
Multi-Swap k -Means++

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The k -means++ algorithm of Arthur and Vassilvitskii (SODA 2007) is often the
2 practitioners' choice algorithm for optimizing the popular k -means clustering ob-
3 jective and is known to give an $O(\log k)$ -approximation in expectation. To obtain
4 higher quality solutions, Lattanzi and Sohler (ICML 2019) proposed augmenting
5 k -means++ with $O(k \log \log k)$ local search steps obtained through the k -means++
6 sampling distribution to yield a c -approximation to the k -means clustering problem,
7 where c is a large absolute constant. Here we generalize and extend their local
8 search algorithm by considering larger and more sophisticated local search neigh-
9 borhoods hence allowing to swap multiple centers at the same time. Our algorithm
10 achieves a $9 + \varepsilon$ approximation ratio, which is the best possible for local search.
11 Importantly we show that our approach yields substantial practical improvements,
12 we show significant quality improvements over the approach of Lattanzi and Sohler
13 (ICML 2019) on several datasets.

14 1 Introduction

15 Clustering is a central problem in unsupervised learning. In clustering one is interested in grouping
16 together "similar" object and separate "dissimilar" one. Thanks to its popularity many notions of
17 clustering have been proposed overtime. In this paper, we focus on metric clustering and on one of
18 the most studied problem in the area: the Euclidean k -means problem.

19 In the Euclidean k -means problem one is given in input a set of points P in \mathbb{R}^d . The goal of
20 the problem is to find a set of k centers so that the sum of the square distances to the centers
21 is minimized. More formally, we are interested in finding a set C of k points in \mathbb{R}^d such that
22 $\sum_{p \in P} \min_{c \in C} \|p - c\|^2$, where with $\|p - c\|$ we denote the Euclidean distance between p and c .

23 The k -means problem has a long history, in statistics and operations research. For Euclidean k -
24 means with running time polynomial in both n, k and d , a 5.912-approximation was recently shown
25 in Cohen-Addad et al. [2022a], improving upon Kanungo et al. [2004], Ahmadian et al. [2019],
26 Grandoni et al. [2022] by leveraging the properties of the Euclidean metric. In terms of lower bounds,
27 the first to show that the high-dimensional k -means problems were APX-hard were Guruswami and
28 Indyk [2003], and later Awasthi et al. [2015] showed that the APX-hardness holds even if the centers
29 can be placed arbitrarily in \mathbb{R}^d . The inapproximability bound was later slightly improved by Lee et al.
30 [2017] until the recent best known bounds of Cohen-Addad and Karthik C. S. [2019], Cohen-Addad
31 et al. [2022d] that showed that it is NP-hard to achieve a better than 1.06-approximation and hard to
32 approximate it better than 1.36 assuming a stronger conjecture. From a more practical point of view,
33 Arthur and Vassilvitskii [2009] showed that the widely-used popular heuristic of Lloyd Lloyd [1957]
34 can lead to solutions with arbitrarily bad approximation guarantees, but can be improved by a simple
35 seeding strategy, called k -means++, so as to guarantee that the output is within an $O(\log k)$ factor of
36 the optimum Arthur and Vassilvitskii [2007].

37 Thanks to its simplicity k -means++ is widely adopted in practice. In an effort to improve its
 38 performances Lattanzi and Sohler [2019], Choo et al. [2020] combine k -means++ and local search to
 39 efficiently obtain a constant approximation algorithm with good practical performance. These two
 40 studies show that one can use the k -means++ distribution in combination with a local search algorithm
 41 to get the best of both worlds: a practical algorithm with constant approximation guarantees.

42 However, the constant obtained in Lattanzi and Sohler [2019], Choo et al. [2020] is very large (several
 43 thousands in theory) and the question as whether one could obtain a practical algorithm that would
 44 efficiently match the $9 + \varepsilon$ -approximation obtained by the $n^{O(d/\varepsilon)}$ algorithm of Kanungo et al. [2004]
 45 has remained open. Bridging the gap between the theoretical approach of Kanungo et al. [2004] and
 46 k -means++ has thus been a long standing goal.

47 **Our Contributions.** We make significant progress on the above line of work.

- 48 • We adapt techniques from the analysis of Kanungo et al. [2004] to obtain a tighter analysis
 49 of the algorithm in Lattanzi and Sohler [2019]. In particular in Corollary 4, we show that
 50 their algorithm achieves an approximation of ratio of ≈ 26.64 .
- 51 • We extend this approach to multi-swaps, where we allow swapping more than one center at
 52 each iteration of local search, improving significantly the approximation to ≈ 10.48 in time
 53 $O(nd \cdot \text{poly}(k))$.
- 54 • Leveraging ideas from Cohen-Addad et al. [2021], we design a better local search swap that
 55 improves the approximation further to $9 + \varepsilon$ (see Theorem 12). This new algorithm matches
 56 the $9 + \varepsilon$ -approximation achieved by the local search algorithm in Kanungo et al. [2004], but
 57 it is significantly more efficient. Notice that 9 is the best approximation achievable through
 58 local search algorithms, as proved in Kanungo et al. [2004].
- 59 • We provide experiments where we compare against k -means++ and Lattanzi and Sohler
 60 [2019]. We study a variant of our algorithm that performs very competitively with our
 61 theoretically sound algorithm. The variant is very efficient and still outperforms previous
 62 work in terms of solution quality, even after the standard postprocessing using Lloyd.

63 **Additional Related Work.** We start by reviewing the approach of Kanungo et al. [2004] and a
 64 possible adaptation to our setting. The bound of $9 + \varepsilon$ on the approximation guarantee shown by
 65 Kanungo et al. [2004] is for the following algorithm: Given a set S of k centers, if there is a set
 66 S^+ of at most $2/\varepsilon$ points in \mathbb{R}^d together with a set S^- of $|S^+|$ points in S such that $S \setminus S^- \cup S^+$
 67 achieves a better k -means cost than S , then set $S := S \setminus S^- \cup S^+$ and repeat until convergence.
 68 The main drawback of the algorithm is that it asks whether there exists a set S^+ of points in \mathbb{R}^d that
 69 could be swapped with elements of S to improve the cost. Identifying such a set, even of constant
 70 size, is already non-trivial. The best way of doing so is through the following path: First compute
 71 a coresset using the state-of-the-art coresset construction of Cohen-Addad et al. [2022b] and apply
 72 the dimensionality reduction of Becchetti et al. [2019], Makarychev et al. [2019], hence obtaining
 73 a set of $\tilde{O}(k/\varepsilon^4)$ points in dimension $O(\log k/\varepsilon^2)$. Then, compute grids using the discretization
 74 framework of Matousek [2000] to identify a set of $\varepsilon^{-O(d)} \sim k^{O(\varepsilon^{-2} \log(1/\varepsilon))}$ grid points that contains
 75 nearly-optimum centers. Now, run the local search algorithm where the sets S^+ are chosen from
 76 the grid points by brute-force enumeration over all possible subsets of grid points of size at most,
 77 say s . The running time of the whole algorithm with swaps of magnitude s , i.e.: $|S^+| \leq s$, hence
 78 becomes $k^{O(s \cdot \varepsilon^{-2} \log(1/\varepsilon))}$ for an approximation of $(1 + \varepsilon)(9 + 2/s)$, meaning a dependency in k
 79 of $k^{O(\varepsilon^{-3} \log(1/\varepsilon))}$ to achieve a $9 + \varepsilon$ -approximation. Our results improves upon this approach in
 80 two ways: (1) it improves over the above theoretical bound and (2) does so through an efficient and
 81 implementable, i.e.: practical, algorithm.

82 Recently, Grunau et al. [2023] looked at how much applying a greedy rule on top of the k -means++
 83 heuristic improves its performance. The heuristic is that at each step, the algorithm samples ℓ centers
 84 and only keeps the one that gives the best improvement in cost. Interestingly the authors prove that
 85 from a theoretical standpoint this heuristic does not improve the quality of the output. Local search
 86 algorithms for k -median and k -means have also been studied by Gupta and Tangwongsan [2008] who
 87 drastically simplified the analysis of Arya et al. [2004]. Cohen-Addad and Schwiegelshohn [2017]
 88 demonstrated the power of local search for stable instances. Friggstad et al. [2019], Cohen-Addad
 89 et al. [2019] showed that local search yields a PTAS for Euclidean inputs of bounded dimension (and
 90 doubling metrics) and minor-free metrics. Cohen-Addad [2018] showed how to speed up the local
 91 search algorithm using k d -trees (i.e.: for low dimensional inputs).

92 For fixed k , there are several known approximation schemes, typically using small coresets Becchetti
 93 et al. [2019], Feldman and Langberg [2011], Kumar et al. [2010]. The state-of-the-art approaches
 94 are due to Bhattacharya et al. [2020], Jaiswal et al. [2014]. The best known coreset construction
 95 remains Cohen-Addad et al. [2022c,b].

96 If the constraint on the number of output centers is relaxed, then we talk about bicriteria approxima-
 97 tions and k -means has been largely studied Bandyapadhyay and Varadarajan [2016], Charikar and
 98 Guha [2005], Cohen-Addad and Mathieu [2015], Korupolu et al. [2000], Makarychev et al. [2016].

99 2 Preliminaries

100 **Notation.** We denote with $P \subseteq \mathbb{R}^d$ the set of input points and let $n = |P|$. Given a point set $Q \subseteq P$
 101 we use $\mu(Q)$ to denote the mean of points in Q . Given a point $p \in P$ and a set of centers A we denote
 102 with $A[p]$ the closest center in A to p (ties are broken arbitrarily). We denote with \mathcal{C} the set of centers
 103 currently found by our algorithm and with \mathcal{O}^* an optimal set of centers. Therefore, given $p \in P$,
 104 we denote with $\mathcal{C}[p]$ and $\mathcal{O}^*[p]$ its closest ALG-center and OPT-center respectively. We denote by
 105 $\text{cost}(Q, A)$ the cost of points in $Q \subseteq P$ w.r.t. the centers in A , namely

$$\text{cost}(Q, A) = \sum_{q \in Q} \min_{c \in A} \|q - c\|^2.$$

106 We use ALG and OPT as a shorthand for $\text{cost}(P, \mathcal{C})$ and $\text{cost}(P, \mathcal{O}^*)$ respectively. When we sample
 107 points proportionally to their current cost (namely, sample q with probability $\text{cost}(q, \mathcal{C}) / \text{cost}(P, \mathcal{C})$)
 108 we call this the D^2 distribution. When using $O_\varepsilon(\cdot)$ and $\Omega_\varepsilon(\cdot)$ we mean that ε is considered constant.
 109 We use $\tilde{O}(f)$ to hide polylogarithmic factors in f . The following lemma is folklore.

110 **Lemma 1.** *Given a point set $Q \subseteq P$ and a point $p \in P$ we have*

$$\text{cost}(Q, p) = \text{cost}(Q, \mu(Q)) + |Q| \cdot \|p - \mu(Q)\|^2.$$

111 Let O_i^* be an optimal cluster, we define the *radius* of O_i^* as ρ_i such that $\rho_i^2 \cdot |O_i^*| = \text{cost}(O_i^*, o_i)$,
 112 where $o_i = \mu(O_i^*)$. We define the δ -core of the optimal cluster O_i^* as the set of points $p \in O_i^*$
 113 that lie in a ball of radius $(1 + \delta)\rho_i$ centered in o_i . In symbols, $\text{core}(O_i^*) = P \cap B(o_i, (1 + \delta)\rho_i)$.
 114 Throughout the paper, δ is always a small constant fixed upfront, hence we omit it.

115 **Lemma 2.** *Let O_i^* be an optimal cluster and sample $q \in O_i^*$ according to the D^2 -distribution
 116 restricted to O_i^* . If $\text{cost}(O_i^*, \mathcal{C}) > (2 + 3\delta) \cdot \text{cost}(O_i^*, o_i)$ then $\Pr[q \in \text{core}(O_i^*)] = \Omega_\delta(1)$.*

117 *Proof.* Define $\alpha := \text{cost}(O_i^*, \mathcal{C}) / \text{cost}(O_i^*, o_i) > 2 + 3\delta$. Thanks to Lemma 1, for each $c \in \mathcal{C}$ we
 118 have $\|c - o_i\|^2 \geq (\alpha - 1)\rho_i^2$. Therefore, for each $y \in \text{core}(O_i^*)$ and every $c \in \mathcal{C}$ we have

$$\text{cost}(y, c) = \|y - c\|^2 \geq (\sqrt{\alpha - 1} - (1 + \delta))^2 \cdot \rho_i^2 = \Omega_\delta(\alpha \rho_i^2).$$

119 Moreover, by a Markov's inequality argument we have $|O_i^* \setminus \text{core}(O_i^*)| \leq \frac{1}{1 + \delta} \cdot |O_i^*|$ and thus
 120 $|\text{core}(O_i^*)| \geq \Omega_\delta(|O_i^*|)$. Combining everything we get

$$\text{cost}(\text{core}(O_i^*), \mathcal{C}) \geq |\text{core}(O_i^*)| \cdot \min_{\substack{c \in \mathcal{C} \\ y \in \text{core}(O_i^*)}} \text{cost}(y, c) = \Omega_\delta(|O_i^*|) \cdot \Omega_\delta(\alpha \rho_i^2)$$

121 and $|O_i^*| \cdot \alpha \rho_i^2 = \text{cost}(O_i^*, \mathcal{C})$, hence $\text{cost}(\text{core}(O_i^*), \mathcal{C}) = \Omega_\delta(\text{cost}(O_i^*, \mathcal{C}))$. \square

122 3 Multi-Swap k -Means++

123 The single-swap local search (SSLS) k -means++ algorithm in Lattanzi and Sohler [2019] works
 124 as follows. First, k centers are sampled using k -means++ (namely, they are sampled one by one
 125 according to the D^2 distribution, updated for every new center). Then, $O(k \log \log k)$ steps of local
 126 search follow. In each local search step a point $q \in P$ is D^2 -sampled, then let c be the center among
 127 the current centers \mathcal{C} such that $\text{cost}(P, (\mathcal{C} \setminus \{c\}) \cup \{q\})$ is minimum. If $\text{cost}(P, (\mathcal{C} \setminus \{c\}) \cup \{q\}) <$
 128 $\text{cost}(P, \mathcal{C})$ then we swap c and q , or more formally we set $\mathcal{C} \leftarrow (\mathcal{C} \setminus \{c\}) \cup \{q\}$.

129 We extend the SSLS so that we allow to swap multiple centers simultaneously and call this algorithm
 130 multi-swap local search (MSLS) k -means++. Swapping multiple centers at the same time achieves a
 131 lower approximation ratio, in exchange for a higher time complexity. In this section, we present and
 132 analyse the p -swap local search (LS) algorithm for a generic number of p centers swapped at each
 133 step. For any constant $\delta > 0$, we obtain an approximation ratio $\text{ALG}/\text{OPT} = \eta^2 + \delta$ where

$$\eta^2 - (2 + 2/p)\eta - (4 + 2/p) = 0. \quad (1)$$

134 **The Algorithm.** First, we initialize our set of centers using k -means++. Then, we run $O(ndk^{p-1})$
 135 local search steps, where a local search step works as follows. We D^2 -sample a set $In = \{q_1 \dots q_p\}$
 136 of points from P (without updating costs). Then, we iterate over all possible sets $Out = \{c_1 \dots c_p\}$
 137 of p distinct elements in $\mathcal{C} \cup In$ and select the set Out such that performing the swap (In, Out)
 138 maximally improves the cost¹. If this choice of Out improves the cost, then we perform the swap
 139 (In, Out) , else we do not perform any swap for this step.

140 **Theorem 3.** For any $\delta > 0$, the p -swap local search algorithm above runs in $\tilde{O}(ndk^{2p})$ time and,
 141 with constant probability, finds an $(\eta^2 + \delta)$ -approximation of k -means, where η satisfies Equation (1).

142 Notice that the SLS algorithm of Lattanzi and Sohler [2019] is exactly the p -swap LS algorithm
 143 above for $p = 1$.

144 **Corollary 4.** The single-swap local search in Lattanzi and Sohler [2019], Choo et al. [2020] achieves
 145 an approximation ratio < 26.64 .

146 **Corollary 5.** For $p = O(1)$ large enough, multi-swap local search achieves an approximation ratio
 147 < 10.48 in time $O(nd \cdot \text{poly}(k))$.

148 3.1 Analysis of Multi-Swap k -means++

149 In this section we prove Theorem 3. Our main stepping stone is the following lemma.

150 **Lemma 6.** Let ALG denote the cost at some point in the execution of $MMLS$. As long as $ALG/OPT >$
 151 $\eta^2 + \delta$, a local search step improves the cost by a factor $1 - \Omega(1/k)$ with probability $\Omega(1/k^{p-1})$.

152 *Proof of Theorem 3.* First, we show that $O(k^p \log \log k)$ local steps suffice to obtain the desired
 153 approximation ratio, with constant probability. Notice that a local search step can only improve the
 154 cost function, so it is sufficient to show that the approximation ratio is achieved at some point in time.
 155 We initialize our centers using k -means++, which gives a $O(\log k)$ -approximation in expectation.
 156 Thus, using Markov's inequality the approximation guarantee $O(\log k)$ holds with arbitrary high
 157 constant probability. We say that a local-search step is *successful* if it improves the cost by a factor of
 158 at least $1 - \Omega(1/k)$. Thanks to Lemma 6, we know that unless the algorithm has already achieved the
 159 desired approximation ratio then a local-search step is successful with probability $\Omega(1/k^{p-1})$. To go
 160 from $O(\log k)$ to $\eta^2 + \delta$ we need $O(k \log \log k)$ successful local search steps. Standard concentration
 161 bounds on the value of a Negative Binomial random variable show that, with high probability, the
 162 number of trial to obtain $O(k \log \log k)$ successful local-search steps is $O(k^p \log \log k)$. Therefore,
 163 after $O(k^p \log \log k)$ local-search steps we obtain an approximation ratio of $\eta^2 + \delta$.

164 To prove the running time bound it is sufficient to show that a local search step can be performed in
 165 time $\tilde{O}(ndk^{p-1})$. This is possible if we maintain, for each point $x \in P$, a dynamic sorted dictionary²
 166 storing the pairs $(\text{cost}(x, c_i), c_i)$ for each $c_i \in \mathcal{C}$. Then we can combine the exhaustive search
 167 over all possible size- p subsets of $\mathcal{C} \cup In$ and the computation of the new cost function using time
 168 $O(ndk^{p-1} \log k)$. To do so, we iterate over all possible size- $(p-1)$ subsets Z of $\mathcal{C} \cup In$ and update
 169 all costs as if these centers were removed, then for each point $x \in P$ we compute how much its cost
 170 increases if we remove its closest center c_x in $(\mathcal{C} \cup In) \setminus Z$ and charge that amount to c_x . In the end,
 171 we consider $Out = Z \cup \{c\}$ where c is the cheapest-to-remove center found in this way. \square

172 The rest of this section is devoted to proving Lemma 6. For convenience, we prove that Lemma 6
 173 holds whenever $ALG/OPT > \eta^2 + O(\delta)$, which is wlog by rescaling δ . Recall that we now focus on
 174 a given step of the algorithm, and when we say current cost, current centers and current clusters we
 175 refer to the state of these objects at the end of the last local-search step before the current one. Let
 176 $O_1^* \dots O_k^*$ be an optimal clustering of P and let $\mathcal{O}^* = \{o_i = \mu(O_i^*) \mid \text{for } i = 1 \dots k\}$ be the set of
 177 optimal centers of these clusters. We denote with $C_1 \dots C_k$ the current set of clusters at that stage of
 178 the local search and with $\mathcal{C} = \{c_1 \dots c_k\}$ the set of their respective current centers.

179 We say that c_i captures o_j if c_i is the closest current center to o_j , namely $c_i = \mathcal{C}[o_j]$. We say that
 180 c_i is *busy* if it captures more than p optimal centers, and we say it is *lonely* if it captures no optimal
 181 center. Let $\tilde{\mathcal{O}} = \{o_i \mid \text{cost}(O_i^*, \mathcal{C}) > \delta \cdot ALG/k\}$ and $\tilde{\mathcal{C}} = \mathcal{C} \setminus \{\mathcal{C}[o_i] \mid o_i \in \mathcal{O}^* \setminus \tilde{\mathcal{O}}\}$. For ease of
 182 notation, we simply assume that $\tilde{\mathcal{O}} = \{o_1 \dots o_h\}$ and $\tilde{\mathcal{C}} = \{c_1 \dots c_{h'}\}$. Notice that $h' > h$.

¹If $In \cap Out \neq \emptyset$ then we are actually performing the swap $(In \setminus Out, Out \setminus In)$ of size $< p$.

²Also known as dynamic predecessor search data structure.

183 **Weighted ideal multi-swaps.** Given $In \subseteq P$ and $Out \subseteq \tilde{\mathcal{C}}$ of the same size we say that the
184 swap (In, Out) is an *ideal swap* if $In \subseteq \tilde{\mathcal{O}}$. We now build a set of *weighted ideal multi-swaps*
185 \mathcal{S} . First, suppose wlog that $\{c_1 \dots c_t\}$ is the set of current centers in $\tilde{\mathcal{C}}$ that are neither lonely nor
186 busy. Let \mathcal{L} be the set of lonely centers in $\tilde{\mathcal{C}}$. For each $i = 1 \dots t$, we do the following. Let In
187 be the set of optimal centers in $\tilde{\mathcal{O}}$ captured by c_i . Choose a set \mathcal{L}_i of $|In| - 1$ centers from \mathcal{L} , set
188 $\mathcal{L} \leftarrow \mathcal{L} \setminus \mathcal{L}_i$ and define $Out = \mathcal{L}_i \cup \{c_i\}$. Assign weight 1 to (In, Out) and add it to \mathcal{S} . For each
189 busy center $c_i \in \{c_{t+1} \dots c_{h'}\}$ let A be the set of optimal centers in $\tilde{\mathcal{O}}$ captured by c_i , pick a set \mathcal{L}_i
190 of $|A| - 1$ lonely current centers from \mathcal{L} (a counting argument shows that this is always possible).
191 Set $\mathcal{L} \leftarrow \mathcal{L} \setminus \mathcal{L}_i$. For each $o_j \in A$ and $c_\ell \in \mathcal{L}_i$ assign weight $1/(|A| - 1)$ to (o_j, c_ℓ) and add it to \mathcal{S} .
192 **Observation 7.** *The process above generates a set of weighted ideal multi-swaps such that: (i) Every*
193 *swap has size at most p ; (ii) The combined weights of swaps involving an optimal center o_i is 1; (iii)*
194 *The combined weights of swaps involving a current center c_i is at most $1 + 1/p$.*

195 Consider an ideal swap (In, Out) . Let $O_{In}^* = \bigcup_{o_i \in In} O_i^*$ and $C_{Out} = \bigcup_{c_j \in Out} C_j$. Define the
196 reassignment cost $\text{Reassign}(In, Out)$ as the increase in cost of reassigning points in $C_{Out} \setminus O_{In}^*$ to
197 centers in $\mathcal{C} \setminus Out$. Namely,

$$\text{Reassign}(In, Out) = \text{cost}(C_{Out} \setminus O_{In}^*, \mathcal{C} \setminus Out) - \text{cost}(C_{Out} \setminus O_{In}^*, \mathcal{C}).$$

198 We take the increase in cost of the following reassignment as an upper bound to the reassignment
199 cost. For each $p \in C_{Out} \setminus O_{In}^*$ we consider its closest optimal center $\mathcal{O}^*[p]$ and reassign p to the
200 current center that is closest to $\mathcal{O}^*[p]$, namely $\mathcal{C}[\mathcal{O}^*[p]]$. In formulas, we have

$$\begin{aligned} \text{Reassign}(In, Out) &\leq \sum_{p \in C_{Out} \setminus O_{In}^*} \text{cost}(p, \mathcal{C}[\mathcal{O}^*[p]]) - \text{cost}(p, \mathcal{C}[p]) \\ &\leq \sum_{p \in C_{Out}} \text{cost}(p, \mathcal{C}[\mathcal{O}^*[p]]) - \text{cost}(p, \mathcal{C}[p]). \end{aligned}$$

201 Indeed, by the way we defined our ideal swaps we have $\mathcal{C}[\mathcal{O}^*[p]] \notin Out$ for each $p \notin O_{In}^*$ and this
202 reassignment is valid. Notice that the right hand side in the equation above does not depend on In .

203 **Lemma 8.** $\sum_{p \in P} \text{cost}(p, \mathcal{C}[\mathcal{O}^*[p]]) \leq 2OPT + ALG + 2\sqrt{ALG}\sqrt{OPT}$.

204 *Proof.* Deferred to the supplementary material. \square

205 **Lemma 9.** *The combined weighted reassignment costs of all ideal multi-swaps in \mathcal{S} is at most*
206 $(2 + 2/p) \cdot (OPT + \sqrt{ALG}\sqrt{OPT})$.

207 *Proof.* Denote by $w(In, Out)$ the weight associated with the swap (In, Out) .

$$\begin{aligned} &\sum_{(In, Out) \in \mathcal{S}} w(In, Out) \cdot \text{Reassign}(In, Out) \leq \\ &\sum_{(In, Out) \in \mathcal{S}} w(In, Out) \cdot \sum_{p \in C_{Out}} \text{cost}(p, \mathcal{C}[\mathcal{O}^*[p]]) - \text{cost}(p, \mathcal{C}[p]) \leq \\ &(1 + 1/p) \cdot \sum_{c_j \in \mathcal{C}} \sum_{p \in C_j} \text{cost}(p, \mathcal{C}[\mathcal{O}^*[p]]) - \text{cost}(p, \mathcal{C}[p]) \leq \\ &(1 + 1/p) \cdot \left(\sum_{p \in P} \text{cost}(p, \mathcal{C}[\mathcal{O}^*[p]]) - ALG \right). \end{aligned}$$

208 The second inequality uses (iii) from Observation 7. Applying Lemma 8 completes the proof. \square

209 Recall the notions of radius and core of an optimal cluster introduced in Section 2. We say that a
210 swap (In, Out) is *strongly improving* if $\text{cost}(P, (\mathcal{C} \cup In) \setminus Out) \leq (1 - \delta/k) \cdot \text{cost}(P, \mathcal{C})$. Let
211 $In = \{o_1 \dots o_s\} \subseteq \tilde{\mathcal{O}}$ and $Out = \{c_1 \dots c_s\} \subseteq \tilde{\mathcal{C}}$ we say that an ideal swap (In, Out) is *good*
212 if for every $q_1 \in \text{core}(o_1) \dots q_s \in \text{core}(o_s)$ the swap (\mathcal{Q}, Out) is strongly improving, where
213 $\mathcal{Q} = \{q_1 \dots q_s\}$. We call an ideal swap *bad* otherwise. We say that an optimal center $o_i \in \tilde{\mathcal{O}}$ is good
214 if that's the case for at least one of the ideal swaps it belongs to. Recalling how we constructed our

215 set of weighted ideal swaps \mathcal{S} , if o_i belongs to a s -swap $(In, Out) \in \mathcal{S}$ for $s > 1$ then it is good if
 216 and only if (In, Out) is good; else o_i might belong to up to $p - 1$ single swaps and then any of them
 217 being good would suffice. Denote with G the union of cores of good optimal centers in $\tilde{\mathcal{O}}$.

218 **Lemma 10.** *If an ideal swap (In, Out) is bad, then we have*

$$\text{cost}(O_{In}^*, \mathcal{C}) \leq (2 + \delta)\text{cost}(O_{In}^*, \mathcal{O}^*) + \text{Reassign}(In, Out) + \delta\text{ALG}/k. \quad (2)$$

219 *Proof.* Let $In = \{o_1 \dots o_s\}$, $\mathcal{Q} = \{q_1 \dots q_s\}$ such that $q_1 \in \text{core}(o_1) \dots q_s \in \text{core}(o_s)$.
 220 Then, by Lemma 1 $\text{cost}(O_{In}^*, \mathcal{Q}) \leq (2 + \delta)\text{cost}(O_{In}^*, \mathcal{O}^*)$. Moreover, $\text{Reassign}(In, Out) =$
 221 $\text{cost}(P \setminus O_{In}^*, \mathcal{C} \setminus Out) - \text{cost}(P \setminus O_{In}^*, \mathcal{C})$ because points in $P \setminus C_{Out}$ are not affected by the
 222 swap. Therefore, $\text{cost}(P, (\mathcal{C} \cup \mathcal{Q}) \setminus Out) \leq (2 + \delta)\text{cost}(O_{In}^*, \mathcal{O}^*) + \text{Reassign}(In, Out) +$
 223 $\text{cost}(P \setminus O_{In}^*, \mathcal{C})$. Suppose by contradiction that Equation (4) does not hold, then

$$\begin{aligned} & \text{cost}(P, \mathcal{C}) - \text{cost}(P, (\mathcal{C} \cup \mathcal{Q}) \setminus Out) = \\ & \text{cost}(P \setminus O_{In}^*, \mathcal{C}) + \text{cost}(O_{In}^*, \mathcal{C}) - \text{cost}(P, (\mathcal{C} \cup \mathcal{Q}) \setminus Out) \geq \delta\text{ALG}/k. \end{aligned}$$

224 Hence, (\mathcal{Q}, Out) is strongly improving and this holds for any choice of \mathcal{Q} , contradiction. \square

225 **Lemma 11.** *If $\text{ALG}/\text{OPT} > \eta^2 + \delta$ then $\text{cost}(G, \mathcal{C}) = \Omega_\delta(\text{cost}(P, \mathcal{C}))$. Thus, if we D^2 -sample q
 226 we have $P[q \in G] = \Omega_\delta(1)$.*

227 *Proof.* First, we observe that the combined current cost of all optimal clusters in $\mathcal{O}^* \setminus \tilde{\mathcal{O}}$ is at most
 228 $k \cdot \delta\text{ALG}/k = \delta\text{ALG}$. Now, we prove that the combined current cost of all O_i^* such that o_i is bad is
 229 $\leq (1 - 2\delta)\text{ALG}$. Suppose, by contradiction, that it is not the case, then we have:

$$\begin{aligned} (1 - 2\delta)\text{ALG} &< \sum_{\text{Bad } o_i \in \tilde{\mathcal{O}}} \text{cost}(O_i^*, \mathcal{C}) \leq \sum_{\text{Bad } (In, Out) \in \mathcal{S}} w(In, Out) \cdot \text{cost}(O_{In}^*, \mathcal{C}) \leq \\ & \sum_{\text{Bad } (In, Out)} w(In, Out) \cdot ((2 + \delta)\text{cost}(O_{In}^*, \mathcal{O}^*) + \text{Reassign}(In, Out) + \delta\text{ALG}/k) \leq \\ & (2 + \delta)\text{OPT} + (2 + 2/p)\text{OPT} + (2 + 2/p)\sqrt{\text{ALG}}\sqrt{\text{OPT}} + \delta\text{ALG}. \end{aligned}$$

230 The second and last inequalities make use of Observation 7. The third inequality uses Lemma 10.

231 Setting $\eta^2 = \text{ALG}/\text{OPT}$ we obtain the inequality $\eta^2 - (2 + 2/p \pm O(\delta))\eta - (4 + 2/p \pm O(\delta)) \leq 0$.
 232 Hence, we obtain a contradiction in the previous argument as long as $\eta^2 - (2 + 2/p \pm O(\delta))\eta -$
 233 $(4 + 2/p \pm O(\delta)) > 0$. A contradiction there implies that at least an δ -fraction of the current cost is
 234 due to points in $\bigcup_{\text{Good } o_i \in \tilde{\mathcal{O}}} O_i^*$. We combine this with Lemma 2 and conclude that the total current
 235 cost of $G = \bigcup_{\text{Good } o_i \in \tilde{\mathcal{O}}} \text{core}(O_i^*)$ is $\Omega_\delta(\text{cost}(P, \mathcal{C}))$. \square

236 Finally, we prove Lemma 6. Whenever $q_1 \in G$ we have that $q_1 \in \text{core}(o_1)$ for some good o_1 .
 237 Then, for some $s \leq p$ we can complete o_1 with $o_2 \dots o_s$ such that $In = \{o_1 \dots o_s\}$ belongs
 238 to a good swap. Concretely, there exists $Out \subseteq \mathcal{C}$ such that (In, Out) is a good swap. Since
 239 $In \subset \tilde{\mathcal{O}}$ we have $\text{cost}(O_i^*, \mathcal{C}) > \delta\text{OPT}/k$ for all $o_i \in In$, which combined with Lemma 2 gives
 240 that for $i = 2 \dots s$ $P[q_i \in \text{core}(o_i)] \geq \Omega_\delta(1/k)$. Hence, we have $P[q_i \in \text{core}(o_i)$ for $i =$
 241 $1 \dots s] \geq \Omega_{\delta, p}(1/k^{p-1})$. Whenever we sample $q_1 \dots q_s$ from $\text{core}(o_1) \dots \text{core}(o_s)$, we have that
 242 (\mathcal{Q}, Out) is strongly improving. Notice, however, that (\mathcal{Q}, Out) is a s -swap and we may have $s < p$.
 243 Nevertheless, whenever we sample $q_1 \dots q_s$ followed by any sequence $q_{s+1} \dots q_p$ it is enough to
 244 choose $Out' = Out \cup \{q_{s+1} \dots q_p\}$ to obtain that $(\{q_1 \dots q_p\}, Out')$ is an improving p -swap.

245 4 A Faster $(9 + \varepsilon)$ -Approximation Local Search Algorithm

246 The MSLS algorithm from Section 3 achieves an approximation ratio of $\eta^2 + \varepsilon$, where $\eta^2 - (2 +$
 247 $2/p)\eta - (4 + 2/p) = 0$ and $\varepsilon > 0$ is an arbitrary small constant. For large p we have $\eta \approx 10.48$. On
 248 the other hand, employing p simultaneous swaps, Kanungo et al. [2004] achieve an approximation
 249 factor of $\xi^2 + \varepsilon$ where $\xi^2 - (2 + 2/p)\xi - (3 + 2/p) = 0$. If we set $p \approx 1/\varepsilon$ this yields a $(9 + O(\varepsilon))$ -
 250 approximation. In the same paper, they prove that 9-approximation is indeed the best possible for
 251 p -swap local search, if p is constant (see Theorem 3.1 in Kanungo et al. [2004]). They showed that 9
 252 is the right locality gap for local search, but they matched it with a very slow algorithm. To achieve a
 253 $(9 + \varepsilon)$ -approximation, they discretize the space reducing to $O(n\varepsilon^{-d})$ candidate centers and perform
 254 an exhaustive search over all size- $(1/\varepsilon)$ subsets of candidates at every step. As we saw in the related

255 work section, it is possible to combine techniques from coresets and dimensionality reduction to
 256 reduce the number of points to $n' = k \cdot \text{poly}(\varepsilon^{-1})$ and the number of dimensions to $d' = \log k \cdot \varepsilon^{-2}$.
 257 This reduces the complexity of Kanungo et al. [2004] to $k^{O(\varepsilon^{-3} \log \varepsilon^{-1})}$.

258 In this section, we leverage techniques from Cohen-Addad et al. [2021] to achieve a $(9 + \varepsilon)$ -
 259 approximation faster³. In particular, we obtain the following.

260 **Theorem 12.** *Given a set of n points in \mathbb{R}^d with aspect ratio Δ , there exists an algorithm that*
 261 *computes a $9 + \varepsilon$ -approximation to k -means in time $ndk^{O(\varepsilon^{-2})} \log^{O(\varepsilon^{-1})}(\Delta) \cdot 2^{-\text{poly}(\varepsilon^{-1})}$.*

262 Notice that, besides being asymptotically slower, the pipeline obtained combining known techniques
 263 is highly impractical and thus it did not make for an experimental test-bed. Moreover, it is not obvious
 264 how to simplify such an ensemble of complex techniques to obtain a practical algorithm.

265 **Limitations of MSLS.** The barrier we need to overcome in order to match the bound in Kanungo
 266 et al. [2004] is that, while we only consider points in P as candidate centers, the discretization
 267 they employ considers also points in $\mathbb{R}^d \setminus P$. In the analysis of MSLS we show that we sample
 268 each point q_i from $\text{core}(O_i^*)$ or equivalently that $q_i \in B(o_i, (1 + \varepsilon)\rho_i)$, where ρ_i is such that O_i^*
 269 would have the same cost w.r.t. o_i if all its points were moved on a sphere of radius ρ_i centered
 270 in o_i . This allows us to use a Markov's inequality kind of argument and conclude that there must
 271 be $\Omega_\varepsilon(|O_i^*|)$ points in $O_i^* \cap B(o_i, (1 + \varepsilon)\rho_i)$. However, we have no guarantee that there is any
 272 point at all in $O_i^* \cap B(o_i, (1 - \varepsilon)\rho_i)$. Indeed, all points in O_i^* might lie on $\partial B(o_i, \rho_i)$. The fact
 273 that potentially all our candidate centers q are at distance at least ρ_i from o_i yields (by Lemma 1)
 274 $\text{cost}(O_i^*, q) \geq 2\text{cost}(O_i^*, o_i)$, which causes the zero-degree term in $\xi^2 - (2 + 2/p)\xi - (3 + 2/p) = 0$
 275 from Kanungo et al. [2004] to become a 4 in our analysis.

276 **Improving MSLS by taking averages.** First, we notice that, in order to achieve $(9 + \varepsilon)$ -
 277 approximation we need to set $p = \Theta(1/\varepsilon)$. The main hurdle to achieve a $(9 + \varepsilon)$ -approximation is
 278 that we need to replace the q_i in MSLS with a better approximation of o_i . We design a subroutine
 279 that computes, with constant probability, an ε -approximation \hat{o}_i of o_i (namely, $\text{cost}(O_i^*, \hat{o}_i) \leq$
 280 $(1 + \varepsilon)\text{cost}(O_i^*, o_i)$). The key idea is that, if sample uniformly $O(1/\varepsilon)$ points from O_i^* and define
 281 \hat{o}_i to be the average of our samples then $\text{cost}(O_i^*, \hat{o}_i) \leq (1 + \varepsilon)\text{cost}(O_i^*, o_i)$

282 Though, we do not know O_i^* , so sampling uniformly from it is non-trivial. To achieve that, for each
 283 q_i we identify a set N of nice candidate points in P such that a $\text{poly}(\varepsilon)/k$ fraction of them are from
 284 O_i^* . We sample $O(1/\varepsilon)$ points uniformly from N and thus with probability $(\varepsilon/k)^{O(1/\varepsilon)}$ we sample
 285 only points from O_i^* . Thus far, we sampled $O(1/\varepsilon)$ points uniformly from $N \cap O_i^*$. What about
 286 the points in $O_i^* \setminus N$? We can define N so that all points in $O_i^* \setminus N$ are either very close to some
 287 of the $(q_j)_j$ or they are very far from q_i . The points that are very close to points $(q_j)_j$ are easy to
 288 treat. Indeed, we can approximately locate them and we just need to guess their mass, which is
 289 matters only when $\geq \text{poly}(\varepsilon)\text{ALG}$, and so we pay only a $\log^{O(1/\varepsilon)}(1/\varepsilon)$ multiplicative overhead to
 290 guess the mass close to q_j for $j = 1 \dots p = \Theta(1/\varepsilon)$. As for a point f that is very far from q_i (say,
 291 $\|f - q_i\| \gg \rho_i$) we notice that, although f 's contribution to $\text{cost}(O_i^*, o_i)$ may be large, we have
 292 $\text{cost}(f, o) \approx \text{cost}(f, o_i)$ for each $o \in B(q_i, \rho_i) \subseteq B(o_i, (2 + \varepsilon)\rho_i)$ assuming $q_i \in \text{core}(o_i)$.

293 5 Experiments

294 In this section, we show that our new algorithm using multi-swap local search can be employed to
 295 design an efficient seeding algorithm for Lloyd's which outperforms both the classical k -means++
 296 seeding and the single-swap local search from Lattanzi and Sohler [2019].

297 **Algorithms.** The multi-swap local search algorithm that we analysed above performs very well
 298 in terms of solution quality. This empirically verifies the improved approximation factor of our
 299 algorithm, compared to the single-swap local search of Lattanzi and Sohler [2019].

300 Motivated by practical considerations, we heuristically adapt our algorithm to make it very competitive
 301 with SSLS in terms of running time and still remain very close, in terms of solution quality, to the
 302 theoretically superior algorithm that we analyzed. The adaptation of our algorithm replaces the phase
 303 where it selects the p centers to swap-out by performing an exhaustive search over $\binom{k+p}{p}$ subsets of

³The complexity in Theorem 12 can be improved by applying the same preprocessing techniques using coresets and dimensionality reduction, similar to what can be used to speed up the approach of Kanungo et al. [2004]. Our algorithm hence becomes asymptotically faster.

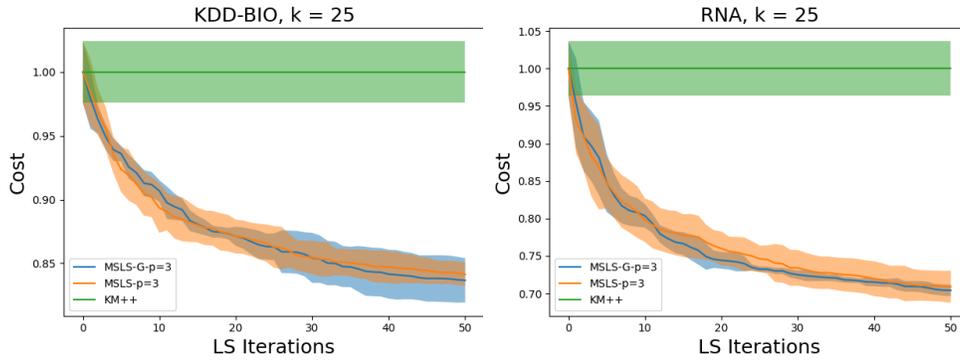


Figure 1: Comparison between MSLS and MSLS-G, for $p = 3$, for $k = 25$, on the datasets KDD-BIO and RNA. The y axis shows the solution cost divided by the means solution cost of KM++.

304 centers. Instead, we use an efficient heuristic procedure for selecting the p centers to swap-out, by
 305 greedily selecting one by one the centers to swap-out. Specifically, we select the first center to be the
 306 cheapest one to remove (namely, the one that increases the cost by the least amount once the points
 307 in its cluster are reassigned to the remaining centers). Then, we update all costs and select the next
 308 center iteratively. After p repetitions we are done. We perform an experimental evaluation of the
 309 “greedy” variant of our algorithm compared to the theoretically-sound algorithm from Section 3 and
 310 show that employing the greedy heuristic does not measurably impact performance.

311 The four algorithms that we evaluate are the following: 1) **KM++**: The k -means++ from Arthur
 312 and Vassilvitskii [2007], 2) **SSLS**: The Single-swap local search method from Lattanzi and Sohler
 313 [2019], 3) **MSLS**: The multi-swap local search from Section 3, and 4) **MSLS-G**: The greedy variant
 314 of multi-swap local search as described above.

315 We use $\text{MSLS-G-}p = x$ and $\text{MSLS-}p = x$, to denote MSLS-G and MSLS with $p = x$, respectively.
 316 Notice that $\text{MSLS-G-}p = 1$ is exactly SSLS. Our experimental evaluation explores the effect of
 317 p -swap LS, for $p > 1$, in terms of solution cost and running time.

318 **Datasets.** We consider the three datasets used in Lattanzi and Sohler [2019] to evaluate the perfor-
 319 mance of SSLS: 1) KDD-PHY – 100, 000 points with 78 features representing a quantum physic task
 320 kdd [2004], 2) RNA - 488, 565 points with 8 features representing RNA input sequence pairs Uzilov
 321 et al. [2006], and 3) KDD-BIO – 145, 751 points with 74 features measuring the match between a
 322 protein and a native sequence kdd [2004]. We discuss the results for two of our datasets, namely
 323 KDD-BIO and RNA. We defer the results on KDD-PHY to the appendix and note that the results are
 324 very similar to the results on RNA.

325 We performed a preprocessing step to clean-up the datasets. We observed that the standard deviation
 326 of some features was disproportionately high. This causes all costs being concentrated in few
 327 dimensions making the problem, in some sense, lower-dimensional. Thus, we apply min-max scaling
 328 to all datasets and observed that this causes all our features’ standard deviations to be comparable.

329 **Experimental setting.** All our code is written in Python. The code will be made available upon
 330 publication of this work. We did not make use of parallelization techniques. To run our experiments,
 331 we used a personal computer with 8 cores, a 1.8 Ghz processor, and 15.9 GiB of main memory
 332 We run all experiments 5 times and report the mean and standard deviation in our plots. All our
 333 plots report the progression of the cost either w.r.t local search steps, or Lloyd’s iterations. We run
 334 experiments on all our datasets for $k = 10, 25, 50$. The main body of the paper reports the results for
 335 $k = 25$, while the rest can be found in the appendix. We note that the conclusions of the experiments
 336 for $k = 10, 50$ are similar to those of $k = 25$.

337 **Removing centers greedily.** We first we compare MSLS-G with MSLS. To perform our experi-
 338 ment, we initialize $k = 25$ centers using k -means++ and then run 50 iterations of local search for both
 339 algorithms, for $p = 3$ swaps. Due to the higher running of the MSLS we perform this experiments on
 340 1% uniform sample of each of our datasets. We find out that the performance of the two algorithms is
 341 comparable on all our instances, while they both perform roughly 15%-27% at convergence. Figure 1
 342 shows the aggregate results, over 5 repetitions of our experiment.

343 It may happen that MSLS, which considers all possible swaps of size p at each LS iteration, performs
 344 worse than MSLS-G as a sub-optimal swap at intermediate iterations may still lead to a better local

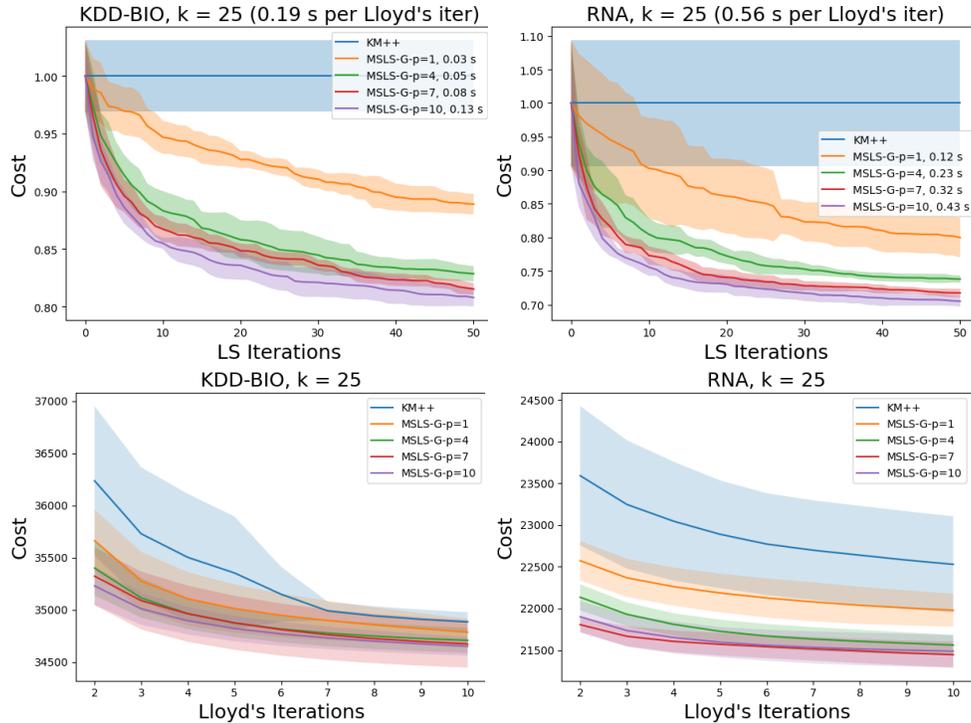


Figure 2: The first row compares the cost of MSLS-G, for $p \in \{1, 4, 7, 10\}$, divided by the mean cost of KM++ at each LS step, for $k = 25$. The legend reports also the running time of MSLS-G per LS step (in seconds). The second row compares the cost after each of the 10 iterations of Lloyd with seeding from MSLS-G, for $p \in \{1, 4, 7, 10\}$ and 15 local search steps and KM++, for $k = 25$.

345 optimum by coincidence. Given that MSLS-G performs very comparably to MSLS, while it is much
 346 faster in practice, we use MSLS-G for the rest of our experiments where we compare to baselines.
 347 This allows us to consider higher values of p , without compromising much the running time.

348 **Results: Evaluating the quality and performance of the algorithms.** In our first experiment we
 349 run KM++ followed by 50 iterations of MSLS-G with $p = 1, 4, 7, 10$ and plot the relative cost w.r.t.
 350 KM++ at each iteration, for $k = 25$. The first row of Figure 2 plots the results. Our experiment shows
 351 that, after 50 iterations MSLS-G for $p = 4, 7, 10$ achieves improvements of roughly 10% compared
 352 to MSLS-G- $p = 1$ and of the order of 20% – 30% compared to KM++. We also report the time per
 353 iteration that each algorithm takes. For comparison, we report the running time of a single iteration of
 354 Lloyd’s next to the dataset’s name. It is important to notice that, although MSLS-G- $p = 1$ is faster,
 355 running more iterations MSLS-G- $p = 1$ is not sufficient to compete with MSLS-G when $p > 1$.

356 **Results: Evaluating the quality after postprocessing using Lloyd.** In our second experiment,
 357 we use KM++ and MSLS-G as a seeding algorithm for Lloyd’s and measure how much of the
 358 performance improvement measured in the first experiment is retained after running Lloyd’s. First,
 359 we initialize our centers using KM++ and the run 15 iterations of MSLS-G for $p = 1, 4, 7$. We
 360 measure the cost achieved by running 10 iterations of Lloyd’s starting from the solutions found by
 361 MSLS-G as well as KM++. In Figure 2 (second row) we plot the results. Notice that, according to
 362 the running times from the first experiment, 15 iterations of MSLS-G take less than 10
 363 iterations of Lloyd’s for $p = 4, 7$ (and also for $p = 10$, except on RNA). We observe that MSLS-G
 364 for $p > 1$ performs at least as good as SLS from Lattanzi and Sohler [2019] and in some cases
 365 maintains non-trivial improvements.

366 Conclusion and Future Directions

367 We present a new algorithm for the k -means problem and we show that it outperforms theoretically
 368 and experimentally state-of-the-art practical algorithms with provable guarantees in terms of solution
 369 quality. A very interesting open question is to improve our local search procedure by avoiding the
 370 exhaustive search over all possible size- p subsets of centers to swap out, concretely an algorithm with
 371 running time $\tilde{O}(2^{poly(1/\varepsilon)} ndk)$.

372 **References**

- 373 Kdd cup. 2004. URL <http://osmot.cs.cornell.edu/kddcup/datasets.html>.
- 374 Sara Ahmadian, Ashkan Norouzi-Fard, Ola Svensson, and Justin Ward. Better guarantees for k-
375 means and Euclidean k-median by primal-dual algorithms. *SIAM Journal on Computing*, 49(4):
376 FOCS17–97–FOCS17–156, 2019.
- 377 David Arthur and Sergei Vassilvitskii. K-means++ the advantages of careful seeding. In *Proceedings*
378 *of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035, 2007.
- 379 David Arthur and Sergei Vassilvitskii. Worst-case and smoothed analysis of the ICP algorithm, with
380 an application to the k-means method. *SIAM J. Comput.*, 39(2):766–782, 2009. doi: 10.1137/
381 070683921. URL <http://dx.doi.org/10.1137/070683921>.
- 382 Vijay Arya, Naveen Garg, Rohit Khandekar, Adam Meyerson, Kamesh Munagala, and Vinayaka
383 Pandit. Local search heuristics for k-median and facility location problems. *SIAM J. Comput.*,
384 33(3):544–562, 2004. doi: 10.1137/S0097539702416402. URL <https://doi.org/10.1137/S0097539702416402>.
- 386 Pranjali Awasthi, Moses Charikar, Ravishankar Krishnaswamy, and Ali Kemal Sinop. The hardness of
387 approximation of euclidean k-means. In Lars Arge and János Pach, editors, *31st International Sym-*
388 *posium on Computational Geometry, SoCG 2015, June 22-25, 2015, Eindhoven, The Netherlands*,
389 volume 34 of *LIPICs*, pages 754–767. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2015. doi:
390 10.4230/LIPICs.SOCG.2015.754. URL <https://doi.org/10.4230/LIPICs.SOCG.2015.754>.
- 391 Sayan Bandyapadhyay and Kasturi Varadarajan. On variants of k-means clustering. In *32nd*
392 *International Symposium on Computational Geometry (SoCG 2016)*. Schloss Dagstuhl-Leibniz-
393 Zentrum fuer Informatik, 2016.
- 394 Luca Becchetti, Marc Bury, Vincent Cohen-Addad, Fabrizio Grandoni, and Chris Schwiegelshohn.
395 Oblivious dimension reduction for k -means: beyond subspaces and the johnson-lindenstrauss
396 lemma. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing,*
397 *STOC 2019, Phoenix, AZ, USA, June 23-26, 2019*, pages 1039–1050, 2019. doi: 10.1145/3313276.
398 3316318. URL <https://doi.org/10.1145/3313276.3316318>.
- 399 Anup Bhattacharya, Dishant Goyal, Ragesh Jaiswal, and Amit Kumar. On sampling based algorithms
400 for k-means. In Nitin Saxena and Sunil Simon, editors, *40th IARCS Annual Conference on*
401 *Foundations of Software Technology and Theoretical Computer Science, FSTTCS 2020, December*
402 *14-18, 2020, BITS Pilani, K K Birla Goa Campus, Goa, India (Virtual Conference)*, volume 182
403 of *LIPICs*, pages 13:1–13:17. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2020. doi: 10.
404 4230/LIPICs.FSTTCS.2020.13. URL <https://doi.org/10.4230/LIPICs.FSTTCS.2020.13>.
- 405 Moses Charikar and Sudipto Guha. Improved combinatorial algorithms for facility location problems.
406 *SIAM J. Comput.*, 34(4):803–824, 2005. doi: 10.1137/S0097539701398594. URL <https://doi.org/10.1137/S0097539701398594>.
- 408 Davin Choo, Christoph Grunau, Julian Portmann, and Vaclav Rozhon. k-means++: few more steps
409 yield constant approximation. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th*
410 *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*
411 *Research*, pages 1909–1917. PMLR, 2020. URL <https://proceedings.mlr.press/v119/choo20a.html>.
- 413 Vincent Cohen-Addad. A fast approximation scheme for low-dimensional k -means. In Artur Czumaj,
414 editor, *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms,*
415 *SODA 2018, New Orleans, LA, USA, January 7-10, 2018*, pages 430–440. SIAM, 2018. doi:
416 10.1137/1.9781611975031.29. URL <https://doi.org/10.1137/1.9781611975031.29>.
- 417 Vincent Cohen-Addad and Karthik C. S. Inapproximability of clustering in l_p metrics. In David
418 Zuckerman, editor, *60th IEEE Annual Symposium on Foundations of Computer Science, FOCS*
419 *2019, Baltimore, Maryland, USA, November 9-12, 2019*, pages 519–539. IEEE Computer Society,
420 2019. doi: 10.1109/FOCS.2019.00040. URL <https://doi.org/10.1109/FOCS.2019.00040>.

- 421 Vincent Cohen-Addad and Claire Mathieu. Effectiveness of local search for geometric optimization.
 422 In *31st International Symposium on Computational Geometry, SoCG 2015, June 22-25, 2015,*
 423 *Eindhoven, The Netherlands*, pages 329–343, 2015. doi: 10.4230/LIPIcs.SOCG.2015.329. URL
 424 <http://dx.doi.org/10.4230/LIPIcs.SOCG.2015.329>.
- 425 Vincent Cohen-Addad and Chris Schwiegelshohn. On the local structure of stable clustering instances.
 426 In Chris Umans, editor, *58th IEEE Annual Symposium on Foundations of Computer Science, FOCS*
 427 *2017, Berkeley, CA, USA, October 15-17, 2017*, pages 49–60. IEEE Computer Society, 2017. doi:
 428 10.1109/FOCS.2017.14. URL <https://doi.org/10.1109/FOCS.2017.14>.
- 429 Vincent Cohen-Addad, Philip N. Klein, and Claire Mathieu. Local search yields approximation
 430 schemes for k-means and k-median in euclidean and minor-free metrics. *SIAM J. Comput.*, 48(2):
 431 644–667, 2019. doi: 10.1137/17M112717X. URL <https://doi.org/10.1137/17M112717X>.
- 432 Vincent Cohen-Addad, David Saulpic, and Chris Schwiegelshohn. Improved coresets and sublinear
 433 algorithms for power means in euclidean spaces. *Advances in Neural Information Processing*
 434 *Systems*, 34:21085–21098, 2021.
- 435 Vincent Cohen-Addad, Hossein Esfandiari, Vahab S. Mirrokni, and Shyam Narayanan. Improved
 436 approximations for euclidean k -means and k -median, via nested quasi-independent sets. In Stefano
 437 Leonardi and Anupam Gupta, editors, *STOC '22: 54th Annual ACM SIGACT Symposium on*
 438 *Theory of Computing, Rome, Italy, June 20 - 24, 2022*, pages 1621–1628. ACM, 2022a. doi:
 439 10.1145/3519935.3520011. URL <https://doi.org/10.1145/3519935.3520011>.
- 440 Vincent Cohen-Addad, Kasper Green Larsen, David Saulpic, and Chris Schwiegelshohn. Towards
 441 optimal lower bounds for k-median and k-means coresets. In Stefano Leonardi and Anupam
 442 Gupta, editors, *STOC '22: 54th Annual ACM SIGACT Symposium on Theory of Computing, Rome,*
 443 *Italy, June 20 - 24, 2022*, pages 1038–1051. ACM, 2022b. doi: 10.1145/3519935.3519946. URL
 444 <https://doi.org/10.1145/3519935.3519946>.
- 445 Vincent Cohen-Addad, Kasper Green Larsen, David Saulpic, Chris Schwiegelshohn,
 446 and Omar Ali Sheikh-Omar. Improved coresets for euclidean k-means. In
 447 *NeurIPS*, 2022c. URL [http://papers.nips.cc/paper_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/120c9ab5c58ba0fa9dd3a22ace1de245-Abstract-Conference.html)
 448 [120c9ab5c58ba0fa9dd3a22ace1de245-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/120c9ab5c58ba0fa9dd3a22ace1de245-Abstract-Conference.html).
- 449 Vincent Cohen-Addad, Euiwoong Lee, and Karthik C. S. Johnson coverage hypothesis: Inapproxima-
 450 bility of k -means and k -median in ℓ_p metrics. In *Proceedings of the 2022 ACM-SIAM Symposium*
 451 *on Discrete Algorithms, SODA 2022*. SIAM, 2022d.
- 452 D. Feldman and M. Langberg. A unified framework for approximating and clustering data. In *STOC*,
 453 pages 569–578, 2011.
- 454 Zachary Friggstad, Mohsen Rezapour, and Mohammad R. Salavatipour. Local search yields a
 455 PTAS for k-means in doubling metrics. *SIAM J. Comput.*, 48(2):452–480, 2019. doi: 10.1137/
 456 17M1127181. URL <https://doi.org/10.1137/17M1127181>.
- 457 Fabrizio Grandoni, Rafail Ostrovsky, Yuval Rabani, Leonard J. Schulman, and Rakesh Venkat.
 458 A refined approximation for euclidean k-means. *Inf. Process. Lett.*, 176:106251, 2022. doi:
 459 10.1016/j.ipl.2022.106251. URL <https://doi.org/10.1016/j.ipl.2022.106251>.
- 460 Christoph Grunau, Ahmet Alper Özüdogru, Václav Rozhon, and Jakub Tetek. A nearly tight analysis
 461 of greedy k-means++. In Nikhil Bansal and Viswanath Nagarajan, editors, *Proceedings of the*
 462 *2023 ACM-SIAM Symposium on Discrete Algorithms, SODA 2023, Florence, Italy, January 22-*
 463 *25, 2023*, pages 1012–1070. SIAM, 2023. doi: 10.1137/1.9781611977554.ch39. URL [https://](https://doi.org/10.1137/1.9781611977554.ch39)
 464 doi.org/10.1137/1.9781611977554.ch39.
- 465 Anupam Gupta and Kanat Tangwongsan. Simpler analyses of local search algorithms for facility
 466 location. *CoRR*, abs/0809.2554, 2008. URL <http://arxiv.org/abs/0809.2554>.
- 467 Venkatesan Guruswami and Piotr Indyk. Embeddings and non-approximability of geometric problems.
 468 In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms, January*
 469 *12-14, 2003, Baltimore, Maryland, USA*, pages 537–538, 2003. URL [http://dl.acm.org/](http://dl.acm.org/citation.cfm?id=644108.644198)
 470 [citation.cfm?id=644108.644198](http://dl.acm.org/citation.cfm?id=644108.644198).

- 471 Mary Inaba, Naoki Katoh, and Hiroshi Imai. Applications of weighted voronoi diagrams and
472 randomization to variance-based k-clustering. In *Proceedings of the tenth annual symposium on*
473 *Computational geometry*, pages 332–339, 1994.
- 474 Ragesh Jaiswal, Amit Kumar, and Sandeep Sen. A simple D 2-sampling based PTAS for k-means and
475 other clustering problems. *Algorithmica*, 70(1):22–46, 2014. doi: 10.1007/s00453-013-9833-9.
476 URL <https://doi.org/10.1007/s00453-013-9833-9>.
- 477 Tapas Kanungo, David M. Mount, Nathan S. Netanyahu, Christine D. Piatko, Ruth Silverman,
478 and Angela Y. Wu. A local search approximation algorithm for k-means clustering. *Com-*
479 *putational Geometry*, 28(2):89–112, 2004. ISSN 0925-7721. doi: [https://doi.org/10.1016/](https://doi.org/10.1016/j.comgeo.2004.03.003)
480 [j.comgeo.2004.03.003](https://doi.org/10.1016/j.comgeo.2004.03.003). URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0925772104000215)
481 [S0925772104000215](https://www.sciencedirect.com/science/article/pii/S0925772104000215). Special Issue on the 18th Annual Symposium on Computational Geometry
482 - SoCG2002.
- 483 Madhukar R. Korupolu, C. Greg Plaxton, and Rajmohan Rajaraman. Analysis of a local search
484 heuristic for facility location problems. *J. Algorithms*, 37(1):146–188, 2000. doi: 10.1006/jagm.
485 2000.1100. URL <http://dx.doi.org/10.1006/jagm.2000.1100>.
- 486 Amit Kumar, Yogish Sabharwal, and Sandeep Sen. Linear-time approximation schemes for clustering
487 problems in any dimensions. *J. ACM*, 57(2):5:1–5:32, 2010. doi: 10.1145/1667053.1667054. URL
488 <https://doi.org/10.1145/1667053.1667054>.
- 489 Silvio Lattanzi and Christian Sohler. A better k-means++ algorithm via local search. In Kamalika
490 Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on*
491 *Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 3662–3671.
492 PMLR, 2019. URL <https://proceedings.mlr.press/v97/lattanzi19a.html>.
- 493 Euiwoong Lee, Melanie Schmidt, and John Wright. Improved and simplified inapproximability
494 for k-means. *Inf. Process. Lett.*, 120:40–43, 2017. doi: 10.1016/j.ipl.2016.11.009. URL <https://doi.org/10.1016/j.ipl.2016.11.009>.
- 495
- 496 SP Lloyd. Least square quantization in pcm. bell telephone laboratories paper. published in journal
497 much later: Lloyd, sp: Least squares quantization in pcm. *IEEE Trans. Inform. Theor.*(1957/1982),
498 18, 1957.
- 499 Konstantin Makarychev, Yury Makarychev, Maxim Sviridenko, and Justin Ward. A bi-criteria
500 approximation algorithm for k-means. In Klaus Jansen, Claire Mathieu, José D. P. Rolim, and Chris
501 Umans, editors, *Approximation, Randomization, and Combinatorial Optimization. Algorithms and*
502 *Techniques, APPROX/RANDOM 2016, September 7-9, 2016, Paris, France*, volume 60 of *LIPIcs*,
503 pages 14:1–14:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2016. doi: 10.4230/LIPIcs.
504 APPROX-RANDOM.2016.14. URL [https://doi.org/10.4230/LIPIcs.APPROX-RANDOM.](https://doi.org/10.4230/LIPIcs.APPROX-RANDOM.2016.14)
505 2016.14.
- 506 Konstantin Makarychev, Yury Makarychev, and Ilya P. Razenshteyn. Performance of johnson-
507 lindenstrauss transform for k-means and k-medians clustering. In Moses Charikar and Edith Cohen,
508 editors, *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing, STOC*
509 *2019, Phoenix, AZ, USA, June 23-26, 2019*, pages 1027–1038. ACM, 2019. doi: 10.1145/3313276.
510 3316350. URL <https://doi.org/10.1145/3313276.3316350>.
- 511 Jirí Matousek. On approximate geometric k-clustering. *Discrete & Computational Geometry*,
512 24(1):61–84, 2000. doi: 10.1007/s004540010019. URL [http://dx.doi.org/10.1007/](http://dx.doi.org/10.1007/s004540010019)
513 [s004540010019](http://dx.doi.org/10.1007/s004540010019).
- 514 Andrew V Uzilov, Joshua M Keegan, and David H Mathews. Detection of non-coding rnas on the
515 basis of predicted secondary structure formation free energy change. *BMC bioinformatics*, 7(1):
516 1–30, 2006.

517 **Supplementary Material**

518 **Proofs from Section 3**

519 **Lemma 8.** $\sum_{p \in P} \text{cost}(p, \mathcal{C}[\mathcal{O}^*[p]]) \leq 2OPT + ALG + 2\sqrt{ALG}\sqrt{OPT}$.

Proof.

$$\begin{aligned}
& \sum_{p \in P} \text{cost}(p, \mathcal{C}[\mathcal{O}^*[p]]) = \\
& \sum_{o_i \in \mathcal{O}^*} \sum_{p \in \mathcal{O}_i^*} \text{cost}(p, \mathcal{C}[o_i]) = \\
& \sum_{o_i \in \mathcal{O}^*} |O_i^*| \cdot \text{cost}(o_i, \mathcal{C}[o_i]) + \text{cost}(O_i^*, o_i) = \\
& OPT + \sum_{p \in P} \text{cost}(\mathcal{O}^*[p], \mathcal{C}[\mathcal{O}^*[p]]) \leq \\
& OPT + \sum_{p \in P} \text{cost}(\mathcal{O}^*[p], \mathcal{C}[p]) \leq \\
& OPT + \sum_{p \in P} (\|\mathcal{O}^*[p] - p\| + \|p - \mathcal{C}[p]\|)^2 = \\
& 2OPT + ALG + 2 \sum_{p \in P} \|\mathcal{O}^*[p], p\| \cdot \|p, \mathcal{C}[p]\| \leq 2OPT + ALG + 2\sqrt{ALG}\sqrt{OPT}.
\end{aligned}$$

520 The second equality is due to Lemma 1 and the last inequality is due to Cauchy-Schwarz. \square

521 **Proofs from Section 4**

522 In this section, we prove the following.

523 **Theorem 12.** *Given a set of n points in \mathbb{R}^d with aspect ratio Δ , there exists an algorithm that*
524 *computes a $9 + \varepsilon$ -approximation to k -means in time $ndk^{O(\varepsilon^{-2})} \log^{O(\varepsilon^{-1})}(\Delta) \cdot 2^{-\text{poly}(\varepsilon^{-1})}$.*

525 We start with a key lemma showing that a sample of size $O(1/\varepsilon)$ is enough to approximate 1-mean.

526 **Lemma 13** (Form Inaba et al. [1994]). *Given an instance $P \subseteq \mathbb{R}^d$, sample $m = 1/(\varepsilon\delta)$ points*
527 *uniformly at random from P and denote the set of samples with S . Then $\text{cost}(P, \mu(S)) \leq (1 +$
528 $\varepsilon) \text{cost}(P, \mu(P))$ with probability at least $1 - \delta$.*

529 *Proof.* We want to prove that with probability $1 - \delta$ we have $\|\mu(S) - \mu(P)\|^2 \leq$
530 $\varepsilon \text{cost}(P, \mu(P)) / |P|$. Then, applying Lemma 1 gives the desired result. First, we notice that
531 $\mu(P)$ is an unbiased estimator of $\mu(P)$, namely $E[\mu(S)] = \mu(P)$. Then, we have

$$E[\|\mu(S) - \mu(P)\|^2] = \frac{1}{m} \sum_{i=1}^{|S|} E[\|s_i - \mu(P)\|^2] = \frac{\text{cost}(P, \mu(P))}{m \cdot |P|}$$

532 where s_i are uniform independent samples from P . Applying Markov's inequality concludes the
533 proof. \square

534 The algorithm that verifies Theorem 12 is very similar to the MSLS algorithm from Section 3 and
535 we use the same notation to describe it. The intuition is that in MSLS we sample $\mathcal{Q} = \{q_1 \dots q_p\}$
536 hoping that $q_i \in \text{core}(o_i)$ for each i ; here we refine q_i to a better approximation \hat{o}_i of o_i and swap
537 the points $(\hat{o}_i)_i$ rather than $(q_i)_i$. Our points \hat{o}_i are generated taking the average of some sampled
538 point, thus we possibly have $\hat{o}_i \notin P$ while, on the other hand, $q_i \in P$.

539 **A $(9 + \varepsilon)$ -approximation MSLS algorithm.** First, we initialize our set of centers using k -means++.
540 Then, we run $ndk^{O(\varepsilon^{-2})} \cdot 2^{\text{poly}(\varepsilon^{-1})}$ local search steps, where a local search step works as follows.
541 Set $p = \Theta(\varepsilon^{-1})$. We D^2 -sample a set $\mathcal{Q} = \{q_1 \dots q_p\}$ of points from P (without updating costs).
542 Then, we iterate over all possible sets $Out = \{c_1 \dots c_p\}$ of p distinct elements in $\mathcal{C} \cup \mathcal{Q}$. We define
543 the set of *temporary centers* $\mathcal{T} = (\mathcal{C} \cup \mathcal{Q}) \setminus Out$ and run a subroutine $\text{APX-CENTERS}(\mathcal{T})$ which
544 returns a list of $\text{poly}(\varepsilon^{-1}) \cdot \log^{O(\varepsilon^{-1})}(\Delta)$ size- s sets $\widehat{In} = \{\hat{o}_1 \dots \hat{o}_s\}$ (where $s = |\mathcal{Q} \setminus Out|$). We
545 select the set \widehat{In} in this list such that the swap $(\widehat{In}, Out \setminus \mathcal{Q})$ yields the maximum cost reduction.
546 Then we select the set Out that maximizes the cost reduction obtained in this way. If $(\widehat{In}, Out \setminus \mathcal{Q})$
547 actually reduces the cost then we perform that swap.

548 **A subroutine to approximate optimal centers.** Here we describe the subroutine
549 $\text{APX-CENTERS}(\mathcal{T})$. Let $\mathcal{Q} \setminus Out = \{q_1 \dots q_s\}$. Recall that $s \leq p = O(\varepsilon^{-1})$. This subroutine out-
550 puts a list of $2^{\text{poly}(\varepsilon^{-1})} \cdot \log^{O(\varepsilon^{-1})}(\Delta)$ size- s sets $\widehat{In} = \{\hat{o}_1 \dots \hat{o}_s\}$. Here we describe how to find a list
551 of $2^{\text{poly}(\varepsilon^{-1})} \cdot \log(\Delta)$ values for \hat{o}_1 . The same will apply for $\hat{o}_2 \dots \hat{o}_s$ and taking the Cartesian product
552 yields a list of $2^{\text{poly}(\varepsilon^{-1})} \cdot \log^{O(\varepsilon^{-1})}(\Delta)$ size- s sets. Assume wlog that the pairwise distances between
553 points in P lie in $[1, \Delta]$. We iterate over all possible values of $\rho_1 \in \{1, (1 + \varepsilon) \dots (1 + \varepsilon)^{\lceil \log_{1+\varepsilon} \Delta \rceil}\}$.
554 We partition P in three sets: the set of *far points* $F = \{x \in P \mid \text{cost}(x, q_1) > \rho_1^2 / \varepsilon^3\}$, the set of
555 *close points* $C = \{x \in P \setminus F \mid \text{cost}(x, \mathcal{T}) \leq \varepsilon^3 \rho_1^2\}$ and the set of *nice points* $N = P \setminus (C \cup F)$.
556 Then, we sample uniformly from N a set S of size $\Theta(\varepsilon^{-1})$. For each $(s + 1)$ -tuple of coefficients
557 $\alpha_0, \alpha_1 \dots \alpha_s \in \left\{1, (1 - \varepsilon), (1 - \varepsilon)^2, \dots, (1 - \varepsilon)^{\lceil \log_{1-\varepsilon}(\varepsilon^7) \rceil}\right\} \cup \{0\}$ we output the candidate solution
558 given by the convex combination

$$\hat{o}_1 = \hat{o}_1(\alpha_0 \dots \alpha_s) = \frac{\alpha_0 \mu(S) + \sum_{i=1}^s \alpha_i q_i}{\sum_{i=0}^s \alpha_i} \quad (3)$$

559 so, for each value of ρ_1 , we output $2^{\text{poly}(\varepsilon^{-1})}$ values for \hat{o}_1 . Hence, $2^{\text{poly}(\varepsilon^{-1})} \cdot \log(\Delta)$ values in total.

560 Analysis

561 The key insight in the analysis of the MSLS algorithm from Section 3 was that every q_i was a proxy
562 for o_i because $q_i \in \text{core}(o_i)$, and thus q_i provided a good center for O_i^* . In the analysis of this
563 improved version of MSLS we replace q_i with \hat{o}_i which makes a better center for O_i^* . Formally, fixed
564 Out , we say that a point \hat{o}_i is a *perfect approximation* of o_i when $\text{cost}(O_i^*, (\mathcal{C} \cup \{\hat{o}_i\}) \setminus Out) \leq$
565 $(1 + \varepsilon)\text{OPT}_i + \varepsilon\text{OPT}/k$. We define $\widetilde{\mathcal{O}}$ and $\widetilde{\mathcal{C}}$ as in Section 3, except that we replace δ with ε (which
566 here is not assumed to be a constant). Likewise, we build the set \mathcal{S} of ideal multi-swaps as in Section 3.
567 Recall that we say that a multi-swap (In, Out) is *strongly improving* if $\text{cost}(P, (\mathcal{C} \cup In) \setminus Out) \leq$
568 $(1 - \varepsilon/k) \cdot \text{cost}(P, \mathcal{C})$. Let $In = \{o_1 \dots o_s\} \subseteq \widetilde{\mathcal{O}}$ and $Out = \{c_1 \dots c_s\} \subseteq \widetilde{\mathcal{C}}$, we overload
569 the definition from Section 3 and say that the ideal multi-swap (In, Out) is *good* if for every
570 $\widehat{In} = \{\hat{o}_1 \dots \hat{o}_s\}$ such that each \hat{o}_i is a perfect approximation of o_i for each $i = 1 \dots s$ the swap
571 (\widehat{In}, Out) is strongly improving. We call an ideal swap *bad* otherwise. As in Section 3, we define
572 the *core* of an optimal center; once again we replace δ with ε , which is no longer constant. The two
573 following lemmas are our stepping stones towards Theorem 12.

574 **Lemma 14.** *If $\text{ALG}/\text{OPT} > 9 + O(\varepsilon)$ then, with probability $k^{-O(\varepsilon^{-1})} \cdot 2^{-\text{poly}(\varepsilon^{-1})}$, there exists*
575 *$Out \subseteq \mathcal{C} \cup \mathcal{Q}$ such that:*

- 576 (i) *If $\mathcal{Q} \setminus Out = \{q_1 \dots q_s\}$ then $q_1 \in \text{core}(o_1) \dots q_s \in \text{core}(o_s)$ for some $o_1 \dots o_s \in \mathcal{O}^*$*
- 577 (ii) *If we define $In = \{o_1 \dots o_s\}$ then $(In, Out \setminus \mathcal{Q})$ is a good ideal swap.*

578 **Lemma 15.** *If (i) from Lemma 14 holds, then with probability $k^{-O(\varepsilon^{-2})} \cdot 2^{-\text{poly}(\varepsilon^{-1})}$, the list*
579 *returned by APX-CENTERS contains $\widehat{In} = \{\hat{o}_1 \dots \hat{o}_s\}$ such that \hat{o}_i is a perfect approximation of o_i*
580 *for each $i = 1 \dots s$.*

581 *Proof of Theorem 12.* Here we prove that our improved MSLS algorithm achieves a $(9 + O(\varepsilon))$ -
582 approximation, which is equivalent to Theorem 12 up to rescaling ε . Combining Lemma 14 and

583 Lemma 15 we obtain that, as long as $\text{ALG}/\text{OPT} > 9 + O(\varepsilon)$, with probability at least $k^{-O(\varepsilon^{-2})}$.
584 $2^{-\text{poly}(\varepsilon^{-1})}$, the list returned by APX-CENTERS contains $\widehat{In} = \{\hat{o}_1 \dots \hat{o}_s\}$ such that $(\widehat{In}, \text{Out} \setminus \mathcal{Q})$
585 is strongly improving. If this happens, we call such a local step *successful*. Now the proof goes
586 exactly as the proof of Theorem 3. Indeed, We show that $k^{O(\varepsilon^{-2})} \cdot 2^{\text{poly}(\varepsilon^{-1})}$ local steps suffice to
587 obtain $\Omega(k \log \log k/\varepsilon)$ successful local steps, and thus to obtain the desired approximation ratio,
588 with constant probability.

589 To prove the running time bound it is sufficient to notice that a local search step can be performed in
590 time $nd \log^{O(\varepsilon^{-1})}(\Delta) \cdot 2^{\text{poly}(\varepsilon^{-1})}$. \square

591 In the rest of this section, we prove Lemma 14 and Lemma 15.

592 **Observation 16.** *If we assume $\delta = \varepsilon$ non-constant in Lemma 2, then performing the computations*
593 *explicitly we obtain $\Pr[q \in \text{core}(O_i^*)] \geq \text{poly}(\varepsilon)$.*

594 In order to prove Lemma 14, we first prove the two lemmas. Lemma 17 is the analogous of Lemma 10
595 and Lemma 18 is the analogous of Lemma 11. Overloading once again the definition from Section 3,
596 we define G as the union of cores of good optimal centers in $\widetilde{\mathcal{O}}$, where an optimal center is defined to
597 be good if at least one of the ideal multi-swaps in \mathcal{S} it belongs to is good (exactly as in Section 3).

598 **Lemma 17.** *If an ideal swap (In, Out) is bad, then we have*

$$\text{cost}(O_{In}^*, \mathcal{C}) \leq (1 + \varepsilon) \text{cost}(O_{In}^*, \mathcal{O}^*) + \text{Reassign}(In, Out) + \varepsilon \text{ALG}/k. \quad (4)$$

599 *Proof.* Let $In = \{o_1 \dots o_s\}$, $\widehat{In} = \{\hat{o}_1 \dots \hat{o}_s\}$ such that \hat{o}_i is a perfect approximation of o_i for each
600 $i = 1 \dots s$. Recall that $O_{In}^* := \bigcup_{i=1}^s O_i^*$, then

$$\text{cost}(O_{In}^*, (\mathcal{C} \cup \widehat{In}) \setminus Out) \leq \sum_{i=1}^s \text{cost}(O_i^*, (\mathcal{C} \cup \{\hat{o}_i\}) \setminus Out) \leq (1 + \varepsilon) \text{cost}(O_{In}^*, \mathcal{O}^*). \quad (5)$$

601 Moreover, $\text{Reassign}(In, Out) = \text{cost}(P \setminus O_{In}^*, \mathcal{C} \setminus Out) - \text{cost}(P \setminus O_{In}^*, \mathcal{C})$ because points
602 in $P \setminus C_{Out}$ are not affected by the swap. Therefore, $\text{cost}(P, (\mathcal{C} \cup \widehat{In}) \setminus Out) \leq (1 +$
603 $\varepsilon) \text{cost}(O_{In}^*, \mathcal{O}^*) + \text{Reassign}(In, Out) + \text{cost}(P \setminus O_{In}^*, \mathcal{C})$. Suppose by contradiction that Equa-
604 tion (4) does not hold, then

$$\begin{aligned} & \text{cost}(P, \mathcal{C}) - \text{cost}(P, (\mathcal{C} \cup \widehat{In}) \setminus Out) = \\ & \text{cost}(P \setminus O_{In}^*, \mathcal{C}) + \text{cost}(O_{In}^*, \mathcal{C}) - \text{cost}(P, (\mathcal{C} \cup \widehat{In}) \setminus Out) \geq \varepsilon \text{ALG}/k. \end{aligned}$$

605 Hence, (\widehat{In}, Out) is strongly improving and this holds for any choice of \widehat{In} , contradiction. \square

606 **Lemma 18.** *If $\text{ALG}/\text{OPT} > 9 + O(\varepsilon)$ then $\text{cost}(G, \mathcal{C}) \geq \text{cost}(P, \mathcal{C}) \cdot \text{poly}(\varepsilon)$. Thus, if we*
607 *D^2 -sample q we have $P[q \in G] \geq \text{poly}(\varepsilon)$.*

608 *Proof.* First, we observe that the combined current cost of all optimal clusters in $\mathcal{O}^* \setminus \widetilde{\mathcal{O}}$ is at most
609 $k \cdot \varepsilon \text{ALG}/k = \varepsilon \text{ALG}$. Now, we prove that the combined current cost of all O_i^* such that o_i is bad is
610 $\leq (1 - 2\varepsilon) \text{ALG}$. Suppose, by contradiction, that it is not the case, then we have:

$$\begin{aligned} (1 - 2\varepsilon) \text{ALG} &< \sum_{\text{Bad } o_i \in \widetilde{\mathcal{O}}} \text{cost}(O_i^*, \mathcal{C}) \leq \sum_{\text{Bad } (In, Out) \in \mathcal{S}} w(In, Out) \cdot \text{cost}(O_{In}^*, \mathcal{C}) \leq \\ & \sum_{\text{Bad } (In, Out)} w(In, Out) \cdot ((1 + \varepsilon) \text{cost}(O_{In}^*, \mathcal{O}^*) + \text{Reassign}(In, Out) + \varepsilon \text{ALG}/k) \leq \\ & (1 + \varepsilon) \text{OPT} + (2 + 2/p) \text{OPT} + (2 + 2/p) \sqrt{\text{ALG}} \sqrt{\text{OPT}} + \varepsilon \text{ALG}. \end{aligned}$$

611 The second and last inequalities make use of Observation 7. The third inequality uses Lemma 17.

612 Setting $\eta^2 = \text{ALG}/\text{OPT}$ we obtain the inequality $\eta^2 - (2 + 2/p \pm O(\varepsilon))\eta - (3 + 2/p \pm O(\varepsilon)) \leq 0$.
613 Hence, we obtain a contradiction in the previous argument as long as $\eta^2 - (2 + 2/p \pm O(\varepsilon))\eta - (3 +$

614 $2/p \pm O(\varepsilon) > 0$, which holds for $p = \Theta(\varepsilon^{-1})$ and $\eta^2 = 9 + O(\varepsilon)$. A contradiction there implies that
615 at least an ε -fraction of the current cost is due to points in $\bigcup_{\text{Good } o_i \in \tilde{\mathcal{O}}} O_i^*$. Thanks to Observation 16,
616 we have $P_{q \sim \text{cost}(q, \mathcal{C})}[q \in \text{core}(O_i^*) \mid q \in O_i^*] \geq \text{poly}(\varepsilon)$. Therefore, we can conclude that the
617 current cost of $G = \bigcup_{\text{Good } o_i \in \tilde{\mathcal{O}}} \text{core}(O_i^*)$ is at least a $\text{poly}(\varepsilon)$ -fraction of the total current cost. \square

618 *Proof of Lemma 14.* Thanks to Lemma 18, we have that $P[q_1 \in G] \geq \text{poly}(\varepsilon)$. Whenever $q_1 \in G$
619 we have that $q_1 \in \text{core}(o_1)$ for some good o_1 . Then, for some $s \leq p$ we can complete o_1 with
620 $o_2 \dots o_s$ such that $In = \{o_1 \dots o_s\}$ belongs to a good swap. Concretely, there exists $Out \subseteq \mathcal{C}$ such
621 that (In, Out) is a good swap. Since $In \subset \tilde{\mathcal{O}}$ we have $\text{cost}(O_i^*, \mathcal{C}) > \varepsilon \text{OPT}/k$ for all $o_i \in In$,
622 which combined with Observation 16 gives that, for each $i = 2 \dots s$, $P[q_i \in \text{core}(o_i)] \geq \text{poly}(\varepsilon)/k$.
623 Hence, we have $P[q_i \in \text{core}(o_i) \text{ for } i = 1 \dots s] \geq 2^{-\text{poly}(\varepsilon^{-1})} k^{-O(\varepsilon^{-1})}$. Notice, however, that
624 (\widehat{In}, Out) is a s -swap and we may have $s < p$. Nevertheless, whenever we sample $q_1 \dots q_s$ followed
625 by any sequence $q_{s+1} \dots q_p$ it is enough to choose $Out' = Out \cup \{q_{s+1} \dots q_p\}$ to obtain that
626 $(\{q_1 \dots q_p\}, Out')$ is an improving p -swap. \square

627 In order to prove Lemma 15 we first need a few technical lemmas.

628 **Lemma 19** (Lemma 2 from Lattanzi and Sohler [2019]). *For each $x, y, z \in \mathbb{R}^d$ and $\varepsilon > 0$,*
629 $\text{cost}(x, y) \leq (1 + \varepsilon) \text{cost}(x, z) + (1 + 1/\varepsilon) \text{cost}(z, y)$.

630 **Lemma 20.** *Given $q \in \mathbb{R}^d$ and $Z \subseteq \mathbb{R}^d$ such that $\text{cost}(Z, q) \leq \varepsilon^2 \Gamma$ then, for each $o \in \mathbb{R}^d$*

$$(1 - O(\varepsilon)) \text{cost}(Z, o) - O(\varepsilon) \Gamma \leq |Z| \text{cost}(q, o) \leq (1 + O(\varepsilon)) \text{cost}(Z, o) + O(\varepsilon) \Gamma$$

631 *Proof.* To obtain the first inequality, we apply Lemma 19 to bound $\text{cost}(z, o) \leq (1 + \varepsilon) \text{cost}(z, o) +$
632 $(1 + 1/\varepsilon) \text{cost}(z, q)$ for each $z \in Z$. To obtain the second inequality, we bound $\text{cost}(q, o) \leq$
633 $(1 + \varepsilon) \text{cost}(z, o) + (1 + 1/\varepsilon) \text{cost}(z, q)$ for each $z \in Z$. \square

634 **Lemma 21.** *Let $X = \{x_1 \dots x_\ell\}$ be a weighted set of points in \mathbb{R}^d such that x_i has weight w_i .*
635 *Let μ be the weighted average of X . Let $\hat{\mu} = \hat{\mu}(\alpha_1 \dots \alpha_\ell)$ be the weighted average of X where x_i*
636 *has weight α_i . If $w_i \leq \alpha_i \leq w_i/(1 - \varepsilon)$ for each $i = 1 \dots \ell$, then if we interpret $\text{cost}(X, \mathcal{C})$ as*
637 $\sum_{x_i \in X} w_i \cdot \text{cost}(x_i, \mathcal{C})$ *we have $\text{cost}(X, \hat{\mu}) \leq (1 + O(\varepsilon)) \text{cost}(X, \mu)$.*

638 *Proof.* We note that μ minimizes the expression $\text{cost}(X, \mu)$. Moreover, $\text{cost}(X, z) \leq \sum_{i=1}^{\ell} \alpha_i \cdot$
639 $\text{cost}(x_i, z) \leq \text{cost}(X, z)/(1 - \varepsilon)$. Since $\hat{\mu}$ minimizes the expression $\sum_{i=1}^{\ell} \alpha_i \cdot \text{cost}(x_i, z)$ it
640 must be $\text{cost}(X, \hat{\mu}) \leq \text{cost}(X, \mu)/(1 - \varepsilon)$. \square

641 Adopting the same proof strategy, we obtain the following.

642 **Observation 22.** *Thanks to Lemma 20, we can assume that the points in Z are concentrated in q for*
643 *the purpose of computing a $(1 + O(\varepsilon))$ -approximation to the 1-means problem on Z , whenever an*
644 *additive error Γ is tolerable. Indeed, moving all points in Z to q introduces a $1 + O(\varepsilon)$ multiplicative*
645 *error on $\text{cost}(Z, \cdot)$ and a $O(\varepsilon)\Gamma$ additive error.*

646 The next lemma shows that a point z that is far from a center o experiences a small variation of
647 $\text{cost}(z, o)$ when the position of o is slightly perturbed.

648 **Lemma 23.** *Given $o, z \in \mathbb{R}^d$ such that $\|o - z\| \geq r/\varepsilon$ we have that for every $o' \in B(o, r)$,*
649 $\text{cost}(z, o') = (1 \pm O(\varepsilon)) \text{cost}(z, o)$.

650 *Proof.* It is enough to prove it for all o' that lie on the line L passing through o and z , any other
651 point in $o'' \in B(o, r)$ admits a point $o' \in B(o, r) \cap L$ with $\|o' - z\| = \|o'' - z\|$. It is enough to
652 compute the derivative of $\text{cost}(z, \cdot)$ with respect to the direction of L and see that $\frac{\partial \text{cost}(z, \cdot)}{\partial L} \Big|_{B(o, r)} =$
653 $(1 \pm O(\varepsilon))r/\varepsilon$. Thus, $\text{cost}(z, o') = \text{cost}(z, o) \pm (1 \pm O(\varepsilon))r^2/\varepsilon = (1 \pm O(\varepsilon)) \text{cost}(z, o)$. \square

654 *Proof of Lemma 15.* Here we prove that for each $o_1 \dots o_s$ there exist coefficients $\alpha_0^{(i)} \dots \alpha_s^{(i)} \in$
655 $\left\{1, (1 - \varepsilon) \dots (1 - \varepsilon)^{\lceil \log_{1-\varepsilon}(\varepsilon^7) \rceil}\right\} \cup \{0\}$ such that the convex combination $\hat{o}_i = \hat{o}_i(\alpha_0^{(i)} \dots \alpha_s^{(i)})$
656 is a perfect approximation of o_i , with probability $k^{-O(\varepsilon^{-2})} \cdot 2^{-\text{poly}(\varepsilon^{-1})}$. Wlog, we show this

657 for o_1 only. Concretely, we want to show that, with probability $k^{-O(\varepsilon^{-1})} \cdot 2^{-poly(\varepsilon^{-1})}$, there
658 exist coefficients $\alpha_0 \dots \alpha_s$ such that $\hat{o}_1 = \hat{o}_1(\alpha_0 \dots \alpha_s)$ satisfies $\text{cost}(O_1^*, (\mathcal{C} \cup \{\hat{o}_1\}) \setminus \text{Out}) \leq$
659 $(1 + O(\varepsilon))\text{OPT}_1 + O(\varepsilon)\text{OPT}/k$. Taking the joint probability of these events for each $i = 1 \dots s$
660 we obtain the success probability $k^{-O(\varepsilon^{-2})} \cdot 2^{-poly(\varepsilon^{-1})}$. Note that we are supposed to prove that
661 $\text{cost}(O_1^*, (\mathcal{C} \cup \{\hat{o}_1\}) \setminus \text{Out}) \leq (1 + \varepsilon)\text{OPT}_1 + \varepsilon\text{OPT}/k$, however we prove a weaker version where
662 ε is replaced by $O(\varepsilon)$, which is in fact equivalent up to rescaling ε .

663 Similarly to $\mathcal{C}[\cdot]$ and $\mathcal{O}^*[\cdot]$ define $\mathcal{T}[p]$ as the closest center to p in \mathcal{T} . Denote with C_1, F_1 and N_1 the
664 intersections of O_1^* with C, F and N respectively. In what follows we define the values of $\alpha_0 \dots \alpha_s$
665 that define $\hat{o}_1 = \hat{o}_1(\alpha_0 \dots \alpha_s)$ and show an assignment of points in O_1^* to centers in $(\mathcal{C} \cup \{\hat{o}_1\}) \setminus \text{Out}$
666 with cost $(1 + O(\varepsilon))\text{OPT}_1 + O(\varepsilon)\text{OPT}/k$. Recall that we assume that $q_i \in \text{core}(o_i)$ for each
667 $i = 1 \dots s$.

668 In what follows, we assign values to the coefficients $(\alpha_i)_i$. It is understood that if the final value
669 we choose for α_i is v then we rather set α_i to the smallest power of $(1 - \varepsilon)$ which is larger than v ,
670 if $v > \varepsilon^7$. Else, set α_i to 0. We will see in the end that this restrictions on the values of α_i do not
671 impact our approximation.

672 In what follows, we will assign the points in O_1^* to $\mathcal{C} \setminus \text{Out}$, if this can be done inexpensively. If it
673 cannot, then we will assign points to \hat{o}_1 . In order to compute a good value for \hat{o}_1 we need an estimate
674 of the average of points assigned to \hat{o}_1 . For points in N_1 , computing this average is doable (leveraging
675 Lemma 13) while for points in $O_1^* \setminus N_1$ we show that either their contribution is negligible or we
676 can collapse them so as to coincide with some $q_i \in \mathcal{Q}$ without affecting our approximation. The
677 coefficients $(\alpha_i)_{i \geq 1}$ represent the fraction of points in O_i^* which is collapsed to q_i . α_0 represents the
678 fraction of points in O_1^* which average we estimate as $\mu(S)$. Thus, Equation (3) defines \hat{o}_1 as the
679 weighted average of points q_i , where the weights are the (approximate) fractions of points collapsed
680 onto q_i , together with the the average $\mu(S)$ and its associated weight α_0 .

681 **Points in C_1 .** All points $p \in C_1$ such that $\mathcal{T}[p] \notin \mathcal{Q}$ can be assigned to $\mathcal{T}[p] \in \mathcal{C} \setminus \text{Out}$ incurring a
682 total cost of at most $\varepsilon^6\text{OPT}_1$, by the definition of C_1 . Given a point $p \in C_1$ with $\mathcal{T}[p] \in \mathcal{Q}$ we might
683 have $\mathcal{T}[p] \notin \mathcal{C} \setminus \text{Out}$ and thus we cannot assign p to $\mathcal{T}[p]$. Denote with W the set of points p with
684 $\mathcal{T}[p] \in \mathcal{Q}$. Our goal is now to approximate $\mu(W)$. In order to do that, we will move each $p \in W$
685 to coincide with $q_i = \mathcal{T}[p]$. We can partition W into $W_1 \dots W_s$ so that for each $z \in W_i$ $\mathcal{T}[z] = q_i$.
686 If $p \in Z_i$ then we have $\|p - q_i\|^2 \leq \varepsilon^3 \rho_1^2$. Hence, thanks to Observation 22, we can consider
687 points in W_i as if they were concentrated in q_i while losing at most an additive factor $O(\varepsilon)\text{OPT}_1$
688 and a multiplicative factor $(1 + \varepsilon)$ on their cost. For $i = 1 \dots s$, set $\alpha_i \leftarrow |W_i|/|O_1^*|$. In this way,
689 $\sum_{i=1}^s \alpha_i \cdot q_i / \sum_{i=1}^s \alpha_i$ is an approximates solution to 1-mean on W up to a multiplicative factor
690 $(1 + \varepsilon)$ and an additive factor $O(\varepsilon)\text{OPT}_1$.

691 **Points in N_1 .** Consider the two cases: (i) $\text{cost}(N_1, \mathcal{T}) > \varepsilon^2\text{OPT}/k$; (ii) $\text{cost}(N_1, \mathcal{T}) \leq$
692 $\varepsilon^2\text{OPT}/k$.

693 Case (i). We show that in this case $\mu(S)$ is a $(1 + \varepsilon)$ -approximation for 1-mean on N_1 , with
694 probability $k^{-O(\varepsilon^{-1})} \cdot 2^{-poly(\varepsilon^{-1})}$. First, notice that if we condition on $S \subseteq N_1$ then Lemma 13
695 gives that $\mu(S)$ is a $(1 + \varepsilon)$ -approximation for 1-mean on N_1 with constant probability. Thus, we are
696 left to prove that $S \subseteq N_1$ with probability $k^{-O(\varepsilon^{-1})} \cdot 2^{-poly(\varepsilon^{-1})}$. We have that the $P_{p \sim \text{cost}(p, \mathcal{T})}[p \in$
697 $N_1 \mid p \in N] \geq \varepsilon^2/k$, however the costs w.r.t. \mathcal{T} of points in N varies of at most a factor
698 $poly(\varepsilon^{-1})$, thus $P_{p \sim \text{Unif}}[p \in N_1 \mid p \in N] \geq poly(\varepsilon)/k$. The probability of $S \subseteq N_1$ is thus
699 $(poly(\varepsilon)/k)^{|S|} = k^{-O(\varepsilon^{-1})} \cdot 2^{-poly(\varepsilon^{-1})}$. In this case, we set $\alpha_0 \leftarrow |N_1|/|O_1^*|$ because $\mu(S)$
700 approximates the mean of the entire set N_1 .

701 Case (ii). Here we give up on estimating the mean of N_1 and set $\alpha_0 \leftarrow 0$. The point $x \in N_1$ such that
702 $\mathcal{T}[x] \notin \mathcal{Q}$ can be assigned to $\mathcal{T}[x]$ incurring a combined cost of $\varepsilon^2\text{OPT}/k$. We partition the remaining
703 points in N_1 into $Z_1 \cup \dots \cup Z_s$ where each point x is placed in Z_i if $\mathcal{T}[x] = q_i$. Now, we collapse the
704 points in Z_i so as to coincide with q_i and show that this does not worsen our approximation factor. In
705 terms of coefficients $(\alpha_i)_i$, this translates into the updates $\alpha_i \leftarrow \alpha_i + |Z_i|/|O_i^*|$ for each $i = 1 \dots s$.

706 Indeed, using Observation 22 we can move all points in Z_i to q_i incurring an additive combined cost
707 of $\varepsilon\text{OPT}/k$ and a multiplicative cost of $1 + O(\varepsilon)$.

708 **Points in F_1 .** Points in F_1 are very far from q_1 and thus far from o_1 , hence even if their contribution
709 to $\text{cost}(O_1^*, o_1)$ might be large, we have $\text{cost}(F_1, o_1) = (1 \pm O(\varepsilon))\text{cost}(F_1, o')$ for all o' in a ball
710 of radius ρ_1/ε centered in o_1 , thanks to Lemma 23.

711 Let H be the set of points that have not been assigned to centers in $C \setminus \text{Out}$. In particular, $H = W \cup N_1$
712 if points in N_1 satisfy case (i) and $H = W \cup Z_1 \dots Z_s$ if points in N_1 satisfy case (ii). We consider
713 two cases.

714 If $\|\mu(H) - q_1\| \leq \rho/\varepsilon$, then $\|\mu(H) - o_1\| \leq \rho(1 + \varepsilon + 1/\varepsilon)$ because $q_1 \in \text{core}(o_1)$. Since
715 for each $f \in F_1$ we have $\|f - o_1\| \geq \|f - q_1\| - (1 + \varepsilon)\rho \geq \Omega(\rho/\varepsilon^3)$ then $\text{cost}(f, o') =$
716 $(1 \pm O(\varepsilon))\text{cost}(f, o_1)$ for each o' in a ball of radius $O(\rho/\varepsilon)$ centered in o_1 , and so in particular for
717 $o' = \mu(H)$. Thus in this case we can simply disregard all points in F_1 and computing \hat{o}_1 according
718 to the $(\alpha_i)_i$ defined above yields a perfect approximation of o_i .

719 Else, if $\|\mu(H) - q_1\| > \rho/\varepsilon$, a similar argument applies to show that $\text{cost}(H, o') = (1 \pm$
720 $\varepsilon)\text{cost}(H, o)$ for each o' in ball of radius $O(\rho)$ centered in o_1 . Indeed, we can rewrite $\text{cost}(H, o')$
721 as $|H| \cdot \text{cost}(\mu(H), o') + \text{cost}(\mu(H), H)$. If $\|\mu(H) - q_1\| < \rho/\varepsilon$ the first term varies of at most a
722 factor $(1 + \varepsilon)$ and the second term is constant. Thus in this case $\hat{o}_1 = q_1$ is a perfect approximation
723 of o_1 and we simply set $\alpha_1 = 1$ and $\alpha_j = 0$ for $j \neq 1$. In other words, here $\mu(N_1 \cup H)$ is too far
724 from q_1 (and thus o_1) to significantly influence the position of \hat{o}_1 and the same holds for any point in
725 F_1 . This works, of course, because we assumed $q_1 \in \text{core}(o_1)$. \square

726 **Discussing the limitations on the coefficients values.** The proof above would work smoothly if
727 we were allowed to set α_i to exactly the values discussed above, representing the fractions of points
728 from O_i^* captured by different q_i s. However, to make the algorithm efficient we limit ourselves to
729 values in $\{1, (1 - \varepsilon) \dots (1 - \varepsilon)^{\lceil \log_{1-\varepsilon}(\varepsilon^7) \rceil}\} \cup \{0\}$. Lemma 21 shows that as long as the values of
730 $(\alpha_i)_i$ estimate the frequencies described above up to a factor $1 \pm O(\varepsilon)$ then the approximation error
731 is within a multiplicative factor $1 \pm O(\varepsilon)$.

732 We are left to take care of the case in which α_i is set to a value $< \varepsilon^7$. We set α_i when dealing
733 with points in $C_1 \cup N_1$ and for each $x \in C_1 \cup N_1$ we have, for each $o' \in B(q_1, (1 + \varepsilon)\rho)$,
734 $\text{cost}(x, o') \leq 2\text{cost}(q_1, o') + 2\text{cost}(x, q_1) = O(\rho_1 \varepsilon^{-6})$. Thus, if we simply set $\alpha_i \leftarrow 0$ whenever
735 we have $\alpha_i < \varepsilon^7$ then the combined cost of points in O_1^* with respect to o' varies by $\varepsilon^7 |O_1^*| \cdot \rho_1 \varepsilon^{-6} =$
736 $O(\varepsilon)\text{OPT}_1$. Effectively, ignoring these points does not significantly impact the cost. hence solving
737 1-mean ignoring these points finds a $(1 + O(\varepsilon))$ -approximate solution to the original problem.

738 Additional Experimental Evaluation

739 In this section we report additional experiments which presentation did not fit in the main body. In
740 particular, we run experiments on the dataset KDD-PHY and for $k = 10, 50$.

741 In Figure 3 we compare MSLS-G with MSLS. To perform our experiment, we initialize $k = 25$
742 centers using KM++ and then run 50 iterations of local search for both algorithms, for $p \in \{2, 3\}$
743 swaps. We repeat each experiment 5 times. For ease of comparison, we repeat the plot for the
744 KDD-BIO and RNA datasets that we present in the main body of the paper. Due to the higher running
745 of the MSLS we perform this experiments on 1% uniform sample of each of our datasets. We find
746 out that the performance of the two algorithms is comparable on all our instances, while they both
747 perform roughly 15%-27% better than k -means++ at convergence.

748 In Figure 4 we run KM++ followed by 50 iterations of MSLS-G with $p = 1, 4, 7, 10$ and $k =$
749 $10, 25, 50$ (excluding the degenerate case $p = k = 10$) and plot the relative cost w.r.t. KM++ at
750 each iteration. The results for $k = 25$ on KDD-BIO and RNA can be found in Figure 2. We repeat
751 each experiment 5 times. Our experiment shows that, after 50 iterations MSLS-G for $p = 4, 7, 10$
752 achieves improvements of roughly 5 – 10% compared to MSLS-G- $p = 1$ and of the order of
753 20% – 40% compared to KM++. These improvements are more prominent for $k = 25, 50$. We also
754 report the time per iteration that each algorithm takes. For comparison, we report the running time
755 of a single iteration of Lloyd's next to the dataset's name. Notice that the experiment on RNA for
756 $k = 50$ is performed on a 10% uniform sample of the original dataset, due to the high running time.

757 In Figure 5, we use KM++ and MSLS-G as a seeding algorithm for Lloyd's and measure how much
758 of the performance improvement measured is retained after running Lloyd's. First, we initialize

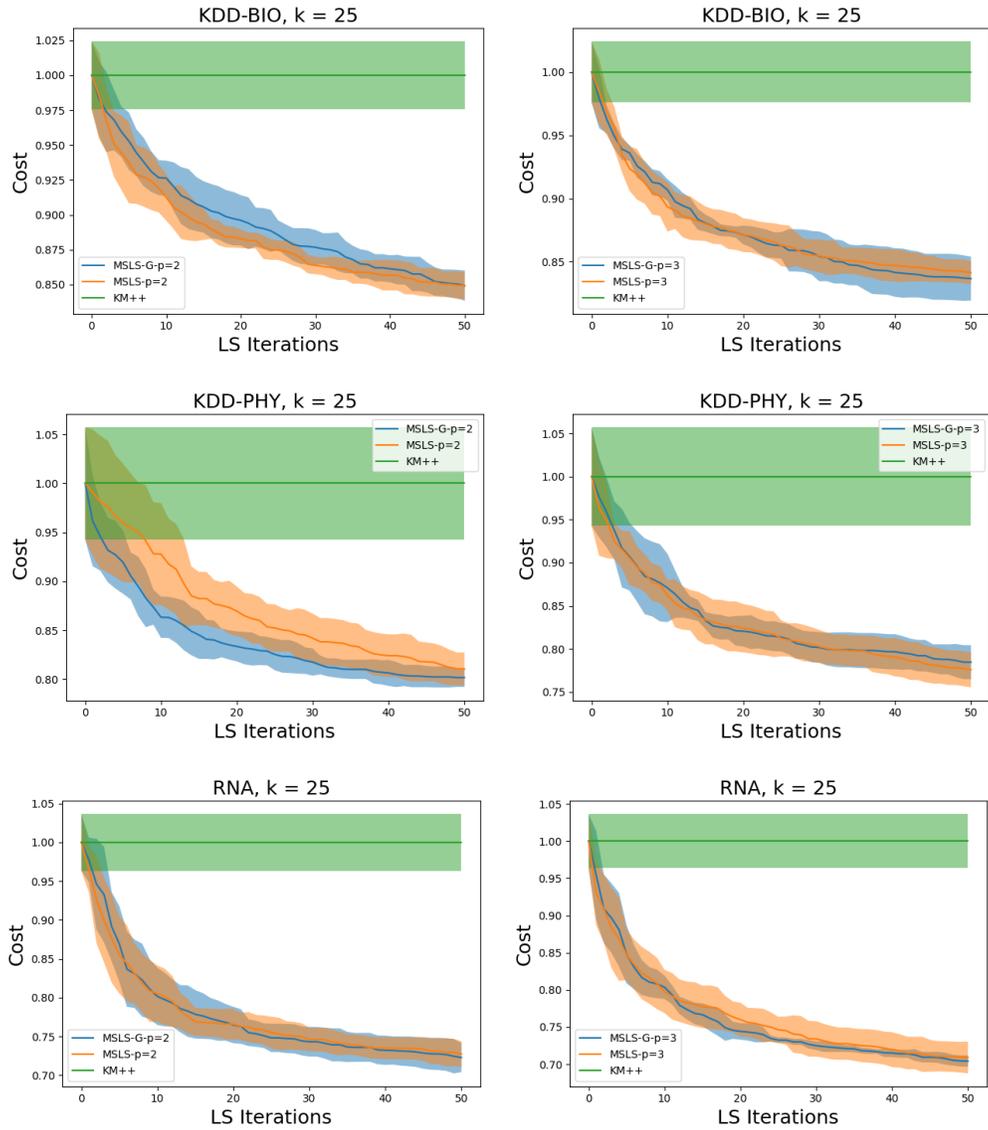


Figure 3: Comparison between MSLS and MSLS-G, for $p = 2$ (left column) and $p = 3$ (right column), for $k = 25$, on the datasets KDD-BIO (first row), KDD-PHY (second row) and RNA (third row). The y axis shows the mean solution cost, over the 5 repetitions of the experiment, divided by the means solution cost of KM++.

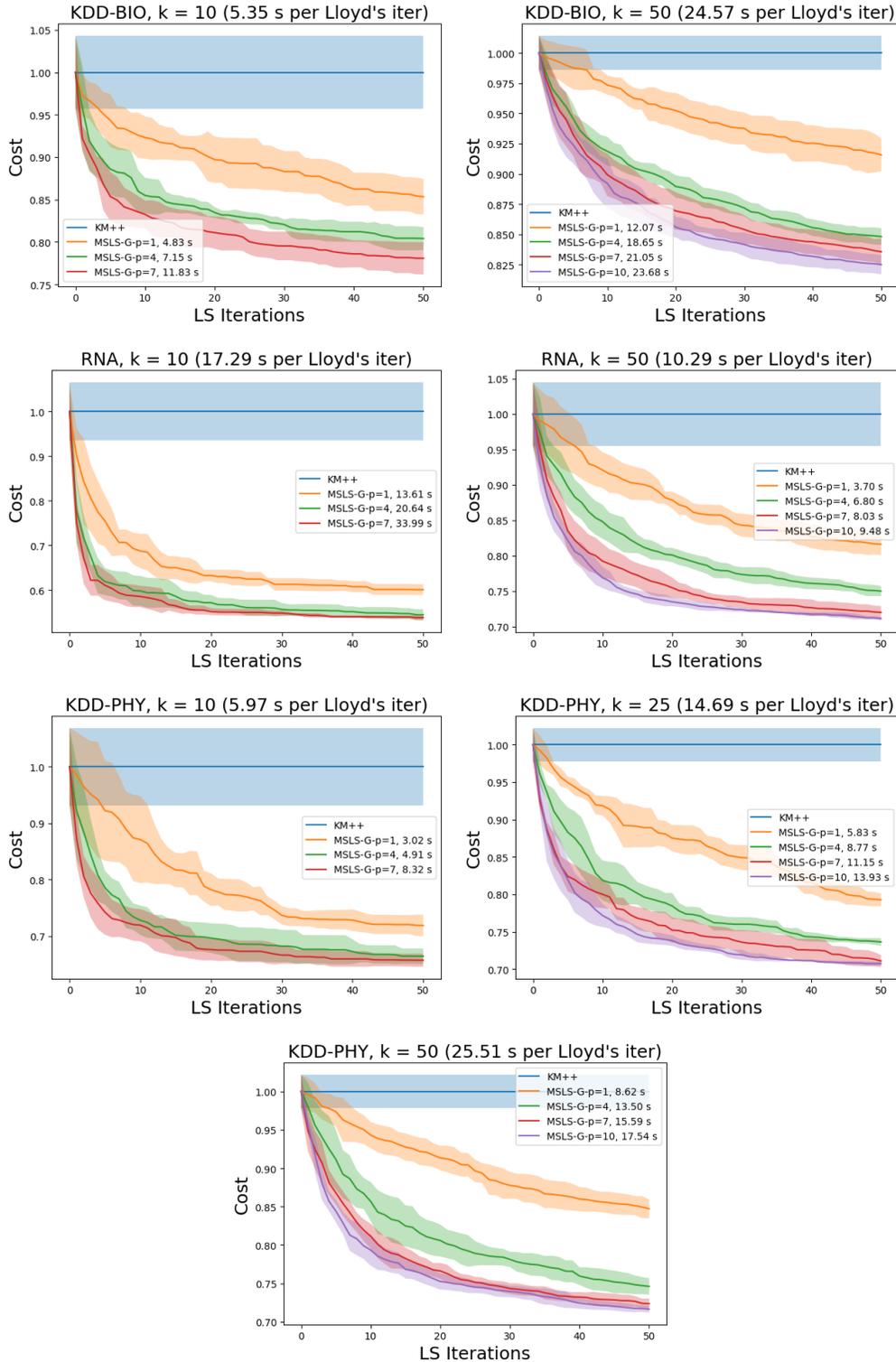


Figure 4: We compare the cost of MSLS-G, for $p \in \{1, 4, 7, 10\}$, divided by the mean cost of KM++ at each LS step, for $k \in \{10, 25, 50\}$, excluding the degenerate case $p = k = 10$. The legend reports also the running time of MSLS-G per LS step (in seconds). The experiments were run on all datasets: KDD-BIO, RNA and KDD-PHY, excluding the case of $k = 25$ for KDD-BIO and RNA which are reported in the main body of the paper.

759 our centers using KM++ and the run 15 iterations of MSLS-G for $p = 1, 4, 7$. We measure the
760 cost achieved by running 10 iterations of Lloyd’s starting from the solutions found by MSLS-G as
761 well as KM++. We run experiments for $k = 10, 25, 50$ and we repeat each experiment 5 times. We
762 observe that for $k = 25, 50$ MSLS-G for $p > 1$ performs at least as good as SLS from Lattanzi
763 and Sohler [2019] and in some cases maintains non-trivial improvements. These improvements are
764 not noticeable for $k = 10$; however, given how Lloyd’s behave for $k = 10$ we conjecture that $k = 10$
765 might be an “unnatural” number of clusters for our datasets.

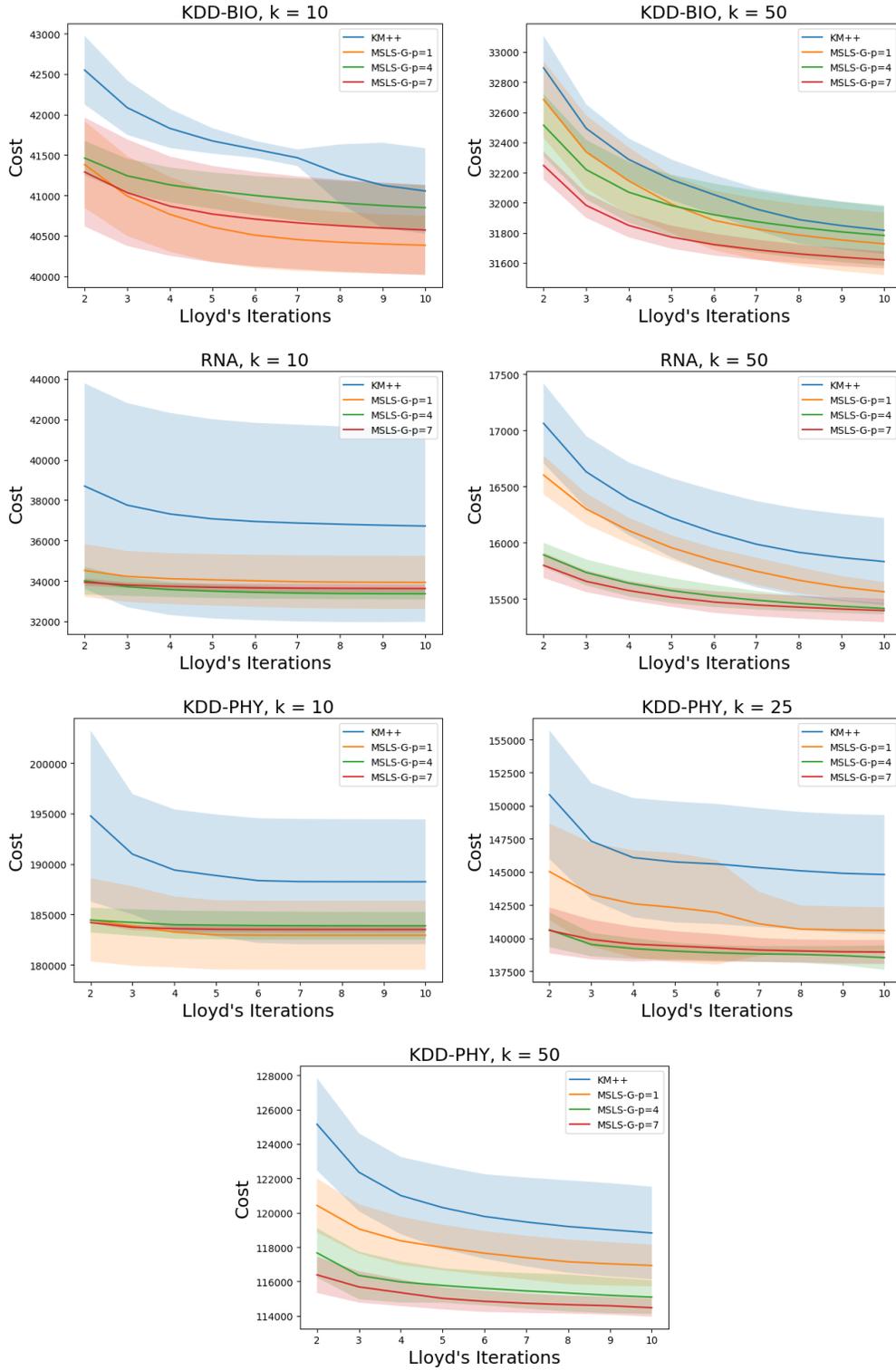


Figure 5: We compare the cost after each of the 10 iterations of Lloyd with seeding from MSLS-G, for $p \in \{1, 4, 7, 10\}$ and 15 local search steps and KM++, for $k \in \{10, 25, 50\}$. We excluded the degenerate case $p = k = 10$, and the experiments reported in the main body of the paper.