies the condition that any $c$-approximation to the given clustering objective $\Phi$ is $\epsilon$-close to the target—then we can efficiently produce clusterings that are $O(\epsilon)$-close to the target, *even for values $c$ for which obtaining a $c$-approximation is NP-hard.* We call this condition $(c, \epsilon)$-*approximation-stability* and show that we can achieve this guarantee for any constant $c > 1$, for three important and widely studied objective functions, namely $k$-median, $k$-means, and min-sum objectives. Thus, we can perform nearly as well in terms of distance to the target clustering *as if* we could approximate the objective to this NP-hard value. That is, under the assumption that $c$-approximations are good solutions, we can approximate the target even if in general it is hard to achieve a $c$-approximation to the objective.

Categories and Subject Descriptors: F.2.0 [**Analysis of Algorithms and Problem Complexity**]: General; I.5.3 [**Pattern Recognition**]: Clustering

General Terms: Algorithms, Theory.

Additional Key Words and Phrases: Clustering, Approximation Algorithms, $k$-median, $k$-means, min-sum, Clustering Accuracy.

## 1.  INTRODUCTION

Problems of clustering data are ubiquitous throughout science. They arise in many different fields, from computational biology where such problems include clustering protein sequences by function, to computer vision where one may want to cluster images by subject, to information retrieval, including problems of clustering documents or search results by topic, just to name a few. A wide range of algorithms have been developed, and the study of clustering methods continues to be a highly active area of research both experimentally and theoretically.

From a theoretical perspective, one of the primary approaches to attacking the problem of clustering has been the development of algorithms for (approximately) optimizing various natural distance-based objective functions. These include the $k$-median problem (partition data into $k$ clusters $C_i$, giving each a center $c_i$, to minimize the sum of distances of all datapoints to the centers of their cluster), the $k$-means problem (find $k$ clusters $C_i$ and centers $c_i$ that minimize the sum of squares of distances of all datapoints to the centers of their cluster) and min-sum clustering (find $k$ clusters $C_i$ that minimize the sum of all intra-cluster pairwise distances). The field of approximation algorithms for such objectives is a very active one, with a large number of algorithms having been developed. The $k$-median problem has a $3 + \epsilon$-approximation [Arya et al. 2004], building on a long line of work initiated by [Charikar et al. 1999], and it is known to be NP-hard to approximate to better than $1 + 1/e$ [Jain et al. 2002]. 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 $k$ [Kumar et al. 2004]. The min-sum clustering problem admits an $O(\log^{1+\delta} n)$-approximation for general metric spaces, a constant-factor approximation when $k = o(\log n / \log \log n)$ [Czumaj and Sohler 2007], and a PTAS when $k$ is a constant [de la Vega et al. 2003]. 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 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 close to the true answer in terms of how the data points are clustered, 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 our instance satisfies the condition that any $c$-approximation to the given objective $\Phi$ is $\epsilon$-close to the target clustering in terms of how points are clustered,

then we can in fact use this to efficiently 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.* We call this condition $(c, \epsilon)$-*approximation-stability*, and in particular we achieve this guarantee for any constant $c > 1$ for the $k$-median and $k$-means objectives, as well as for min-sum when the target clusters are sufficiently large compared to $\frac{\epsilon n}{c-1}$. Moreover, if the target clusters are sufficiently large compared to $\frac{\epsilon n}{c-1}$, for $k$-median we can actually get $\epsilon$-close (rather than $O(\epsilon)$-close) to the target. Note that one should view $k$ and $\epsilon$ here as parameters and not as constants: our algorithms will run in time polynomial in the number of points $n$ and the number of clusters $k$.

Thus, our results show that by using the implicit assumptions motivating the desire for improved approximation algorithms in the first place when the true goal is to match an unknown target clustering, and in particular by using the structure of the data these assumptions imply, we can in fact perform nearly as well in terms of the true goals *as if* we could approximate the objective to an NP-hard value. That is, under the assumption that $c$-approximations to the given clustering objective are good solutions, we can approximate the target even if it is hard to approximate the objective, which in the end was a proxy for our real aims.

Our results can also be viewed in the context of a line of work (e.g., [Ostrovsky et al. 2006; Balcan et al. 2008; Bilu and Linial 2010; Kumar and Kannan 2010], see discussion below) giving efficient algorithms for solving clustering problems under natural stability conditions on the data. In that view, our stability condition is that small changes in the objective value around the optimum should yield only small changes in the clustering produced. Of course, our motivation is that this is what is implicitly desired in any event when considering approximation algorithms for these objectives.

Subsequent work has already demonstrated the practicality of our approach for real world clustering problems. For example, Voevodski et al. [2010] show that a variant of the algorithm that we propose for the $k$-median problem provides state of the art results for clustering biological datasets. Our work has inspired a number of other both theoretical and practical exciting subsequent developments and we discuss these further in Section 7.

More broadly, we believe that this approach of making implicit assumptions explicit and using properties they imply, may be useful for other classes of problems in which the objective functions being optimized may not be the same as the true underlying problem goals.

## 1.1 Related Work

**Work on approximation algorithms:** For $k$-*median*, $O(1)$-approximations were first given by Charikar et al. [1999], Jain and Vazirani [2001], and Charikar and Guha [1999] and the best approximation guarantee known is $(3 + \epsilon)$ due to Arya et al. [2004]. A reduction from max-$k$-coverage shows an easy $(1 + 1/e)$-hardness of approximation [Guha and Khuller 1999; Jain et al. 2002]. The $k$-median problem on constant-dimensional Euclidean spaces admits a PTAS [Arora et al. 1999].

For $k$-*means* in 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 $1 + 3/e$

follows from the ideas of [Guha and Khuller 1999; Jain et al. 2002]. 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 $\epsilon$* [Kumar et al. 2004]. Lloyd's local search algorithm [Lloyd 1982] is often used in practice, despite having poor worst-case performance [Arthur and Vassilvitskii 2006]. Ostrovsky et al. [2006] study ways of seeding Lloyd's local search algorithm: they show that on instances satisfying an $\epsilon$-*separation* condition, their seeding results in solutions with provable approximation guarantees. Their $\epsilon$-separation condition has an interesting relation to approximation-stability, which we discuss more fully in Section 6. Essentially, it is a stronger assumption than ours; however, their goal is different—they want to approximate the objective whereas we want to approximate the target clustering. An interesting extension of the $k$-means objective to clusters lying in different subspaces is given in [Agarwal and Mustafa 2004].

*Min-sum $k$-clustering* on general metric spaces admits a PTAS for the case of constant $k$ by Fernandez de la Vega et al. [2003] (see also Indyk [1999]). 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. [2001]. The problem has also been studied in geometric spaces for constant $k$ by Schulman [2000] who gave an algorithm for $(R^d, \ell_2^2)$ 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 $d$. More recently, Czumaj and Sohler have developed a $(4 + \epsilon)$-approximation algorithm for the case when $k$ is small compared to $\log n / \log \log n$ [Czumaj and Sohler 2007].

**Related work on error to a target clustering:** There has also been significant work in machine learning and theoretical computer science on clustering or learning with mixture models [Achlioptas and McSherry 2005; Arora and Kannan 2005; Duda et al. 2001; Devroye et al. 1996; Kannan et al. 2005; Vempala and Wang 2004; Dasgupta 1999]. That work, like ours, has an explicit notion of a correct ground-truth clustering; however, it makes strong probabilistic assumptions about how data points are generated.

Balcan et al. [2008] investigate the goal of approximating a desired target clustering without probabilistic assumptions. They analyze 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 at least one clustering in the list. In relation to implicit assumptions about approximation algorithms, Balcan et al. [2008] show that for $k$-median, the assumption that any 2-approximation is $\epsilon$-close to the target can be used 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 are 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 $c > 1$*.

**Clustering under natural conditions:** In addition to results of Ostrovsky et al. [2006] and Balcan et al. [2008] mentioned above, there has been other work

investigating clustering under natural restrictions on the input. Bilu and Linial [2010] give an algorithm for finding the optimal maxcut (which can be viewed as a 2-clustering problem) under the assumption that the optimal solution is stable to large perturbations in the edge weights (roughly of order $O(n^{2/3})$). Kumar and Kannan [2010] and Awasthi et al. [2010b] consider clustering under proximity conditions that extend the notion of $\epsilon$-separability of Ostrovsky et al. [2006], and give several efficient algorithms.

**Other work on clustering:** There is a large body of work on other topics in clustering such as defining measures of clusterability of data sets, on formulating definitions of good clusterings [Gollapudi et al. 2006], on axiomatizing clustering (in the sense of postulating what natural axioms should a "good clustering algorithm" satisfy), both with possibility and impossibility results [Kleinberg 2002], on "comparing clusterings" [Meila 2003; 2005], and on efficiently testing if a given data set has a clustering satisfying certain properties [Alon et al. 2000]. The main difference between this type of work and our work is that we have an explicit notion of a correct ground-truth clustering of the data points, and indeed the results we are trying to prove are quite different.

The work of Meila [2006] is complementary to ours: it shows sufficient conditions under which $k$-means instances satisfy the property that near-optimal solutions are $\epsilon$-close to the *optimal k*-means solution.

## 2. DEFINITIONS, PRELIMINARIES & FORMAL STATEMENT OF MAIN RESULTS

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 $d : \binom{X}{2} \to R_{\geq 0}$ 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 $S$. For a set $S' \subseteq S$, we define $\mu(S')$ as $|S'|/n$. A $k$-clustering $\mathcal{C}$ is a partition of $S$ into $k$ sets $C_1, C_2, \ldots, C_k$. In this paper, we always assume that there is a *true* or *target* $k$-clustering $\mathcal{C}_T$ for the point set $S$.

### 2.1 The Objective Functions

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 aims to find clusters $C_i$ and points $c_i \in X$ to minimize

$$\Phi_2(\mathcal{C}) = \sum_{i=1}^{k} \sum_{x \in C_i} d(x, c_i)^2,$$

and *min-sum clustering* aims to find clusters $C_i$ to minimize

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

Given a function $\Phi$ and instance $(\mathcal{M}, S)$, let $\mathrm{OPT}_\Phi = \min_\mathcal{C} \Phi(\mathcal{C})$, where the minimum is over all $k$-clusterings of $(\mathcal{M}, S)$. We will typically use $\mathcal{C}^*$ to denote the optimal clustering for the given objective.

## 2.2   Distance between Clusterings

We define the distance $dist(\mathcal{C}, \mathcal{C}')$ between two $k$-clusterings $\mathcal{C} = \{C_1, C_2, \ldots, C_k\}$ and $\mathcal{C}' = \{C'_1, C'_2, \ldots, C'_k\}$ of a point set $S$ as the fraction of points in $S$ on which they disagree under the optimal matching of clusters in $\mathcal{C}$ to clusters in $\mathcal{C}'$; i.e.,

$$dist(\mathcal{C}, \mathcal{C}') = \min_{\sigma \in \mathfrak{S}_k} \frac{1}{n} \sum_{i=1}^{k} |C_i - C'_{\sigma(i)}|,$$

where $\mathfrak{S}_k$ is the set of bijections $\sigma : [k] \to [k]$. Equivalently, $dist(\mathcal{C}, \mathcal{C}_T)$ is the number of mistakes, or 0/1-loss, of $\mathcal{C}$ with respect to $\mathcal{C}_T$ if we view each as a $k$-way classifier, under the best matching between their $k$ class labels.

   We say that two clusterings $\mathcal{C}$ and $\mathcal{C}'$ are $\epsilon$-*close* if $dist(\mathcal{C}, \mathcal{C}') < \epsilon$. Note that if $\mathcal{C}$ and $\mathcal{C}'$ are $\epsilon$-close and all clusters $C_i$ have size at least $2\epsilon n$, then the bijection $\sigma$ minimizing $\frac{1}{n} \sum_{i=1}^{k} |C_i - C'_{\sigma(i)}|$ has the property that for all $i$, $|C_i \cap C'_{\sigma(i)}| \geq |C_i| - (\epsilon n - 1) > \frac{1}{2}|C_i|$. This implies for instance that such $\sigma$ is unique, in which case we call this the *optimal bijection* and we say that $\mathcal{C}$ and $\mathcal{C}'$ *agree* on $x$ if $x \in C_i \cap C'_{\sigma(i)}$ for some $i$, and $\mathcal{C}$ and $\mathcal{C}'$ *disagree* on $x$ otherwise.

## 2.3   $(c, \epsilon)$-approximation-stability

We now present our main definition whose implications we study throughout this paper:

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

   The above property can be viewed as what is implicitly assumed when proposing to use a $c$-approximation for objective $\Phi$ to solve a clustering problem in which the true goal is to classify data points correctly. (Otherwise, why aim to approximately optimize that objective?) Similarly, the motivation for improving a $c_2$ approximation to a $c_1 < c_2$ approximation is that perhaps the data satisfies $(c_1, \epsilon)$-approximation-stability but not $(c_2, \epsilon)$-approximation-stability. In fact, as we point out in Theorem 18, for any $c_1 < c_2$, it is possible for any of the above three objectives to construct a dataset that is $(c_1, \epsilon)$-approximation-stable but not even $(c_2, 0.49)$-approximation-stable for that $\Phi$. Since we will be thinking of $c$ as being only slightly larger than 1 (e.g., assuming that all 1.1-approximations to the $k$-median objective are $\epsilon$-close to $\mathcal{C}_T$), we will often write $c$ as $1 + \alpha$ and look at the implications in terms of the parameters $\alpha$ and $\epsilon$.

   It is important to note that $1/\epsilon$, $1/\alpha$, and $k$ need not be constants. For example, we might have that $\mathcal{C}_T$ consists of $n^{0.1}$ clusters of size $n^{0.9}$, $\epsilon = 1/n^{0.2}$ and $\alpha = 1/n^{0.09}$ (this would correspond to the "large clusters case" of Theorem 8).

   Note that for any $c > 1$, $(c, \epsilon)$-approximation-stability does not require that the target clustering $\mathcal{C}_T$ exactly coincide with the optimal clustering $\mathcal{C}^*$ under objective

$\Phi$. However, it does imply the following simple facts:

FACT 1. *If* $(\mathcal{M}, S)$ *satisfies* $(c, \epsilon)$-*approximation-stability for* $\Phi$, *then:*

(a) *The target clustering* $\mathcal{C}_T$, *and the optimal clustering* $\mathcal{C}^*$ *for* $\Phi$ *are* $\epsilon$-*close.*
(b) *The distance between* $k$-*clusterings itself forms a metric space. Hence* $(c, \epsilon')$-*approximation-stability with respect to the target clustering* $\mathcal{C}_T$ *implies* $(c, \epsilon + \epsilon')$-*approximation-stability 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.

Finally, we will often want to take some clustering $\mathcal{C}$, reassign some $\tilde{\epsilon}n$ points to different clusters to produce a new clustering $\mathcal{C}'$, and then argue that $dist(\mathcal{C}, \mathcal{C}') = \tilde{\epsilon}$. As mentioned above, if all clusters of $\mathcal{C}$ have size at least $2\tilde{\epsilon}n$, then it is clear that no matter how $\tilde{\epsilon}n$ points are reassigned, the optimal bijection $\sigma$ between the original clusters and the new clusters is the identity mapping, and therefore $dist(\mathcal{C}, \mathcal{C}') = \tilde{\epsilon}$. However, this need not be so when small clusters are present: for instance, if we reassign all points in $C_i$ to $C_j$ and all points in $C_j$ to $C_i$ then $dist(\mathcal{C}, \mathcal{C}') = 0$. Instead, in this case we will use the following convenient lemma.

LEMMA 2. *Let* $\mathcal{C} = \{C_1, \ldots, C_k\}$ *be a* $k$-*clustering in which each cluster is nonempty, and let* $R = \{(x_1, j_1), (x_2, j_2), \ldots, (x_t, j_t)\}$ *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 exactly* $\frac{1}{n}|R'|$ *from* $\mathcal{C}$.

PROOF. See Appendix A.1. □

## 2.4   Two "Strawman" Approaches

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 $(c, \epsilon)$-approximation-stability for some objective $\Phi$ implied, say, $(2c, 2\epsilon)$-approximation-stability. 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, as noted above, for any $c_1 < c_2$ and any $\epsilon > 0$, 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 that are $(c_1, \epsilon)$-approximation-stable for that objective, and yet that do not satisfy even $(c_2, 0.49)$-approximation-stability (See Appendix, Theorem 18). Thus, the result of a direct application of an arbitrary $c_2$-approximation is nearly as poor as possible.

Second, perhaps $(c, \epsilon)$-approximation-stability implies that finding a factor $c$ approximation is somehow trivial. However, this is not the case either: for any $c > 1$ and $\epsilon > 0$, the problem of finding a $c$-approximation to any of the three objectives we consider under $(c, \epsilon)$-approximation-stability is as hard as finding a $c$-approximation in general (Theorem 19). Note that this reduction requires small clusters. Indeed, as pointed our by [Schalekamp et al. 2010], our $k$-median algorithm for the large-clusters case *is*, as a byproduct, a $c$-approximation.

It is also interesting to note that results of the form we are aiming for are *not possible given only* $(1, \epsilon)$-*approximation-stability*. Indeed, because the standard hardness-of-approximation proof for $k$-median produces a metric in which all pairwise distances lie in a bounded range, the proof also implies that it is NP-hard, given a data set with only the guarantee that the optimal solution is $\epsilon$-close to the target, to find a clustering of error $O(\epsilon)$; see Theorem 20. Thus, our results show that there is a perhaps unexpected conceptual difference between assuming that the *optimal* solution to an objective such as $k$-median is $\epsilon$-close to the target, and assuming that any *approximately optimal* solution is $\epsilon$-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 hard, and yet for the latter case we give efficient algorithms.

### 2.5   Main results and organization of this paper

We present our analysis of the $k$-median objective in Section 3, the $k$-means objective in Section 4, and the min-sum objective in Section 5. Our main results for each of these objectives are as follows.

THEOREM 8 ($k$-Median, Large Clusters Case) *If the given instance* $(\mathcal{M}, S)$ *satisfies* $(1 + \alpha, \epsilon)$-*approximation-stability for the* $k$-*median objective, and each cluster in* $\mathcal{C}_T$ *has size at least* $(4 + 15/\alpha)\epsilon n + 2$, *then we can efficiently produce a clustering that is* $\epsilon$-*close to* $\mathcal{C}_T$.

THEOREM 9 ($k$-Median: General Case) *If the given instance* $(\mathcal{M}, S)$ *satisfies* $(1 + \alpha, \epsilon)$-*approximation-stability for the* $k$-*median objective, then we can efficiently produce a clustering that is* $O(\epsilon + \epsilon/\alpha)$-*close to* $\mathcal{C}_T$.

THEOREM 12 ($k$-Means: General Case) *If the given instance* $(\mathcal{M}, S)$ *satisfies* $(1 + \alpha, \epsilon)$-*approximation-stability for the* $k$-*means objective, then we can efficiently produce a clustering that is* $O(\epsilon + \epsilon/\alpha)$-*close to* $\mathcal{C}_T$.

THEOREM 17 (Min-sum: Large Clusters Case) *If the given instance* $(\mathcal{M}, S)$ *satisfies* $(1 + \alpha, \epsilon)$-*approximation-stability for the min-sum objective, then so long as the smallest cluster in* $\mathcal{C}_T$ *has size greater than* $(6 + 120/\alpha)\epsilon n$, *we can efficiently produce a clustering that is* $O(\epsilon + \epsilon/\alpha)$-*close to* $\mathcal{C}_T$.

We emphasize that our algorithms run in time polynomial in $n$ and $k$ with no dependence on $\alpha$ and $\epsilon$; in particular, $1/\alpha$ and $1/\epsilon$ (and $k$) need not be constants.

For the "large-cluster" case of $k$-means, we also have a weaker version of Theorem 8, where we mark some $O(\epsilon n/\alpha)$ points as "don't know" and cluster the rest with error at most $\epsilon$. That is, while the total error in this case may be more than $\epsilon$, we can explicitly point out all but $\epsilon n$ of the points we may err on (see Theorem 13). As noted earlier, we only give results for the large-cluster case of min-sum clustering, though [Balcan and Braverman 2009] have recently extended Theorem 17 to the case of general cluster sizes.

## 3.   THE $K$-MEDIAN PROBLEM

We now study clustering under $(c, \epsilon)$-approximation-stability for the $k$-median objective. 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 $\epsilon$-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 $(c, \epsilon)$-approximation-stability 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 $(c, \epsilon)$-approximation-stability

Given a clustering instance specified by a metric space $\mathcal{M} = (X, d)$ and a set of points $S \subseteq X$, fix an optimal $k$-median clustering $\mathcal{C}^* = \{C_1^*, \ldots, C_k^*\}$, and let $c_i^*$ be the center point (a.k.a. "median") for $C_i^*$. For $x \in S$, define

$$w(x) = \min_i d(x, c_i^*)$$

to 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 $\{c_1^*, c_2^*, \ldots, c_k^*\}$. Also, define

$$w_{avg} = \frac{1}{n} \sum_{i=1}^{n} w(x) = \frac{\text{OPT}}{n}$$

to be the average weight of the points. Finally, let $\epsilon^* = dist(\mathcal{C}_T, \mathcal{C}^*)$. Note that for any $c \geq 1$, $(c, \epsilon)$-approximation-stability for $k$-median implies that $\epsilon^* < \epsilon$.

LEMMA 3. *If the instance* $(\mathcal{M}, S)$ *satisfies* $(1 + \alpha, \epsilon)$*-approximation-stability for the $k$-median objective, then*

(a) *If each cluster in $\mathcal{C}_T$ has size at least $2\epsilon n$, then less than $(\epsilon - \epsilon^*)n$ points $x \in S$ on which $\mathcal{C}_T$ and $\mathcal{C}^*$ agree have $w_2(x) - w(x) < \frac{\alpha w_{avg}}{\epsilon}$.*

(a') *For the case of general cluster sizes in $\mathcal{C}_T$, less than $6\epsilon n$ points $x \in S$ have $w_2(x) - w(x) < \frac{\alpha w_{avg}}{2\epsilon}$.*

*Also, for any $t > 0$ we have:*

(b) *At most $t(\epsilon n/\alpha)$ points $x \in S$ have $w(x) \geq \frac{\alpha w_{avg}}{t\epsilon}$.*

PROOF. To prove Property (a), assume to the contrary. Then one could take $\mathcal{C}^*$ and move $(\epsilon - \epsilon^*)n$ points $x$ on which $\mathcal{C}_T$ and $\mathcal{C}^*$ agree to their second-closest clusters, increasing the objective by at most $\alpha$OPT. Moreover, this new clustering $\mathcal{C}' = \{C_1', \ldots, C_k'\}$ has distance at least $\epsilon$ from $\mathcal{C}_T$, because we begin at distance $\epsilon^*$ 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 $\mathcal{C}'$ remains the same as the optimal bijection between $\mathcal{C}_T$ and $\mathcal{C}^*$). Hence we have a clustering that is not $\epsilon$-close to $\mathcal{C}_T$ with cost only $(1+\alpha)$OPT, a contradiction.

For Property (a'), we use Lemma 2. Specifically, assuming for contradiction that $6\epsilon n$ points satisfy (a'), Lemma 2 states that we can find a subset of $2\epsilon n$ of them such that starting from $\mathcal{C}^*$, for each one that we move to its second-closest cluster, the distance from $\mathcal{C}^*$ increases by $\frac{1}{n}$. Therefore, we can create a clustering $\mathcal{C}'$ that is distance at least $2\epsilon$ from $\mathcal{C}^*$ while increasing the objective by at most $\alpha$OPT;

by Fact 1(b) this clustering $\mathcal{C}'$ is not $\epsilon$-close to $\mathcal{C}_T$, thus contradicting $(1 + \alpha, \epsilon)$-approximation-stability. Property (b) simply follows from the definition of the average weight $w_{avg}$, and Markov's inequality. $\square$

**Notation:** For the case that each cluster in $\mathcal{C}_T$ has size at least $2\epsilon n$, define the *critical distance* $d_{crit} = \frac{\alpha w_{avg}}{5\epsilon}$, else define $d_{crit} = \frac{\alpha w_{avg}}{10\epsilon}$; note that these quantities are $1/5$ times the values in properties (a) and (a') respectively of Lemma 3.

DEFINITION 2. *Define point $x \in S$ to be* good *if both $w(x) < d_{crit}$ and $w_2(x) - w(x) \geq 5d_{crit}$, else $x$ is called* bad. *Let $X_i \subseteq C_i^*$ be the good points in the optimal cluster $C_i^*$, and let $B = S \setminus (\cup X_i)$ be the bad points.*

PROPOSITION 4. *If the instance $(\mathcal{M}, S)$ satisfies $(1+\alpha, \epsilon)$-approximation-stability for the $k$-median objective, then*

*(i) If each cluster in $\mathcal{C}_T$ has size at least $2\epsilon n$, $\mu(B) < (1 + 5/\alpha)\epsilon$.*
*(ii) For the case of general cluster sizes in $\mathcal{C}_T$, $\mu(B) < (6 + 10/\alpha)\epsilon$.*

PROOF. By Lemma 3(a), the number of points on which $\mathcal{C}^*$ and $\mathcal{C}_T$ agree where $w_2(x) - w(x) < 5d_{crit}$ is at most $(\epsilon - \epsilon^*)n$, and there can be at most $\epsilon^* n$ additional such points where $\mathcal{C}^*$ and $\mathcal{C}_T$ disagree. Setting $t = 5$ in Lemma 3(b) bounds the number of points that have $w(x) \geq d_{crit}$ by $(5\epsilon/\alpha)n$. Now recalling that $\mu(B)$ is defined to be $|B|/n$ proves (i). The proof of (ii) similarly follows from Lemma 3(a'), and applying Lemma 3(b) with $t = 10$. $\square$

DEFINITION 3 (THRESHOLD GRAPH). *Define the $\tau$-threshold graph $G_\tau = (S, E_\tau)$ to be the graph produced by connecting all pairs $\{x, y\} \in \binom{S}{2}$ with $d(x, y) < \tau$.*

LEMMA 5 (THRESHOLD GRAPH LEMMA). *For an instance satisfying $(1+\alpha, \epsilon)$-approximation-stability and $\tau = 2d_{crit}$, the threshold graph $G_\tau$ 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 $G_\tau$.*

PROOF. For part (i), since $x, y$ are both good, they are at distance less than $d_{crit}$ to their common cluster center $c_i^*$, by definition. Hence, by the triangle inequality, the distance

$$d(x, y) \leq d(x, c_i^*) + d(c_i^*, y) < 2 \times d_{crit} = \tau.$$

For part (ii), note that the distance from any good point $x$ to any other cluster center, and in particular to $y$'s cluster center $c_j^*$, is at least $5d_{crit}$. Again by the triangle inequality,

$$d(x, y) \geq d(x, c_j^*) - d(y, c_j^*) \geq 5d_{crit} - d_{crit} = 2\tau.$$

Since each edge in $G_\tau$ is between points at distance less than $\tau$, the points $x, y$ cannot share any common neighbors. $\square$

Hence, the graph $G_\tau$ for the above value of $\tau$ 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 structure to find a clustering of error at most $\epsilon$ 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.

LEMMA 6. *Given a graph $G = (S, E)$ satisfying properties (i), (ii) of Lemma 5 and given $b \geq |B|$ such that each $|X_i| \geq b + 2$, there is a deterministic polynomial-time 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 there is no edge between $x, y$ in this new graph $H$. 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 at most $b$. Moreover, at least one of $x$ and $y$ must itself belong to $B$ for them to have any common neighbors at all (again 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 $H$. Note that the component containing some $X_i$ may also contain some vertices from $B$; morover, there may also be components in $H$ that only contain vertices from $B$. But since the $X_i$'s are larger than $B$, we can obtain the claimed clustering by taking the largest $k$ components in $H$, and adding the vertices of all other smaller components in $H$ to any of these, and using this as the $k$-clustering. □

We now show how we can use Lemma 6 to find a clustering that is $\epsilon$-close to $\mathcal{C}_T$ when all clusters are large. For simplicity, we begin by assuming that we are given the value of $w_{avg} = \frac{\text{OPT}}{n}$, and then we show how this assumption can be removed.

THEOREM 7 (LARGE CLUSTERS, KNOWN $w_{avg}$). *If the given instance $(\mathcal{M}, S)$ satisfies $(1 + \alpha, \epsilon)$-approximation-stability for the $k$-median objective and each cluster in $\mathcal{C}_T$ has size at least $(3 + 10/\alpha)\epsilon n + 2$, then given $w_{avg}$ we can efficiently find a clustering that is $\epsilon$-close to $\mathcal{C}_T$.*

PROOF. Let us define $b := (1 + 5/\alpha)\epsilon n$. By assumption, each cluster in the target clustering has at least $(3 + 10/\alpha)\epsilon n + 2 = 2b + \epsilon n + 2$ points. Since the *optimal $k$-median clustering $\mathcal{C}^*$* differs from the target clustering by at most $\epsilon^* n \leq \epsilon n$ points, each cluster $C_i^*$ in $\mathcal{C}^*$ must have at least $2b + 2$ points. Moreover, by Proposition 4(i), the bad points $B$ have $|B| \leq b$, and hence for each $i$,

$$|X_i| = |C_i^* \setminus B| \geq b + 2.$$

Now, given $w_{avg}$, we can construct the graph $G_\tau$ with $\tau = 2d_{crit}$ (which we can compute from the given value of $w_{avg}$), and apply Lemma 6 to find a $k$-clustering $\mathcal{C}'$ where each $X_i$ is contained within a distinct cluster. Note that this clustering $\mathcal{C}'$ differs from the optimal clustering $\mathcal{C}^*$ only in the bad points, and hence, $dist(\mathcal{C}', \mathcal{C}_T) \leq \epsilon^* + \mu(B) \leq O(\epsilon + \epsilon/\alpha)$. However, our goal is to get $\epsilon$-close to the target, which we do as follows.

Call a point $x$ "red" if it satisfies condition (a) in Lemma 3 (i.e., $w_2(x) - w(x) < 5d_{crit}$), "yellow" if it is not red but satisfies condition (b) in Lemma 3 with $t = 5$ (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 $\mathcal{C}' = \{C'_1, \ldots, C'_k\}$ and recall that $\mathcal{C}'$ agrees with $\mathcal{C}^*$ on the green points, so without loss of generality we may assume $X_i \subseteq C'_i$. We now construct a new clustering $\mathcal{C}''$ that agrees with $\mathcal{C}^*$ on both the green and yellow points. Specifically, for each point $x$ and each cluster $C'_j$, compute the median distance $d_{median}(x, C'_j)$ between $x$ and all points in $C'_j$; then insert $x$ into the cluster $C''_i$ for $i = \operatorname{argmin}_j d_{median}(x, C'_j)$. 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 $\mathcal{C}^*$ 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 $\mathcal{C}^*$ 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 $\mathcal{C}'$ has a strict majority of green points (even with point $x$ removed) all of which are clustered as in $\mathcal{C}^*$, this means that for a non-red point $x$, the median distance to points in its correct cluster with respect to $\mathcal{C}^*$ is less than the median distance to points in any incorrect cluster. Thus, $\mathcal{C}''$ agrees with $\mathcal{C}^*$ on all non-red points. Finally, since there are at most $(\epsilon - \epsilon^*)n$ red points on which $\mathcal{C}_T$ and $\mathcal{C}^*$ agree by Lemma 3—and $\mathcal{C}''$ and $\mathcal{C}_T$ might disagree on all these points—this implies

$$dist(\mathcal{C}'', \mathcal{C}_T) \leq (\epsilon - \epsilon^*) + \epsilon^* = \epsilon$$

as desired. For convenience, the above procedure is given as Algorithm 1 below. □

---

**Algorithm 1** $k$-median Algorithm: Large Clusters (given a guess $w$ of $w_{avg}$)

---

**Input:** $w$, $\epsilon \leq 1$, $\alpha > 0$, $k$.
**Step 1:** Construct the $\tau$-threshold graph $G_\tau$ with $\tau = 2d_{crit} = \frac{1}{5}\frac{\alpha w}{\epsilon}$.
**Step 2:** Apply the algorithm of Lemma 6 to find an initial clustering $\mathcal{C}'$. Specifically, construct graph $H$ by connecting $x, y$ if they share at least $b = (1+5/\alpha)\epsilon n$ neighbors in $G_\tau$ and let $C'_1, \ldots, C'_k$ be the $k$ largest components of $H$.
**Step 3:** Produce clustering $\mathcal{C}''$ by reclustering according to smallest median distance in $\mathcal{C}'$. That is, $\mathcal{C}''_i = \{x : i = \operatorname{argmin}_j d_{median}(x, C'_j)\}$.
**Step 4:** Output the $k$ clusters $\mathcal{C}''_1, \ldots, \mathcal{C}''_k$.

---

We now extend the above argument to the case where we are not given the value of $w_{avg}$.

THEOREM 8 (LARGE CLUSTERS, UNKNOWN $w_{avg}$). *If the given instance $(\mathcal{M}, S)$ satisfies $(1+\alpha, \epsilon)$-approximation-stability for the $k$-median objective, and each cluster in $\mathcal{C}_T$ has size at least $(4 + 15/\alpha)\epsilon n + 2$, then we can efficiently produce a clustering that is $\epsilon$-close to $\mathcal{C}_T$.*

PROOF. The algorithm for the case that we are not given the value $w_{avg}$ is the following: we run Steps 1 and 2 of Algorithm 1 repeatedly for different guesses $w$ of $w_{avg}$, starting with $w = 0$ (so the graph $G_\tau$ is empty) and at each step increasing $w$ to the next value such that $G_\tau$ contains at least one new edge (so we have at most $n^2$ different guesses to try). If the current value of $w$ causes the $k$ largest components of $H$ to miss more than $b := (1 + 5/\alpha)\epsilon n$ points, or if any of these components has size $\leq b$, then we discard this guess $w$, and try again with the next larger guess for $w$. Otherwise, we run Algorithm 1 to completion and let $\mathcal{C}''$ be the clustering produced.

Note that we still might have $w < w_{avg}$, but this just means that the resulting graphs $G_\tau$ and $H$ can only have fewer edges than the corresponding graphs for the correct $w_{avg}$. Hence, some of the $X_i$'s might not have fully formed into connected components in $H$. However, if the $k$ largest components together miss at most $b$ points, then this implies we must have at least one component for each $X_i$, and therefore exactly one component for each $X_i$. So, 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 $\mathcal{C}'_i$ contains at least $|X_i| - b \geq b + 2$ correctly clustered green points (with respect to $\mathcal{C}^*$) and at most $b$ misclassified points. Therefore, as shown in the proof of Theorem 7, the resulting clustering $\mathcal{C}''$ will correctly cluster all non-red points as in $\mathcal{C}^*$ and so is at distance at most $(\epsilon - \epsilon^*) + \epsilon^* = \epsilon$ from $\mathcal{C}_T$. For convenience, this procedure is given as Algorithm 2 below. □

---

**Algorithm 2** $k$-median Algorithm: Large Clusters (unknown $w_{avg}$)

---

**Input:** $\epsilon \leq 1$, $\alpha > 0$, $k$.
**For** $j = 1, 2, 3 \ldots$ do:
    **Step 1:** Let $\tau$ be the $j$th smallest pairwise distance in $S$. Construct $\tau$-threshold graph $G_\tau$.
    **Step 2:** Run Step 2 of Algorithm 1 to construct graph $H$ and clusters $C'_1, \ldots, C'_k$.
    **Step 3:** If $\min(|C'_1|, \ldots, |C'_k|) > b$ and $|C'_1 \cup \ldots \cup C'_k| \geq n(1 - \epsilon - 5\epsilon/\alpha)$, run Step 3 of Algorithm 1 and output the clusters $C''_1, \ldots, C''_k$ produced.

---

### 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.

---

**Algorithm 3** $k$-median Algorithm: General Case

---

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

**Initialization:** Run a constant-factor $k$-median approximation algorithm to compute a value $w \in [w_{avg}, \beta w_{avg}]$ for, say, $\beta = 4$.

**Step 1:** Construct the $\tau$-threshold graph $G_\tau$ with $\tau = \frac{1}{5}\frac{\alpha w}{\beta \epsilon}$.

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

   Pick the vertex $v_j$ of highest degree in $G_\tau$.

   Remove $v_j$ and its neighborhood from $G_\tau$ and call this cluster $C(v_j)$.

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

---

THEOREM 9 ($k$-MEDIAN: GENERAL CASE). *If the given instance $(\mathcal{M}, S)$ satisfies $(1 + \alpha, \epsilon)$-approximation-stability for the $k$-median objective, then we can efficiently produce a clustering that is $O(\epsilon + \epsilon/\alpha)$-close to $\mathcal{C}_T$.*

PROOF. The algorithm is as given above in Algorithm 3. First, if we are not given the value of $w_{avg}$, we run a constant-factor $k$-median approximation algorithm (e.g., [Arya et al. 2004]) to compute an estimate $\hat{w} \in [w_{avg}, \beta w_{avg}]$ for, say, $\beta = 4$.[1] Redefining $\hat{d}_{crit} = \frac{\alpha \hat{w}}{10\beta \epsilon} \leq \frac{\alpha w_{avg}}{10\epsilon}$, we can show that the set $B = \{x \in S \mid w(x) \geq \hat{d}_{crit} \text{ or } w_2(x) - w(x) < 5\hat{d}_{crit}\}$ of bad points has size $|B| \leq (6 + 10\beta/\alpha)\epsilon n$. (The proof is identical to Proposition 4(ii), but uses Lemma 3(b) with $t = 10\beta$.) If we again define $X_i = C_i^* \setminus B$, we note that Lemma 5 continues to hold with $\tau = 2\hat{d}_{crit}$: the graph $G_\tau$ satisfies properties *(i),(ii)* that all pairs of points in the same $X_i$ are connected by an edge and all pairs of points in different $X_i, X_j$ have no edge and no neighbors in common. In summary, the situation is much as if we knew $w_{avg}$ exactly, except that the number of bad points is slighly greater.

We now show that the greedy method of Step 2 above correctly captures most of the cliques $X_1, X_2, \ldots, X_k$ in $G_\tau$—in particular, we show there is a bijection $\sigma : [k] \to [k]$ such that $\sum_i |X_{\sigma(i)} \setminus C(v_i)| = O(b)$, where $b = |B|$. Since the $b$ bad points 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 5*(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, $C(v_j)$ must contain at least $r_i$ elements from $B$. 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 $C(v_j)$ 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 $C(v_j)$'s must be

---

[1]The reason we need to do this, rather than simply increasing an initial low guess of $w_{avg}$ as in the proof of Theorem 8, is that 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).

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, the size of any such $X_{i'}$ must be bounded by the sizes of its matched $C(v_j)$—else we would have picked a vertex from this clique rather than picking $v_j$. Hence the total size of such $X_{i'}$ is bounded by $|B| + \sum_i r_i \leq 2b$; in turn, this shows that $\sum_j |X_{\sigma(j)} \setminus C(v_j)| \leq \sum_j |X_{\sigma(j)}| \leq 2b$, where this sum is over $j$'s that correspond to the second case. 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 $\epsilon^* n$ to account for the distance between the $k$-median optimum and the target clustering yields the claimed $4b + \epsilon^* n = O(\epsilon + \epsilon/\alpha)n$ result. $\square$

## 4. THE $K$-MEANS PROBLEM

Algorithm 3 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 and Lemma 5 (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 $\epsilon^* = dist(\mathcal{C}_T, \mathcal{C}^*)$.

LEMMA 10. *If the instance $(\mathcal{M}, S)$ satisfies $(1+\alpha, \epsilon)$-approximation-stability for the k-means objective, then*

(a) *If each cluster in $\mathcal{C}_T$ has size at least $2\epsilon n$, then less than $(\epsilon - \epsilon^*)n$ points $x \in S$ on which $\mathcal{C}_T$ and $\mathcal{C}^*$ agree have $w_2(x)^2 - w(x)^2 < \frac{\alpha\text{OPT}}{\epsilon n}$.*

(a') *For the case of general cluster sizes in $\mathcal{C}_T$, less than $6\epsilon n$ points $x \in S$ have $w_2(x)^2 - w(x)^2 < \frac{\alpha\text{OPT}}{2\epsilon n}$.*

*Also, for any $t > 0$ we have:*

(b) *at most $t(\epsilon n/\alpha)$ points $x \in S$ have $w(x)^2 \geq \frac{\alpha\text{OPT}}{t\epsilon n}$.*

The proof is identical to the proof for Lemma 3, and is omitted here. We now give some details for what changes are needed to make Algorithm 2 from Section 3.3 work here. Again, we use a $\beta$-approximation to $k$-means for some constant $\beta$ to get $\widehat{OPT} \in [\text{OPT}, \beta\text{OPT}]$. Define the critical distance $\hat{d}_{crit}$ as $(\frac{\alpha\widehat{OPT}}{25\epsilon\beta n})^{1/2}$ in the case of large clusters, or $(\frac{\alpha\widehat{OPT}}{50\epsilon\beta n})^{1/2}$ in the case of general cluster sizes—these are at most $1/5$ times the square-roots of the expressions in (a) and (a') above. Call point $x \in S$ to be *good* if both $w(x) < d_{crit}$ and $w_2(x) \geq 5d_{crit}$, else *bad* otherwise; let $B$ be the bad points. The following proposition has a proof very similar to Proposition 4(b).

PROPOSITION 11. *If the instance $(\mathcal{M}, S)$ satisfies $(1+\alpha, \epsilon)$-approximation-stability for the k-means objective, then $\mu(B) < (6 + 50\beta/\alpha)\epsilon$.*

Now the rest of the proof for Theorem 9 goes through unchanged in the $k$-means case as well; indeed, first we note that Lemma 5 is true, because it only relies on the

good points being at distance $< d_{crit}$ to their center, and being at distance $\geq 5d_{crit}$ to any other center, and the rest of the proof only relies on the structure of the threshold graph. The fraction of points we err on is again $\epsilon^* + 4\mu(B) = O(\epsilon + \epsilon/\alpha)$. Summarizing, we have the following result.

THEOREM 12 (k-MEANS: GENERAL CASE). *If the given instance $(\mathcal{M}, S)$ satisfies $(1 + \alpha, \epsilon)$-approximation-stability for the k-means objective, we can efficiently produce a clustering which is $O(\epsilon + \epsilon/\alpha)$-close to $\mathcal{C}_T$.*

## 4.1 An Algorithm for Large Clusters

Unfortunately, the argument for exact $\epsilon$-closeness for $k$-median in the case of large target clusters does not extend directly, because Lemma 10(a) is weaker than Lemma 3(a)—the latter gives us bounds on the difference in distances, whereas the former only gives us bounds on the difference in the squared distances. Instead, however, we will use the same algorithm style to identify most of the bad points (by outputting "don't know" on some $O(\epsilon/\alpha)$ of the points) and output a clustering on the remaining $1 - O(\epsilon/\alpha)$ fraction of the points which makes at most $\epsilon n$ errors on these points.

THEOREM 13. *If the given instance $(\mathcal{M}, S)$ satisfies $(1 + \alpha, \epsilon)$-approximation-stability for the k-means objective, and each cluster in $\mathcal{C}_T$ has size at least $(4 + 75/\alpha)\epsilon n + 2$, then we can efficiently produce a clustering in which at most $O(\epsilon n/\alpha)$ points are labeled as "don't know", and on the remainder the clustering is $\epsilon$-close to $\mathcal{C}_T$.*

PROOF. Let us first assume that we know the value of OPT; we will discharge this assumption later. Define the *critical distance* $d_{crit} := \frac{1}{5}(\frac{\alpha \text{OPT}}{\epsilon n})^{1/2}$. As in Theorem 7, we categorize the points in a more nuanced fashion: a point $x \in S$ is called "red" if it satisfies condition (a) of Lemma 10 (i.e., if $w_2(x)^2 - w(x)^2 < 25d_{crit}^2$), "yellow" if it is not red and has $w(x) \in [d_{crit}, 5d_{crit}]$, "orange" if it is not red and has $w(x) > 5d_{crit}$, and "green" otherwise. that the orange and yellow Hence, Lemma 10(a) tells us that there are at most $(\epsilon - \epsilon^*)n$ points on which $\mathcal{C}^*$ and $\mathcal{C}_T$ agree, and that are red; at most $25\epsilon/\alpha$ fraction are either yellow or orange (by setting $t = 25$); at most $\epsilon/\alpha$ fraction of the points are orange (by setting $t = 1$); the rest are green. Let all the non-green points be called bad, and denoted by the set $B$. Let us define $b := (1 + 25/\alpha)\epsilon n$; note that $|B| \leq b$.

Now, as in Theorem 7, if the cluster sizes in the target clustering are at least $2b + \epsilon n + 2$, then constructing the threshold graph $G_\tau$ with $\tau = 2d_{crit}$ and applying Lemma 6 we can find a $k$-clustering $\mathcal{C}'$ where each $X_i := C_i^* \setminus B$ is contained with a distinct cluster, and only the $O(\epsilon + \epsilon/\alpha)$ bad (i.e., non-green) points are possibly in the wrong clusters. We now want to label some points as "don't knows", and construct another clustering $\mathcal{C}''$ where we correctly cluster the green and yellow points.

Again, this is done as in the $k$-median case: for each point $x$ and each cluster $C_j'$, compute the median distance $d_{med}(x, C_j')$ from $x$ to the points in $C_j'$. If the minimum median distance $\min_{j \in [k]} d_{med}(x, C_j')$ is greater than $4d_{crit}$, then

label the point $x$ as "don't know"; else insert $x$ into the cluster $C_i''$ for $i = \arg\min_j d_{med}(x, C_j')$.

First, we claim that the points labeled "don't know" contain all the orange points. Indeed, for any orange point $x$, the distance to each optimal cluster center is at least $5d_{crit}$; moreover, since the target clusters are large, a majority of the points in each cluster $C_j'$ are green, which are all within distance $d_{crit}$ of the optimal cluster center. Using the triangle inequality, the median distance of an orange point to every cluster center will be at least $4d_{crit}$, and hence it will be classified as "don't know". There may be more points classified this, but using a similar argument we can deduce that all such points must have $w(x) \geq 3d_{crit}$, and Lemma 10(b) implies that there are at most $\frac{25\epsilon n}{9\alpha}$ such "don't know" points.

Next, we show that the yellow and green points will be correctly classified: note that each non-red point $x$ satisfies $w_2(x)^2 - w(x)^2 \geq 25d_{crit}^2$, all yellow/green points satisfy $w(x)^2 \leq 25d_{crit}^2$, and all green points $g$ satisfy $w(g) < d_{crit}$, this means that any yellow/green point $x$ must satisfy the property that for a green point $g_1$ in the *same* cluster as $x$ in $\mathcal{C}^*$, and for a green point $g_2$ in a *different* cluster than $x$ in $\mathcal{C}^*$, we have $d(x, g_1) < d(x, g_2)$. Indeed, $d(x, g_1) < w(x) + d_{crit}$ and $d(x, g_2) > w_2(x) + d_{crit}$, and hence it suffices to show that $w_2(x) \geq w(x) + 2d_{crit}$ for $x$ being yellow or green. To show this, note that

$$\begin{aligned}
w_2(x)^2 &\geq w(x)^2 + 25d_{crit}^2 \\
&\geq w(x)^2 + 4d_{crit}^2 + 4 \cdot d_{crit} \cdot (5d_{crit}) \\
&\geq w(x)^2 + 4d_{crit}^2 + 4 \cdot d_{crit} \cdot w(x) \\
&\geq (w(x) + 2d_{crit})^2
\end{aligned}$$

where we use the fact that $w(x) \leq 5d_{crit}$ for green and yellow points. Again, since each yellow or green point is closer to a strict majority of green points in their "correct" cluster in $\mathcal{C}'$, we will correctly classify them. Finally, we finish the argument as before: ignoring the $O(\epsilon n/\alpha)$ "don't knows", $\mathcal{C}''$ may disagree with $\mathcal{C}^*$ on only the $(\epsilon - \epsilon^*)n$ red points where $\mathcal{C}^*$ and $\mathcal{C}_T$ agree, and the $\epsilon^*n$ points where $\mathcal{C}^*$ and $\mathcal{C}_T$ disagree, which is $\epsilon n$ as claimed.

One loose end remains: we assumed we knew OPT and hence $d_{crit}$. To remove this assumption, we can again try multiple guesses for the value of $d_{crit}$ as in Theorem 8. The argument in that theorem continues to hold, as long as the size of the clusters in the target clustering is at least $3b + \epsilon n + 2 = (4 + 75/\alpha)\epsilon n + 2$, which is what we assumed here.  $\square$

## 5.   THE MIN-SUM CLUSTERING PROBLEM

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

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

In this section, we show how assuming $(1 + \alpha, \epsilon)$-approximation-stability for the min-sum clustering problem, and assuming that all the clusters in the target are "large", allows us to find a clustering that is $O(\epsilon)$-close to the target clustering.

## 5.1   The high-level idea

As one might expect, the general plan is to extend the basic techniques from the previous sections, though the situation is now a bit more delicate. While we can still argue that there cannot be too many points that could be cheaply reassigned to different clusters (since that would violate our basic assumption, though we have to be careful about the somewhat messy issue of multiple reassignments), now the cost of reassigning a point $x$ to cluster $C_j$ is proportional to the number of points in $C_j$. In particular, the net effect of this cost structure is that unlike the $k$-median and $k$-means objectives, there is no longer a uniform threshold or critical distance. Many points in some cluster $C_i$ could be quite close to another cluster $C_j$ if $C_j$ is large. On the other hand, one can show the (good) points in $C_j$ will be even closer to each other. Thus, by slowly growing a threshold distance, we will be able to find the clusters in the order from largest to smallest. We then argue that we can identify points in time when the size of the largest component found is large enough compared to the current threshold to have captured the cluster, allowing us to pull those clusters out before they have had the chance to mistakenly connect to smaller ones. This argument will require an assumption that all clusters are large. (See subsequent work of [Balcan and Braverman 2009] for an algorithm that allows for general cluster sizes).

## 5.2   Properties of Min-Sum Clustering

Let the min-sum optimal clustering be $\mathcal{C}^* = \{C_1^*, \ldots, C_k^*\}$ with objective function value $\mathrm{OPT} = \Phi(\mathcal{C}^*)$. For $x \in C_i^*$, define

$$w(x) = \sum_{y \in C_i^*} d(x, y)$$

so that $\mathrm{OPT} = \sum_x w(x)$, and let $w_{avg} = \mathrm{avg}_x w(x) = \frac{\mathrm{OPT}}{n}$. Define

$$w_2(x) = \min_{j \neq i} \sum_{y \in C_j^*} d(x, y).$$

A useful fact, following immediately from the triangle inequality, is the following:

FACT 14. *For two points $x$ and $y$, and any cluster $C_j^*$,*

$$\sum_{z \in C_j^*} (d(x, z) + d(y, z)) \geq |C_j^*| \, d(x, y).$$

We now prove the following lemma.

LEMMA 15. *If the given instance $(\mathcal{M}, S)$ satisfies $(1+\alpha, \epsilon)$-approximation-stability for the min-sum objective 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 $w_2(x) < \frac{\alpha w_{avg}}{4\epsilon}$, and*

(b) *at most $60\epsilon n/\alpha$ points $x \in S$ have $w(x) > \frac{\alpha w_{avg}}{60\epsilon}$.*

PROOF. To prove Property (a), assume to the contrary. Then one could take $\mathcal{C}^*$ and move a set $S'$ of $(\epsilon - \epsilon^*)n$ points $x$ that have $w_2(x) < \frac{\alpha w_{avg}}{4\epsilon}$ and on which $\mathcal{C}_T$ and $\mathcal{C}^*$ agree to the clusters that define their $w_2$ value. We now argue that the resulting increase in min-sum objective value is less than $\alpha$OPT.

Let the new clustering be $\mathcal{C}' = (C_1', \ldots, C_k')$, where $|C_i' \setminus C_i^*| = \delta_i n$, so that $\sum_i \delta_i = \epsilon - \epsilon^*$. Also, let $C_2(x)$ denote the cluster $C_i^*$ that point $x \in S'$ is moved to. Then, for each point $x \in S'$ moved, the increase to the min-sum objective is at most $2w_2(x) + \sum_{y \in S':C_2(y)=C_2(x)} d(x, y)$—here the factor of two arises because the min-sum objective counts each pair of points in a cluster twice, once from each end. From Fact 14, we know that if $C_2(y) = C_2(x)$ then $d(x, y) \leq \frac{1}{|C_2(x)|}(w_2(x) + w_2(y))$. Thus, we can replace the term $d(x, y)$ in the cost charged to point $x$ with $\frac{2}{|C_2(x)|}w_2(x)$, yielding a total cost charged to point $x$ moved to cluster $C_i^*$ of

$$2w_2(x) + 2w_2(x)\frac{\delta_i n}{|C_i^*|}.$$

Summing over all points $x$ moved to all clusters, and using the fact that $w_2(x) < \frac{\alpha w_{avg}}{4\epsilon}$ for all $x \in S'$, we have a total cost increase of less than

$$\sum_i (\delta_i n)\frac{2\alpha w_{avg}}{4\epsilon}\left[1 + \frac{\delta_i n}{|C_i^*|}\right] \leq \epsilon n \frac{\alpha w_{avg}}{2\epsilon} + \frac{\alpha w_{avg}}{2\epsilon}\sum_i \frac{\delta_i^2 n^2}{|C_i^*|}$$

$$\leq \frac{\alpha}{2}\text{OPT} + \frac{\alpha w_{avg}}{2\epsilon}\frac{\epsilon^2 n^2}{\min_i |C_i^*|}$$

$$\leq \frac{\alpha}{2}\text{OPT} + \frac{\alpha}{4}\text{OPT} < \alpha\text{OPT}.$$

Finally, property (b) follows immediately from Markov's inequality. □

Let us define the *critical value* $v_{crit} := \frac{\alpha w_{avg}}{60\epsilon}$. We call point $x$ *good* if it satisifies both $w(x) \leq v_{crit}$ and $w_2(x) \geq 15v_{crit}$, 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 16 (STRUCTURE OF MIN-SUM OPTIMUM). *If the given instance $(\mathcal{M}, S)$ satisfies $(1 + \alpha, \epsilon)$-approximation-stability for the min-sum objective then as long as the minimum cluster size is at least $2\epsilon n$ we have:*

(i) *For all $x, y$ in the same $X_i$, we have $d(x, y) < \frac{2v_{crit}}{|C_i^*|}$,*

(ii) *For $x \in X_i$ and $y \in X_{j \neq i}$, we have $d(x, y) > \frac{14 v_{crit}}{\min(|C_i^*|, |C_j^*|)}$, and*

(iii) *The number of bad points $|B| = |S \setminus \cup X_i|$ is at most $b := (1 + 60/\alpha)\epsilon n$.*

PROOF. For part (i), note that Fact 14 implies that

$$d(x, y) \leq \frac{1}{|C_i^*|}\sum_{z \in C_i^*}(d(x, z) + d(y, z)) = \frac{1}{|C_i^*|}(w(x) + w(y)).$$

Since $x, y \in X_i$ are both good, we have $w(x), w(y) \leq v_{crit}$, so part (i) follows.

For part (ii), assume without loss of generality that $|C_j^*| \leq |C_i^*|$. Since both $x \in C_i^*, y \in C_j^*$ are good, we have $w_2(x) = \sum_{z \in C_j^*} d(x, z) \geq 15v_{crit}$ and $w(x) =$

$\sum_{z \in C_j^*} d(y, z) \leq v_{crit}$. By the triangle inequality $d(x, y) \geq d(x, z) - d(y, z)$, we have

$$|C_j^*| \, d(x, y) \geq \sum_{x \in C_j^*} (d(x, z) - d(y, z)) = w_2(x) - w(y) \geq 14 v_{crit}.$$

Finally, part (iii) follows from Lemma 15 and a trivial union bound.  □

While Lemma 16 is similar in spirit to Lemma 5, there is a crucial difference: the distance between the good points in $X_i$ and those in $X_j$ is no longer bounded below by some absolute value $\tau$, but rather the bound depends on the sizes of $X_i$ and $X_j$. However, a redeeming feature is that the separation is large compared to the sizes of *both* $X_i$ and $X_j$; we will use this feature crucially in our algorithm.

### 5.3  The Algorithm for Min-Sum Clustering

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

THEOREM 17. *If the given instance $(\mathcal{M}, S)$ satisfies $(1 + \alpha, \epsilon)$-approximation-stability for the min-sum objective and we are given the value of $w_{avg}$, then so long as the smallest correct cluster has size greater than $(5 + 120/\alpha)\epsilon n$, Algorithm 4 produces a clustering that is $O(\epsilon/\alpha)$-close to the target. If we are not given $w_{avg}$, then we can use Algorithm 4 as a subroutine to produce a clustering that is $O(\epsilon/\alpha)$-close to the target.*

---

**Algorithm 4** Min-sum Algorithm

**Input:** $(\mathcal{M}, S)$, $w_{avg}$, $\epsilon \leq 1$, $\alpha > 0$, $k$, $b := (1 + 60/\alpha)\epsilon n$.
Let the initial threshold $\tau = \tau_0$.

**Step 1:** If $k = 0$ or $S = \emptyset$, stop.
**Step 2:** Construct the $\tau$-threshold graph $G_\tau$ on the current set $S$ of points.
**Step 3:** Create a new graph $H$ by connecting two points in $S$ by an edge if they share at least $b$ neighbors in common in $G_\tau$.
**Step 4:** Let $C$ be largest connected component in $H$. If $|C| \geq 3 v_{crit}/\tau$,
   **then** output $C$ as a cluster, set $k \leftarrow k - 1$, $S \leftarrow S \setminus C$, and go to Step 1,
   **else** increase $\tau$ to the next critical threshold and go to Step 1.

---

PROOF. Since each cluster in the target clustering has more than $(5 + 120/\alpha)\epsilon n = 2b + 3\epsilon n$ points by the assumption, and the *optimal min-sum clustering* $\mathcal{C}^*$ must differ from the target clustering by fewer than $\epsilon n$ points, and hence each cluster in $\mathcal{C}^*$ must have more than $2b + 2\epsilon n$ points. Moreover, by Lemma 15(iii), the bad points $B$ constitute at most $b$ points, and hence each $|X_i| = |C_i^* \setminus B| > b + 2\epsilon n \geq b + 2$.

**Analysis under the assumption that $w_{avg}$ is given:** Consider what happens in the execution of the algorithm: as we increase $\tau$, the sizes of the $H$-components increase (since we are adding more edges in $G_\tau$). This happens until the largest $H$-component is "large enough" (i.e., the condition in Step 4 gets satisfied) and we output a component whose size is large enough; and then we go back to raising $\tau$.

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 $\tau$, 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 $|C_1^*| \geq |C_2^*| \geq \ldots \geq |C_k^*|$, and let $n_i = |C_i^*|$. Define $d_i = \frac{2v_{crit}}{n_i}$ and recall that $\max_{x,y \in X_i} d(x,y) \leq d_i$ by Lemma 16(i).

We first claim that as long as $\tau \leq 3\,d_1$, no two points belonging to different $X_i$'s can lie in the same $H$-component. By Lemma 16(ii) the distance between points in any $X_i$ and $X_{j \neq i}$ is strictly greater than $\frac{14v_{crit}}{\min(n_i,n_j)}$, which is strictly greater than $2\tau$ for any $\tau \leq 3\,d_1$. Hence every $x \in X_i$ and $y \in X_j$ share no common neighbors, and by an argument identical to that in Lemma 6, the nodes $x, y$ belong to different components of $H$.

Next, we claim that for values of $\tau < \min\{d_i, 3d_1\}$, the $H$-component containing points from $X_i$ 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 argument in the previous paragraph, and hence any $H$-component containing points from $X_i$ has size at most $|C_i^*| + |B| < \frac{3n_i}{2}$; here we used the fact that each $n_i > 2b$ due to the large cluster assumption. However, the minimum size bound $\frac{3v_{crit}}{\tau}$ in Step 4 is equal to $\frac{3\,d_i,n_i}{2\tau}$ (by the definition of $d_i$, which is at most $\frac{3n_i}{2}$ for values of $\tau < d_i$) hence the condition of Step 4 is not satisfied and the $H$-component will not be output. Moreover, note that when $\tau \geq d_i$, all the points of $X_i$ lie in the same $H$-component.

The above two paragraphs show that nothing bad happens: no incorrect components are constructed or components outputted prematurely. We finally show that something good happens—in particular, 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 minimum-size condition (which at time $\tau = 3d_1$ requires a cluster of size $\frac{3v_{crit}}{\tau} = \frac{v_{crit}}{d_1} = n_1/2$), giving the contradiction.

To recap, we showed that 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 now iterate this argument on remaining set of points to complete the proof for the case when we know $w_{avg}$.

**Analysis if $w_{avg}$ is not given:** In this case, we do not want to use a $\beta$-approximation algorithm for min-sum to obtain a clustering that is $O(\beta\epsilon/\alpha)$-close to the target, because the minsum clustering problem only has 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 $w_{avg}$: we then stop the first time we output $k$ clusters that cover at least $n - b = (1 - O(\epsilon/\alpha))n$ points in $S$. Clearly, if we reached the correct value of $w_{avg}$ 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 clustering all but an $O(\epsilon/\alpha)$-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), such an event would use up an additional cluster and therefore 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 $w_{avg}$. The dangerous case is case (c), but we claim case (c) in fact cannot happen. Indeed, the value of $\tau$ 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 $w_{avg}$ is too small, this $H$-component would certainly satisfy the condition of Step 4 and be output as a cluster instead. ☐

## 6. RELATIONSHIP TO $\epsilon$-SEPARATION CONDITION

Ostrovsky et al. [2006] consider $k$-means clustering in Euclidean spaces, and define and analyze an interesting separation condition that provides a notion of how "naturally clustered" a given dataset is. Specifically, they call a $k$-means instance $\epsilon$-*separated* if the optimal $k$-means cost is at most $\epsilon^2$ times the cost of the optimal $(k-1)$-means solution. Under this assumption on the input, they show how to seed Lloyd's method to obtain a $1+f(\epsilon)$ approximation in $d$-dimensional Euclidean space in time $O(nkd + k^3d)$, and a $(1+\delta)$-PTAS with run-time $nd2^{k(1+\epsilon^2)/\delta}$. This notion of $\epsilon$-separation, namely that any $(k-1)$-means solution is substantially more expensive than the optimal $k$-means solution, is in fact related to $(c, \epsilon)$-approximation-stability. Indeed, in Theorem 5.1 of their paper, they show that their $\epsilon$-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 necessarily the case: an instance could satisfy our condition without being $\epsilon$-separated.[2] We present here a specific example. Consider $k = 2$ where target cluster $C_1$ has $(1-\alpha)n$ points and target cluster $C_2$ has $\alpha n$ points. Suppose that any two points inside the same cluster $C_i$ have distance 1 and any two points inside different clusters have distance $1 + 1/\epsilon$. For any $\alpha \in (\epsilon, 1-\epsilon)$, this satisfies $(2, \epsilon)$-approximation-stability for $k$-median (and is $(2, \epsilon^2)$-approximation-stable for $k$-means for any $\alpha \in (\epsilon^2, 1-\epsilon^2)$). However, it need not satisfy the $\epsilon$-separation property: for $\alpha = 2\epsilon$, the optimal 2-median solution has cost $n - 2$, but the optimal 1-median solution has cost less than $3n$. Likewise for $\alpha = 2\epsilon^2$, the optimal 2-means solution has cost $n - 2$, but the optimal 1-means solutions has cost less than $(3 + 4\epsilon)n$. Thus, in both cases the ratio of costs between $k = 1$ and $k = 2$ is not so large.

---

[2][Ostrovsky et al. 2006] shows an implication in this direction (Theorem 5.2); however, this implication requires a substantially stronger condition, namely that data satisfy $(c, \epsilon)$-approximation-stability for $c = 1/\epsilon^2$ (and that target clusters be large). In contrast, our primary interest is in the case where $c$ is below the threshold for existence of worst-case approximation algorithms.

In fact, for the case that $k$ is much larger than $1/\epsilon$, the difference between the two properties can be more dramatic. Suppose $\epsilon$ is a small constant, and consider a clustering instance in which the target consists of $k = \sqrt{n}$ clusters with $\sqrt{n}$ points each, such that all points in the same cluster have distance 1 and all points in different clusters have distance $D+1$ where $D$ is a large constant. Then, merging two clusters increases cost by $O(\sqrt{n})$, so the optimal $(k-1)$-means/median solution is just a factor $1+O(1/\sqrt{n})$ more expensive than the optimal $k$-means/median clustering. However, for $D$ sufficiently large compared to $1/\epsilon$, this satisfies $(2, \epsilon)$-approximation-stability or even $(1/\epsilon, \epsilon)$-approximation-stability (for proof, see appendix B).

## 7.   SUBSEQUENT WORK

Our work has already inspired a number of both theoretical and practical exciting subsequent developments. We discuss some of these here.

**Further guarantees for min-sum clustering.** In subsequent work, Balcan and Braverman [2009] have further analyzed the min-sum problem and showed how to handle the presence of small target clusters. To achieve this they derive new structural properties implied by $(c, \epsilon)$-approximation-stability. In the case where $k$ is small compared to $\log n / \log \log n$ they output a single clustering which is $O(\epsilon/\alpha)$-close to the target, while in the general case their algorithm outputs a small list of clusterings with the property that the target clustering is close to one of those in the list. Balcan and Braverman [2009] further show that if we do require the clusters to be large (of size at least size greater than $100\epsilon n/\alpha^2$), they can reduce the approximation error from $O(\epsilon/\alpha)$ down to $O(\epsilon)$ – the best one can hope for.

**Further guarantees for $k$-median and $k$-means clustering.** Schalekamp et al. [2010] show that Algorithm 1 for the case of the $k$-median problem when clusters are large, additionally achieves a good approximation to the $k$-median objective. Note that the approximation hardness result shown in Theorem 19 for clustering under $(c, \epsilon)$-approximation-stability requires the target to have small clusters. They also discuss implementation issues and perform a number of experimental comparisons between various algorithms. Awasthi et al. [2010b] go further and provide a PTAS for $k$-median, as well as for $k$-means in Euclidean space, when the target has large clusters. Specifically, they achieve these results under a condition they call "weak deletion stability" which includes $(1 + \alpha, \epsilon)$-approximation-stability in the case that all target clusters have size greater than $\epsilon n$ and $\alpha > 0$ is constant. (The condition also generalizes that of [Ostrovsky et al. 2006], requiring only that the optimum for $k-1$ clusters be a factor $1+\alpha$ more expensive than the optimum for $k$ clusters.) One implication of this is that when $\alpha > 0$ is a constant, they improve the "largeness" condition needed to efficiently get $\epsilon$-close for $k$-median from $O((1 + 1/\alpha)\epsilon n)$ down to $\epsilon n$. Another implication is that they are able to get the same guarantee for $k$-means as well, when points lie in $R^n$, improving on the guarantees in Section 4.1 for points in Euclidean space. Note that while $\alpha$ does not appear in the "largeness" condition, their algorithm has running time that depends exponentially on $1/\alpha$, whereas ours does not depend on $1/\alpha$ at all.

**Correlation clustering.** Balcan and Braverman [2009] also analyze the correlation clustering problem under the $(c, \epsilon)$-approximation-stability assumption. For

correlation clustering, the input is a graph with edges labeled $+1$ or $-1$ and the goal is to find a partition of the nodes that best matches the signs of the edges [Blum et al. 2004]. Usually, two versions of this problem are considered: minimizing disagreements and maximizing agreements. In the former case, the goal is to minimize the number of $-1$ edges inside clusters plus the number of $+1$ edges between clusters, while in the latter case the goal is to maximize the number of $+1$ edges inside the cluster plus the number of $-1$ edges between. These are equivalent at optimality but differ in their difficulty of approximation. [Balcan and Braverman 2009] show that for the objective of minimizing disagreements, $(1 + \alpha, \epsilon)$-approximation-stability implies $(2.5, O(\epsilon/\alpha))$-approximation-stability, so one can use a state-of-the-art 2.5-approximation algorithm for minimizing disagreements in order to get an accurate clustering.[3] This contrasts sharply with the case of objectives such as $k$-median, $k$-means and min-sum (see Theorem 18).

**Stability with noise and outliers.** Balcan, Roeglin, and Teng [Balcan et al. 2009] consider a relaxation of $(c, \epsilon)$-approximation-stability that captures clustering in the presence of outliers. In real data there may well be some data points for which the (heuristic) distance measure does not reflect cluster membership well, causing $(c, \epsilon)$-approximation-stability to be violated. To capture such situations, they define $(\nu, c, \epsilon)$-approximation-stability which requires that the data satisfies $(c, \epsilon)$-approximation-stability only after a $\nu$ fraction of the data points have been removed. Balcan et al. [2009] show that in the case where the target clusters are large (have size $\Omega((\epsilon/\alpha + \nu)n)$ the algorithm we present in this paper for the large clusters case can be used to output a clustering that is $(\nu + \epsilon)$-close to the target clustering. They also show that in the more general case there can be multiple significantly different clusterings that can satisfy $(\nu, c, \epsilon)$-approximation-stability (since two different sets of outliers could result in two different clusterings satisfying the condition). However, if *most* of the points come from large clusters, they show one can in polynomial time output a small list of $k$-clusterings such that any clustering that satisfies the property is close to one of the clusterings in the list.

**The Inductive model.** Balcan et al. [2009] and Balcan and Braverman [2009] also show how to cluster well under approximation-stability in the *inductive* clustering setting. Here, $S$ is a small random subset of points from a much larger set $X$, and our goal is to produce a hypothesis $h : X \to Y$ which implicitly represents a clustering of the whole set $X$ and which has low error on $X$. Balcan et al. [2009] show how in the large clusters case the analysis in our paper can be adapted to the inductive model for $k$-median and $k$-means, and Balcan and Braverman [2009] have shown how to adapt their minsum algorithm to the inductive setting as well.

**Clustering with one-versus all queries.** Motivated by clustering applications in computational biology, Voevodski et al. [2010] analyze $(c, \epsilon)$-approximation-stability in a model with unknown distance information where one can only make a limited number of *one versus all* queries. Voevodski et al. [2010] design an algorithm that given $(c, \epsilon)$-approximation-stability for the $k$-median objective finds a clustering

---

[3]Note that the maximizing agreement version of correlation clustering is less interesting in our framework since it admits a PTAS.

that is very close to the target by using only $O(k)$ one-versus-all queries in the large cluster case, and in addition is faster than the algorithm we present here. In particular, the algorithm for the large clusters case we describe in Section 3 can be implemented in $O(|S|^3)$ time, while the one proposed in [Voevodski et al. 2010] runs in time $O(|S|k(k + \log |S|))$. Voevodski et al. [2010] use their algorithm to cluster biological datasets in the Pfam [Finn et al. 2010] and SCOP [Murzin et al. 1995] databases, where the points are proteins and distances are inversely proportional to their sequence similarity. This setting nicely fits the one-versus all queries model because one can use a fast sequence database search program to query a sequence against an entire dataset. The Pfam [Finn et al. 2010] and SCOP [Murzin et al. 1995] databases are used in biology to observe evolutionary relationships between proteins and to find close relatives of particular proteins. Voevodski et al. [2010] find that for one of these sources they can obtain clusterings that almost exactly match the given classification, and for the other the performance of their algorithm is comparable to that of the best known algorithms using the full distance matrix.

## 8. CONCLUSIONS AND OPEN QUESTIONS

In this work, we provide a new approach to the problem of clustering, motivated by the fact that common clustering objective functions such as $k$-median, $k$-means, and min-sum objectives are often only a proxy for the true goals of partitioning data in a manner that matches some correct or desired solution. From this view, an implicit assumption of the approximation-algorithms approach to clustering is that achieving a good approximation to the objective on the given instance will indeed result in a solution that has low error with respect to the desired answer. In this work, we make this implicit assumption explicit, formalizing it as $(1 + \alpha, \epsilon)$-approximation-stability for the given objective $\Phi$. We then prove that this property, for the case of $k$-median, $k$-means, and min-sum objectives, allows one to design efficient algorithms that provide solutions of error $O(\epsilon)$ for any constant $\alpha > 0$, even values $\alpha$ such that obtaining a $1 + \alpha$-approximation for these objectives is NP-hard. Thus, we are able to approximate the *solution* even if in the worst case it is hard to approximate the objective.

We note that our results can also be viewed in the context of algorithms for clustering under natural stability conditions on the data. In particular, even without reference to the target clustering, we can consider a clustering instance to be stable if all $1 + \alpha$ approximations $\mathcal{C}$ to objective $\Phi$ satisfy $dist(\mathcal{C}, \mathcal{C}^*) < \epsilon$, where $\mathcal{C}^*$ is the $\Phi$-optimal clustering. In this case, our results immediately imply an algorithm to get $O(\epsilon)$-close to $\mathcal{C}^*$. Thus, we still find a clustering that is "nearly as good as $\mathcal{C}^*$" while being agnostic to (i.e., without making assumptions on) the quality of $\mathcal{C}^*$ in terms of its distance to the target.

As noted in Section 7, subsequent work has already demonstrated the practicality of our approach for real world clustering problems.

**Open questions:** One natural open question is whether the $O(\epsilon/\alpha)$ form of the bounds we achieve are intrinsic, or if improved bounds for these objectives are possible. For example, suppose our instance satisfies $(1 + \epsilon, \epsilon)$-approximation-stability for *all* $\epsilon > 0$, say for $k$-median (e.g., achieving a 1.01-approximation would produce a solution of error 1%); is this sufficient to produce a near-optimal solution

of some form? Another natural question is whether one can use this approach for other clustering or partitioning objective functions. For example, the *sparsest cut* problem has been the subject of a substantial body of research, with the best known approximation guarantee a factor of $O(\sqrt{\log n})$ [Arora et al. 2004]. However, in the event this objective is a proxy for a true goal of partitioning a dataset in a nearly-correct manner, it is again natural to consider data satisfying $(c, \epsilon)$-approximation-stability. In this case, given the current state of approximation results, it would be of interest even if $c$ is a large constant. See [Balcan 2009] for more details.

More broadly, there are other types of problems (for example, evolutionary tree reconstruction) where the measurable objectives typically examined may again only be a proxy for the true goals (e.g., to produce a correct evolutionary tree). It would be interesting to examine whether the approach developed here might be of use in those settings as well.

## REFERENCES

ACHLIOPTAS, D. AND MCSHERRY, F. 2005. On spectral learning of mixtures of distributions. In *Proceedings of the Eighteenth Annual Conference on Learning Theory.*

AGARWAL, P. AND MUSTAFA, N. 2004. *k*-means projective clustering. In *Proceedings of the 23rd Annual Symposium on Principles of Database Systems.*

ALON, N., DAR, S., PARNAS, M., AND RON, D. 2000. Testing of clustering. In *Proceedings of the 41st Annual Symposium on Foundations of Computer Science.*

ARORA, S. AND KANNAN, R. 2005. Learning mixtures of arbitrary gaussians. In *Proceedings of the 37th ACM Symposium on Theory of Computing.*

ARORA, S., RAGHAVAN, P., AND RAO, S. 1999. Approximation schemes for Euclidean *k*-medians and related problems. In *Proceedings of the Thirty-First Annual ACM Symposium on Theory of Computing.*

ARORA, S., RAO, S., AND VAZIRANI, U. 2004. Expander flows, geometric embeddings, and graph partitioning. In *Proceedings of the 36th Annual ACM Symposium on Theory of Computing.*

ARTHUR, D. AND VASSILVITSKII, S. 2006. Worst-case and smoothed analyses of the icp algorithm, with an application to the k-means method. In *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science.*

ARYA, V., GARG, N., KHANDEKAR, R., MEYERSON, A., MUNAGALA, K., AND PANDIT., V. 2004. Local search heuristics for k-median and facility location problems. *SIAM J. Comput. 33,* 3, 544–562.

AWASTHI, P., BLUM, A., AND SHEFFET, O. 2010a. Center-based clustering under perturbation stability. Manuscript.

AWASTHI, P., BLUM, A., AND SHEFFET, O. 2010b. Stability yields a PTAS for *k*-median and *k*-means clustering. In *Proceedings of the 51st Annual IEEE Symposium on Foundations of Computer Science.*

BALCAN, M. 2009. Better guarantees for sparsest cut clustering. In *Proceedings of the 22nd Annual Conference on Learning Theory.*

BALCAN, M., BLUM, A., AND VEMPALA, S. 2008. A discrimantive framework for clustering via similarity functions. In *Proceedings of the 40th ACM Symposium on Theory of Computing.*

BALCAN, M. AND BRAVERMAN, M. 2009. Finding low error clusterings. In *Proceedings of the 22nd Annual Conference on Learning Theory.*

BALCAN, M., ROEGLIN, H., AND TENG, S. 2009. Agnostic clustering. In *Proceedings of the 20th International Conference on Algorithmic Learning Theory.*

BARTAL, Y., CHARIKAR, M., AND RAZ, D. 2001. Approximating min-sum k-clustering in metric spaces. In *Proceedings on 33rd Annual ACM Symposium on Theory of Computing.*

BILU, Y. AND LINIAL, N. 2010. Are stable instances easy? In *Proceedings of the First Symposium on Innovations in Computer Science.*

BLUM, A., BANSAL, N., AND CHAWLA, S. 2004. Correlation clustering. *Machine Learning 56*, 89–113.

CHARIKAR, M. AND GUHA, S. 1999. Improved combinatorial algorithms for the facility location and k-median problems. In *Proceedings of the 40th Annual Symposium on Foundations of Computer Science*.

CHARIKAR, M., GUHA, S., TARDOS, E., AND SHMOY, D. B. 1999. A constant-factor approximation algorithm for the k-median problem. In *Proceedings of the Thirty-First Annual ACM Symposium on Theory of Computing*.

CZUMAJ, A. AND SOHLER, C. 2007. Small space representations for metric min-sum k-clustering and their applications. In *Proceedings of the 24th International Symposium on Theoretical Aspects of Computer Science*.

DASGUPTA, S. 1999. Learning mixtures of gaussians. In *Proceedings of The 40th Annual Symposium on Foundations of Computer Science*.

DE LA VEGA, W. F., KARPINSKI, M., KENYON, C., AND RABANI, Y. 2003. Approximation schemes for clustering problems. In *Proceedins of the Thirty-Fifth Annual ACM Symposium on Theory of Computing*.

DEVROYE, L., GYORFI, L., AND LUGOSI, G. 1996. *A Probabilistic Theory of Pattern Recognition*. Springer-Verlag.

DUDA, R. O., HART, P. E., AND STORK, D. G. 2001. *Pattern Classification*. Wiley.

FEIGE, U. 1998. A threshold of $\ln n$ for approximating set cover. *J. ACM 45*, 4, 634–652.

FINN, R., MISTRY, J., TATE, J., COGGILL, P., ANDJ.E. POLLINGTON, A. H., GAVIN, O., GUNE-SEKARAN, P., CERIC, G., FORSLUND, K., HOLM, L., SONNHAMMER, E., EDDY, S., AND BATEMAN, A. 2010. The pfam protein families database. *Nucleic Acids Research 38*, D211–222.

GOLLAPUDI, S., KUMAR, R., AND SIVAKUMAR, D. 2006. Programmable clustering. In *Proceedings of the 25th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*.

GUHA, S. AND KHULLER, S. 1999. Greedy strikes back: Improved algorithms for facility location. *Journal of Algorithms 31*, 1, 228–248.

INDYK, P. 1999. Sublinear time algorithms for metric space problems. In *Proceedings of the Thirty-First Annual ACM Symposium on Theory of Computing*.

JAIN, K., MAHDIAN, M., AND SABERI, A. 2002. A new greedy approach for facility location problems. In *Proceedings of the 34th Annual ACM Symposium on Theory of Computing*.

JAIN, K. AND VAZIRANI, V. V. 2001. Approximation algorithms for metric facility location and k-median problems using the primal-dual schema and lagrangian relaxation. *JACM 48*, 2, 274 – 296.

KANNAN, R., SALMASIAN, H., AND VEMPALA, S. 2005. The spectral method for general mixture models. In *Proceedings of The Eighteenth Annual Conference on Learning Theory*.

KLEINBERG, J. 2002. An impossibility theorem for clustering. In *Proceedings of the Neural Information Processing Systems*.

KUMAR, A. AND KANNAN, R. 2010. Clustering with spectral norm and the $k$-means algorithm. In *Proceedings of the 51st Annual IEEE Symposium on Foundations of Computer Science*.

KUMAR, A., SABHARWAL, Y., AND SEN, S. 2004. 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*.

LLOYD, S. 1982. Least squares quantization in PCM. *IEEE Trans. Inform. Theory 28*, 2, 129–137.

MEILA, M. 2003. Comparing clusterings by the variation of information. In *Proceedings of the The Sixteenth Annual Conference on Learning Theory*.

MEILA, M. 2005. Comparing clusterings – an axiomatic view. In *International Conference on Machine Learning*.

MEILA, M. 2006. The uniqueness of a good clustering for k-means. In *Proceedings of the 23rd International Conference on Machine Learning*.

MURZIN, A., BRENNER, S. E., HUBBARD, T., AND CHOTHIA, C. 1995. Scop: a structural classification of proteins database for the investigation of sequences and structures. *Journal of Molecular Biology 247*, 536–540.

OSTROVSKY, R., RABANI, Y., SCHULMAN, L., AND SWAMY, C. 2006. The effectiveness of lloyd-type methods for the k-means problem. In *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science*.

SCHALEKAMP, F., YU, M., AND A, V. Z. 2010. Clustering with or without the approximation. In *Proceedings of the 16th Annual International Computing and Combinatorics Conference*.

SCHULMAN, L. 2000. Clustering for edge-cost minimization. In *Proceedings of the Thirty-Second Annual ACM Symposium on Theory of Computing*. 547–555.

VEMPALA, S. AND WANG, G. 2004. A spectral algorithm for learning mixture models. *JCSS 68*, 2, 841–860.

VOEVODSKI, K., BALCAN, M., ROEGLIN, H., TENG, S., AND XIA, Y. 2010. Efficient clustering with limited distance information. In *Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence*.

## A. ADDITIONAL PROOFS

THEOREM 18. *For any $1 \le c_1 < c_2$, any $\epsilon, \delta > 0$, for sufficiently large $k$, there exists a family of metric spaces $G$ and target clusterings that satisfy $(c_1, \epsilon)$-approximation-stability for the $k$-median objective (likewise, $k$-means and min-sum) and yet do not satisfy even $(c_2, 1/2 - \delta)$-approximation-stability 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 $n(1 - 2\delta)$ points and $k - 1$ clusters $C_2, \dots, C_k$ each with $\frac{2\delta n}{k-1}$ points. All points in the same cluster have distance 1. The distance between points in any two distinct clusters $C_i, C_j$ for $i, j \ge 2$ 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 $n - k$. 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 $1/2 - \delta$, and furthermore this solution has a cost of $(\frac{2\delta n}{k-1})(D - 1) + n - k$, and we define $D$ to set this equal to $c_2(n - k)) = c_2 \text{OPT}$.

Any $c_1$ approximation, however, must be $\epsilon$-close to the target for $k > 1 + 2\delta/\epsilon$. In particular, by definition of $D$, any $c_1$-approximation must have one median inside each $C_i$. Therefore, it cannot place two medians inside $C_1$ as in the above $c_2$-approximation, and so can have error on fewer than $\frac{2\delta n}{k-1}$ points. This is less than $\epsilon n$ by definition of $k$.

The same construction, with $D$ defined appropriately, applies to $k$-means as well. In particular, we just define $D$ to be the square-root of the value used for $D$ above, and the entire argument proceeds as before.

For min-sum, we modify the construction so that distances in $C_1$ are all equal to 0, so now $\text{OPT} = (k - 1)(\frac{2\delta n}{k-1})(\frac{2\delta n}{k-1} - 1)$. Furthermore, we set points in $C_1$ to be at distance greater than $c_2 \text{OPT}$ from all other points. We again define $D$ so that the cheapest way to use $k - 2$ clusters for the points in $C_2 \cup \dots \cup C_k$ has cost exactly $c_2 \text{OPT}$. However, because of the pairwise nature of the min-sum objective, this is now to equally distribute the points in one of the clusters $C_2, \dots, C_k$ among all the others. This has cost $2(\frac{2\delta n}{k-1})^2 D + \text{OPT} - (\frac{2\delta n}{k-1})^2(\frac{k-3}{k-2})$, which as mentioned above we set to $c_2 \text{OPT}$. Again, because we have defined $D$ such that the cheapest clustering of $C_2 \cup \dots \cup C_k$ using $k - 2$ clusters has cost $c_2 \text{OPT}$, any $c_1$ approximation

must use $k-1$ clusters for these points and therefore again must have error less than $\frac{2\delta n}{k-1} < \epsilon n$.   □

THEOREM 19. *For k-median, k-means, and min-sum objectives, for any $c > 1$, the problem of finding a c-approximation can be reduced to the problem of finding a c-approximation under $(c, \epsilon)$-approximation-stability. Therefore, the problem of finding a c-approximation under $(c, \epsilon)$-approximation-stability 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 that is $(c, \epsilon)$-approximation-stable. In particular we create a new graph $G'$ by adding an extra $n/\epsilon$ nodes that are all at distance $D$ from each other and from the nodes in $G$, where $D$ is chosen to be larger than $c$OPT on $G$ (e.g., $D$ could be the sum of all pairwise distances in $G$). 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 $(c, \epsilon)$-approximation-stability. This is because, by definition of $D$, any solution that does not put each of the new nodes into its own singleton cluster will incur too high a cost to be a $c$-approximation. So a $c$-approximation can only differ from the target on $G$ (which has less than an $\epsilon$ 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 in any of the $k$-median, $k$-means, or min-sum objectives.   □

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

THEOREM 20. *For any contant $c'$, for any $\epsilon < 1/(ec')$, it is NP-hard to find a clustering of error at most $c'\epsilon$ for the k-median and k-means problem under $(1, \epsilon)$-approximation-stability.*

PROOF. First, let us prove a $(1 + 1/e)$-hardness for instances of $k$-median where one is allowed to place centers at any point in the metric space. The proof is very similar to the proof from [Guha and Khuller 1999; Jain et al. 2002] which gives a $(1 + 2/e)$-hardness for the case where one can place centers only at a distinguished set of locations in the metric space. We then show how to alter this hardness result to prove the theorem.

Consider the max-$k$-coverage problem with $n$ elements and $m$ sets: that is, given $m$ subsets of a universe of $n$ elements, find $k$ sets whose union covers as many elements as possible. It is NP-hard to distinguish between instances of this problem where there exist $k$ sets that can cover all the elements, and instances where any $k$ sets cover only $(1 - 1/e)$-fraction of the elements [Feige 1998]. The hard instances have the property that both $m$ and $k$ are a tiny fraction of the number of elements $n$. For some suitably large constant $C$, we construct an instance of $k$-median with $cn+m$ points, one point for each set and $c$ points for each element, assign distance 1 between any two points such that one of them represents an element and the other a set containing that point, and distance 2 to all other pairs.

Note that if there are $k$ sets in the set system that cover all the elements (the "yes" instances), choosing the corresponding $k$ points as centers gives us a solution

of cost $cn + 2(m - k) \leq (1 + \delta)cn$ for some arbitrarily small constant $\delta > 0$. On the other hand, given any solution to $k$-median with cost $C$, if any of these centers are on points corresponding to elements, we can choose a set-point at unit distance from it instead, thus potentially increasing the cost of the solution by at most $m$ to $C + m$. Hence, if this collection of $k$ sets covers at most $(1 - 1/e)$ fraction of the elements (as in a "no" instance of max-$k$-coverage), the cost of this solution would be at least $(1 - 1/e)cn + 2/ecn + 2(m - k) = (1 + 1/e)cn + 2(m - k)$; hence $C$ would be at least $(1 + 1/e - \delta)cn$ in this case. This shows that for every $\delta$, there are instances of $k$-median whose optimal cost is either at most $C$ or at least $(1 + 1/e - \delta)$, such that distinguishing between these two cases is NP-hard.

Let us now add infinitesimal noise to the above instances of $k$-median to make a unique optimal solution and call this the target; the uniqueness of the optimal solution ensures that we satisfy $(1, \epsilon)$-approximation-stability without changing the hardness significantly. Now, in the "yes" case, any clustering with error $c'\epsilon$ will have cost at most $(1 - c'\epsilon)cn + 2c'\epsilon cn + 2(m - k) \leq (1 + c'\epsilon + \delta)cn$. This is less than the cost of the optimal solution in the "no" case (which is still at least $(1 + 1/e - \delta)cn$) as long as $c'\epsilon \leq 1/e - 2\delta$, and would allow us to distingush 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.    □

## A.1    Proof of the Reassignment Lemma

We now prove Lemma 2, which we restate here for convenience.

LEMMA 2.    *Let $\mathcal{C} = \{C_1, \ldots, C_k\}$ be a $k$-clustering in which each cluster is nonempty, and let $R = \{(x_1, j_1), (x_2, j_2), \ldots, (x_t, j_t)\}$ 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 exactly $\frac{1}{n}|R'|$ from $\mathcal{C}$.*

**Note:** 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 $S_1, S_2, \ldots$, 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 $\sigma$ on the indices can have agreement between $\mathcal{C}$ and $\mathcal{C}'$ on at most one point from each $S_i$.

*A simpler factor-8 argument:* We begin for illustration with a simpler factor-8 argument. For this argument we consider two cases. First, suppose that at least half of the reassignments in $R$ involve points $x$ in clusters of size $\geq 2$. In this case, we simply do the following. While there exists some such $(x, j) \in R$, choose some

arbitrary point $y \in C(x)$ and add $\{x, y\}$ to $L$, add $(x, j)$ to $R'$, and then remove $(x, j)$ from $R$ as well as any reassignment involving $y$ if one exists; also remove both $x$ and $y$ from point set $S$. If cluster $C(x)$ has been reduced to a singleton $\{z\}$, then remove $z$ from $R$ and $S$ as well. This process guarantees that all pairs added to $L$ are disjoint, and we remove at most three times as many reassignments from $R$ as we add to $R'$. Thus, since we assumed at least half of $R$ came from such clusters, overall we get a factor of 6. The second case is that at least half of the reassignments in $R$ involve points $x$ in clusters of size 1. In that case, randomly color each cluster red or blue: in expectation, $1/4$ of these reassignments (at least $1/8$ of the total in $R$) go from red clusters to blue clusters. We now simply put all of these reassignments, namely those involving points moving from singleton red clusters to blue clusters, into $R'$. Because all such $(x, j)$ for any given $j$ involve different source clusters, we can create a witness set $S_j$ for each $j$ consisting of those points $x$ together with an arbitrary point $y \in C_j$, guaranteeing the quality of $R'$.

*The factor-3 argument:* For the factor-3 argument, we begin constructing $R'$ and $L$ as in the first case above, but using only clusters of size at least 3. Specifically, 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 any 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 $\leq 1$. 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 $\{x\} \cup \{z\} \cup \{y : (y, j)$ 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 $\{x\} \cup \{z\} \cup \{y : (y, j)$ 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.  $\square$

## B.   ANALYSIS OF EXAMPLE IN SECTION 6

In Section 6, an example is presented of $\sqrt{n}$ clusters of $\sqrt{n}$ points each, with distance 1 between points in the same target cluster, and distance $D + 1$ between points in different target clusters. We prove here that for any $\epsilon < 1/2$, this satisfies $(D\epsilon/2, \epsilon)$-approximation-stability for both $k$-median and $k$-means objectives. Thus, if $D > 4/\epsilon$, then this is $(2, \epsilon)$-approximation-stable.

Let $\mathcal{C}$ be a clustering of distance at least $\epsilon$ from the target clustering $\mathcal{C}_T = \mathcal{C}^*$. Since $\mathcal{C}^*$ has both $k$-median and $k$-means cost equal to $n - \sqrt{n}$, we need to show that $\mathcal{C}$ has $k$-median cost at least $(D\epsilon/2)(n - \sqrt{n})$ (its $k$-means cost can only be larger).

We do this as follows. First, define the "non-permutation distance" from $\mathcal{C}$ to $\mathcal{C}^*$ as $npdist(\mathcal{C}, \mathcal{C}^*) = \frac{1}{n} \sum_{i=1}^{k} \min_j |C_i - C_j^*|$. That is, we remove the restriction that different clusters in $\mathcal{C}$ cannot be mapped to the same cluster in $\mathcal{C}^*$. This is non-symmetric, but clearly satisfies the condition that $npdist(\mathcal{C}, \mathcal{C}^*) \leq dist(\mathcal{C}, \mathcal{C}^*)$. We observe now that the $k$-median cost of $\mathcal{C}$ is equal to $Dn \cdot npdist(\mathcal{C}, \mathcal{C}^*) + (n - \sqrt{n})$. In particular, the optimal median for each cluster $C_i$ in $\mathcal{C}$ is a point in whichever cluster $C_j^*$ of $\mathcal{C}^*$ has the largest intersection with $C_i$. This causes each point in $C_i - C_j^*$ to incur an additional cost of $D$ over its cost in $\mathcal{C}^*$, and so the overall increase over the cost of $\mathcal{C}^*$ is $Dn \cdot npdist(\mathcal{C}, \mathcal{C}^*)$. Thus, it remains just to show that $npdist(\mathcal{C}, \mathcal{C}^*)$ cannot be too much smaller than $dist(\mathcal{C}, \mathcal{C}^*)$.

We now show that $npdist(\mathcal{C}, \mathcal{C}^*) \geq dist(\mathcal{C}, \mathcal{C}^*)/2$. We note that this will rely heavily on the fact that all clusters in $\mathcal{C}^*$ have the same size: if $C^*$ contained clusters of very different sizes, the statement would be false. Since this inequality may be of interest more generally (it is not specific to this example), we formalize it in Lemma 21 below.

LEMMA 21. *For any clustering $\mathcal{C}$, if all clusters of $\mathcal{C}^*$ have size $n/k$, then we have $npdist(\mathcal{C}, \mathcal{C}^*) \geq dist(\mathcal{C}, \mathcal{C}^*)/2$.*

PROOF. Let $p_i = |C_i|/n$ and $p = (p_1, \ldots, p_k)$. Let $u = (1/k, \ldots, 1/k)$ and define $\Delta(p, u) = \sum_{i: p_i > u_i} p_i - u_i$ to be the variation distance between $p$ and $u$. Then, $npdist(\mathcal{C}, \mathcal{C}^*) \geq \Delta(p, u)$ because a cluster in $\mathcal{C}$ of size $p_i n > n/k$ contributes at least $p_i - 1/k$ to the non-permutation distance. Let $\Delta_i = \max(1/k - p_i, 0)$. Since variation distance is symmetric, we have $\Delta(p, u) = \sum_i \Delta_i$.

Now, fix some mapping of clusters $C_i$ to clusters $C_j^*$ yielding the non-permutation distance from $\mathcal{C}$ to $\mathcal{C}^*$. Let $T_j$ denote the set of indices $i$ such that $C_i$ is mapped to $C_j^*$ and let $t_j = |T_j|$. Let $S$ denote the set of indices $j$ such that $t_j \geq 2$ (if this were a permutation then $S$ would be empty). We can now lower-bound the non-permutation distance as

$$
\begin{aligned}
npdist(\mathcal{C}, \mathcal{C}^*) \quad &\geq \quad \sum_{j \in S} \left[ \left( \sum_{i \in T_j} p_i \right) - 1/k \right] \\
&\geq \quad \sum_{j \in S} \left[ \frac{t_j - 1}{k} - \sum_{i \in T_j} \Delta_i \right] \\
&\geq \quad \left( \sum_{j \in S} \frac{t_j - 1}{k} \right) - \Delta(p, u).
\end{aligned}
$$

Therefore, we have

$$
npdist(\mathcal{C}, \mathcal{C}^*) + \Delta(p, u) \geq \sum_{j \in S} \frac{t_j - 1}{k}. \tag{B.1}
$$

We now claim we can convert this mapping into a permutation without increasing the distance by too much. Specifically, for each $j$ such that $t_j \geq 2$, keep only the $i \in T_j$ such that $C_i$ has highest overlap with $C_j^*$ and assign the rest to (arbitrary) unmatched target clusters. This reassignment can increase the distance computation by at most $\frac{1}{k}(\frac{t_j - 1}{t_j}) \leq \frac{1}{k}(\frac{t_j - 1}{2})$. Therefore, we have

$$
dist(\mathcal{C}, \mathcal{C}^*) \leq npdist(\mathcal{C}, \mathcal{C}^*) + \sum_{j \in S} \frac{t_j - 1}{2k}. \tag{B.2}
$$

Combining (B.1) and (B.2) we have $dist(\mathcal{C}, \mathcal{C}^*) - npdist(\mathcal{C}, \mathcal{C}^*) \leq \frac{1}{2}(npdist(\mathcal{C}, \mathcal{C}^*) + \Delta(p, u))$, and since $\Delta(p, u) \leq npdist(\mathcal{C}, \mathcal{C}^*)$, this yields $dist(\mathcal{C}, \mathcal{C}^*) \leq 2npdist(\mathcal{C}, \mathcal{C}^*)$ as desired. □

Finally, as noted above, by Lemma 21 we have that the cost of clustering $\mathcal{C}$ in the construction is at least $Dn\epsilon/2$ as desired.