

Improved Smoothed Analysis of the k -Means Method*

Bodo Manthey[†]

Saarland University
Department of Computer Science
manthey@cs.uni-sb.de

Heiko Röglin[‡]

Boston University
Department of Computer Science
heiko@roeglin.org

September 10, 2008

The k -means method is a widely used clustering algorithm. One of its distinguished features is its speed in practice. Its worst-case running-time, however, is exponential, leaving a gap between practical and theoretical performance. Arthur and Vassilvitskii [3] aimed at closing this gap, and they proved a bound of $\text{poly}(n^k, \sigma^{-1})$ on the smoothed running-time of the k -means method, where n is the number of data points and σ is the standard deviation of the Gaussian perturbation. This bound, though better than the worst-case bound, is still much larger than the running-time observed in practice.

We improve the smoothed analysis of the k -means method by showing two upper bounds on the expected running-time of k -means. First, we prove that the expected running-time is bounded by a polynomial in $n^{\sqrt{k}}$ and σ^{-1} . Second, we prove an upper bound of $k^{kd} \cdot \text{poly}(n, \sigma^{-1})$, where d is the dimension of the data space. The polynomial is independent of k and d , and we obtain a polynomial bound for the expected running-time for $k, d \in O(\sqrt{\log n / \log \log n})$.

Finally, we show that k -means runs in smoothed polynomial time for one-dimensional instances.

1 Introduction

The k -means method is a very popular algorithm for clustering high-dimensional data. It is based on ideas by Lloyd [10]. It is a local search algorithm: Initiated with k arbitrary cluster centers, it assigns every data point to its nearest center, and then readjusts the centers, reassigns the data points, ... until it stabilizes. (In Section 1.1, we describe the algorithm formally.) The k -means method terminates in a local optimum, which might be far worse than the global optimum. However, in practice it works very well. It is particularly popular because of its simplicity and its speed: “In practice, the number of iterations is much less than the

*An extended abstract of this work will appear in *Proc. of the 20th ACM-SIAM Symposium on Discrete Algorithms (SODA 2009)*.

[†]Work done in part at Yale University, Department of Computer Science, supported by the Postdoc-Program of the German Academic Exchange Service (DAAD).

[‡]Supported by a fellowship within the Postdoc-Program of the German Academic Exchange Service (DAAD).

number of samples”, as Duda et al. [6, Section 10.4.3] put it. According to Berkhin [5], the k -means method “is by far the most popular clustering tool used in scientific and industrial applications.”

The practical performance and popularity of the k -means method is at stark contrast to its performance in theory. The only upper bounds for its running-time are based on the observation that no clustering appears twice in a run of k -means: Obviously, n points can be distributed among k clusters in only k^n ways. Furthermore, the number of Voronoi partitions of n points in \mathbb{R}^d into k classes is bounded by a polynomial in n^{kd} [8], which yields an upper bound of $\text{poly}(n^{kd})$. On the other hand, Arthur and Vassilvitskii [2] showed that k -means can run for $2^{\Omega(\sqrt{n})}$ iterations in the worst case.

To close the gap between good practical and poor theoretical performance of algorithms, Spielman and Teng introduced the notion of smoothed analysis [12]: An adversary specifies an instance, and this instance is then subject to slight random perturbations. The smoothed running-time is the maximum over the adversarial choices of the expected running-time. On the one hand, this rules out pathological, isolated worst-case instances. On the other hand, smoothed analysis, unlike average-case analysis, is not dominated by random instances since the instances are not completely random; random instances are usually not typical instances and have special properties with high probability. Thus, smoothed analysis also circumvents the drawbacks of average-case analysis. For a survey of smoothed analysis, we refer to Spielman and Teng [13].

The goal of this paper is to bound the smoothed running-time of the k -means method. There are basically two reasons why the smoothed running-time of the k -means method is a more realistic measure than its worst-case running-time: First, data obtained from measurements is inherently noisy. So even if the original data were a bad instance for k -means, the data measured is most likely a slight perturbation of it. Second, if the data possesses a meaningful k -clustering, then slightly perturbing the data should preserve this clustering. Thus, smoothed analysis might help to obtain a faster k -means method: We take the data measured, perturb it slightly, and then run k -means on the perturbed instance. The bounds for the smoothed running-time carry over to this variant of the k -means method.

1.1 k -Means Method

An instance for k -means clustering is a point set $\mathcal{X} \subseteq \mathbb{R}^d$ consisting of n points. The aim is to find a clustering $\mathcal{C}_1, \dots, \mathcal{C}_k$ of \mathcal{X} , i.e., a partition of \mathcal{X} , as well as cluster centers such that the potential

$$\sum_{i=1}^k \sum_{x \in \mathcal{C}_i} \|x - c_i\|^2$$

is minimized. Given the cluster centers, every data point should obviously be assigned to the cluster whose center is closest to it. The name k -means stems from the fact that, given the clusters, the centers c_1, \dots, c_k should be chosen as the centers of mass, i.e., $c_i = \frac{1}{|\mathcal{C}_i|} \sum_{x \in \mathcal{C}_i} x$ of \mathcal{C}_i . The k -means method proceeds now as follows:

1. Select cluster centers c_1, \dots, c_k .
2. Assign every $x \in \mathcal{X}$ to the cluster \mathcal{C}_i whose cluster center c_i is closest to it.
3. Set $c_i = \frac{1}{|\mathcal{C}_i|} \sum_{x \in \mathcal{C}_i} x$.

4. If clusters or centers have changed, goto 2. Otherwise, terminate.

Since the potential decreases in every step, no clustering occurs twice, and the algorithm eventually terminates.

1.2 Related Work

The problem of finding a good clustering can be approximated arbitrarily well: Bădoiu et al. [4], Matoušek [11], and Kumar et al. [9] devised polynomial time approximation schemes with different dependencies on the approximation ratio $(1+\varepsilon)$ as well as n , k , and d : $O(2^{O(k\varepsilon^{-2} \log k)} \cdot nd)$, $O(n\varepsilon^{-2k^2d} \log^k n)$, and $O(\exp(k/\varepsilon) \cdot nd)$, respectively.

While the polynomial time approximation schemes show that k -means clustering can be approximated arbitrarily well, the method of choice for finding a k -clustering is the k -means method due to its performance in practice. However, the only polynomial bound for k -means holds for $d = 1$, and only for instances with polynomial spread [7], which is the maximum distance of points divided by the minimum distance.

Arthur and Vassilvitskii [3] have analyzed the running-time of the k -means method subject to Gaussian perturbation: The points are drawn according to independent d -dimensional Gaussian distributions with standard deviation σ . Arthur and Vassilvitskii proved that the expected running-time after perturbing the input with Gaussians with standard deviation σ is polynomial in n^k , d , the diameter of the perturbed point set, and $1/\sigma$.

Recently, Arthur [1] showed that the probability that the running-time of k -means subject to Gaussian perturbations exceeds a polynomial in n , d , the diameter of the instance, and $1/\sigma$ is bounded by $O(1/n)$. However, his argument does not yield any significant bound on the expected running-time of k -means: The probability of $O(1/n)$ that the running-time exceeds a polynomial bound is too large to yield an upper bound for the expected running-time, except for the trivial upper bound of $\text{poly}(n^{kd})$.

1.3 New Results

We improve the smoothed analysis of the k -means method by proving two upper bounds on its running-time. First, we show that the smoothed running-time of k -means is bounded by a polynomial in $n^{\sqrt{k}}$ and $1/\sigma$.

Theorem 1. *Let $\mathcal{X} \subseteq \mathbb{R}^d$ be a set of n points drawn according to independent Gaussian distributions whose means are in $[0, 1]^d$. Then the expected running-time of the k -means method on the instance \mathcal{X} is bounded from above by a polynomial in $n^{\sqrt{k}}$ and $1/\sigma$.*

Thus, compared to the previously known bound, we decrease the exponent by a factor of \sqrt{k} . Second, we show that the smoothed running-time of k -means is bounded by $k^{kd} \cdot \text{poly}(n, 1/\sigma)$. In particular, this decouples the exponential part of the bound from the number n of points.

Theorem 2. *Let \mathcal{X} be drawn as described in Theorem 1. Then the expected running-time of the k -means method on the instance \mathcal{X} is bounded from above by $k^{kd} \cdot \text{poly}(n, 1/\sigma)$.*

An immediate consequence of Theorem 2 is the following corollary, which proves that the expected running-time is polynomial in n and $1/\sigma$ if k and d are small compared to n . This result is of particular interest since d and k are usually much smaller than n .

Corollary 3. *Let $k, d \in O(\sqrt{\log n / \log \log n})$. Let \mathcal{X} be drawn as described in Theorem 1. Then the expected running-time of k -means on the instance \mathcal{X} is bounded by a polynomial in n and $1/\sigma$.*

David Arthur [1] presented an insightful proof that k -means runs in time polynomial in n , $1/\sigma$, and the diameter of the instance with a probability of at least $1 - O(1/n)$. It is worth pointing out that his result is orthogonal to our results: neither do our results imply polynomial running time with probability $1 - O(1/n)$, nor does Arthur's result yield any non-trivial bound on the expected running-time (not even $\text{poly}(n^k, 1/\sigma)$) since the success probability of $1 - O(1/n)$ is way too small. The exception is our result for $d = 1$, which yields not only a bound on the expectation, but also a bound that holds with high probability. However, the original definition of smoothed analysis [12] is in terms of expectation, not in terms of bounds that hold with a probability of $1 - o(1)$.

To prove our bounds, we prove a lemma about perturbed point sets (Lemma 5). The lemma bounds the number of points close to the boundaries of Voronoi partitions that arise during the execution of k -means. It might be of independent interest, in particular for smoothed analyses of geometric algorithms and problems.

Finally, we prove a polynomial bound for the running-time of k -means in one dimension.

Theorem 4. *Let $\mathcal{X} \subseteq \mathbb{R}$ be drawn according to 1-dimensional Gaussian distributions as described in Theorem 1. Then the expected running-time of k -means on \mathcal{X} is polynomial in n and $1/\sigma$. Furthermore, the probability that the running-time exceeds a polynomial in n and $1/\sigma$ is bounded by $1/\text{poly}(n)$.*

We remark that this result for $d = 1$ is not implied by the result of Har-Peled and Sadri [7] that the running-time of one-dimensional k -means is polynomial in n and the spread of the instance. The reason is that the expected value of the square of the spread is unbounded.

The restriction of the adversarial points to be in $[0, 1]^d$ is necessary: Without any bound, the adversary can place the points arbitrarily far away, thus diminishing the effect of the perturbation. We can get rid of this restriction and obtain the same results by allowing the bounds to be polynomial in the diameter of the adversarial instance. However, for the sake of clarity and to avoid another parameter, we have chosen the former model.

1.4 Outline

To prove our two main theorems, we first prove a property of perturbed point sets (Section 2): In any step of the k -means algorithm, there are not too many points close to any of the at most k^2 hyperplanes that bisect the centers and that form the Voronoi regions. To put it another way: No matter how k -means partitions the point set \mathcal{X} into k Voronoi regions, the number of points close to any boundary is rather small with overwhelming probability.

We use this lemma in Section 3: First, we use it to prove Lemma 8, which bounds the expected number of iterations in terms of the smallest possible distance of two clusters. Using this bound, we derive a first upper bound for the expected number of iterations (Lemma 9), which will result in Theorem 2 later on.

In Sections 4 and 5, we distinguish between iterations in which at most \sqrt{k} or at least \sqrt{k} clusters gain or lose points. This will result in Theorem 1.

We consider the special case of $d = 1$ in Section 6. For this case, we prove an upper bound polynomial in n and $1/\sigma$ until the potential has dropped by at least 1.

In Sections 3, 4, 5, and 6 we are only concerned with bounding the number of iterations until the potential has dropped by at least 1. Using these bounds and an upper bound on the potential after the first round, we will derive Theorems 1, 2, and 4 as well as Corollary 3 in Section 7.

1.5 Preliminaries

In the following, \mathcal{X} is the perturbed instance on which we run k -means, i.e., $\mathcal{X} = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$ is a set of n points, where each point x_i is drawn according to a d -dimensional Gaussian distribution with mean $\mu_i \in [0, 1]^d$ and standard deviation σ .

Inaba et al. [8] proved that the number of iterations of k -means is $\text{poly}(n^{kd})$ in the worst case. We abbreviate this bound by $W \leq n^{\kappa kd}$ for some constant κ in the following.

Let $D \geq 1$ be chosen such that, with a probability of at least $1 - W^{-1}$, every data point from \mathcal{X} lies in the hypercube $\mathcal{D} := [-D, 1 + D]^d$ after the perturbation. In Section 7, we prove that D can be bounded by a polynomial in n and σ , and we use this fact in the following sections. We denote by \mathcal{F} the *failure event* that there exists one point in \mathcal{X} that does not lie in the hypercube \mathcal{D} after the perturbation. We say that a cluster is *active* in an iteration if it gains or loses at least one point.

We will always assume in the following that $d \leq n$ and $k \leq n$, and we will frequently bound both d and k by n to simplify calculations. Of course, $k \leq n$ holds for every meaningful instance since it does not make sense to partition n points into more than n clusters. Furthermore, we can assume $d \leq n$ for two reasons: First, the dimension is usually much smaller than the number of points, and, second, if $d > n$, then we can project the points to a lower-dimensional subspace without changing anything.

Let $\mathcal{C} = \{\mathcal{C}_1, \dots, \mathcal{C}_k\}$ denote the set of clusters. For a natural number k , let $[k] = \{1, \dots, k\}$. In the following, we will assume that number such as \sqrt{k} are integers. For the sake of clarity, we do not write down the tedious floor and ceiling functions that are actually necessary. Since we are only interested in the asymptotics, this does not affect the validity of the proofs. Furthermore, we assume in the following sections that $\sigma \leq 1$. This assumption is only made to simplify the arguments and we describe in Section 7 how to get rid of it.

2 A Property of Perturbed Point Sets

The following lemma shows that, with high probability, there are not too many points close to the hyperplanes dividing the clusters. It is crucial for our bounds for the smoothed running-time: If not too many points are close to the bisecting hyperplanes, then, eventually, one point that is further away from the bisecting hyperplanes must go from one cluster to another, which causes a significant decrease of the potential.

Lemma 5. *Let $a \in [k]$ be arbitrary. With a probability of at least $1 - 2W^{-1}$, the following holds: In every step of the k -means algorithm (except for the first one) in which at least kd/a points change their assignment, at least one of these points has a distance larger than*

$$\varepsilon := \frac{\sigma^4}{32n^2dD^2} \cdot \left(\frac{\sigma}{3Dn^{3+2\kappa}} \right)^{4a}$$

from the bisector that it crosses.

Proof. We consider a step of the k -means algorithm, and we refer to the configuration before this step as the *first configuration* and to the configuration after this step as the *second configuration*. To be precise, we assume that in the first configuration the positions of the centers are the centers of mass of the points assigned to them in this configuration. The step we consider is the reassignment of the points according to the Voronoi diagram in the first configuration.

Let $B \subseteq \mathcal{X}$ with $|B| = \ell := kd/a$ be a set of points that change their assignment during the step. There are at most n^ℓ choices for the points in B and at most $k^{2\ell} \leq n^{2\ell}$ choices for the clusters they are assigned to in the first and the second configuration. We apply a union bound over all these at most $n^{3\ell}$ choices.

The following sets are defined for all $i, j \in [k]$ and $j \neq i$. Let $B_i \subseteq B$ be the set of points that leave cluster \mathcal{C}_i . Let $B_{i,j} \subseteq B_i$ be the set of points assigned to cluster \mathcal{C}_i in the first and to cluster \mathcal{C}_j in the second configuration, i.e., the points in $B_{i,j}$ leave \mathcal{C}_i and enter \mathcal{C}_j . We have $B = \bigcup_i B_i$ and $B_i = \bigcup_{j \neq i} B_{i,j}$.

Let A_i be the set of points that are in \mathcal{C}_i in the first configuration except for those in B_i . We assume that the positions of the points in A_i are determined by an adversary. Since the sets A_1, \dots, A_k form a partition of the points in $\mathcal{X} \setminus B$ that has been obtained in the previous step on the basis of a Voronoi diagram, there are at most W choices for this partition [8]. We also apply a union bound over the choices for this partition.

In the first configuration, exactly the points in $A_i \cup B_i$ are assigned to cluster \mathcal{C}_i . Let c_1, \dots, c_k denote the positions of the cluster centers in the first configuration, that is, c_i is the center of mass of $A_i \cup B_i$. Since the positions of the points in $\mathcal{X} \setminus B$ are assumed to be fixed by an adversary, and since we apply a union bound over the partition A_1, \dots, A_k , the impact of the set A_i on position c_i is fixed. However, we want to exploit the randomness of the points in B_i in the following. Thus, the positions of the centers are not fixed yet but they depend on the randomness of the points in B . In particular, the bisecting hyperplane $H_{i,j}$ of the clusters \mathcal{C}_i and \mathcal{C}_j is not fixed but depends on B_i and B_j .

In order to complete the proof, we have to estimate the probability of the event

$$\forall i, j: \forall b \in B_{i,j}: \text{dist}(b, H_{i,j}) \leq \varepsilon, \quad (\mathcal{E})$$

where $\text{dist}(x, H) = \min_{y \in H} \|x - y\|$ denotes the shortest distance of a point x to a hyperplane H . In the following, we denote this event by \mathcal{E} . If the hyperplanes $H_{i,j}$ were fixed, the probability of \mathcal{E} could readily be seen to be at most $\left(\frac{2\varepsilon}{\sigma\sqrt{2\pi}}\right)^\ell \leq \left(\frac{\varepsilon}{\sigma}\right)^\ell$. But the hyperplanes are not fixed since their positions and orientations depend on the points in the sets $B_{i,j}$. Therefore, we are only able to prove the following weaker bound in Lemma 6:

$$\Pr[\mathcal{E} \wedge \neg\mathcal{F}] \leq \left(\frac{3D}{\sigma}\right)^{kd} \cdot \left(\frac{32n^2 d D^2 \varepsilon}{\sigma^4}\right)^{\ell/4},$$

where $\neg\mathcal{F}$ denotes the event that, after the perturbation, all points of \mathcal{X} lie in the hypercube $\mathcal{D} = [-D, D + 1]^d$. Now the union bound yields the following upper bound on the probability

that a set B with the stated properties exists:

$$\begin{aligned}
\Pr[\mathcal{E}] &\leq \Pr[\mathcal{E} \wedge \neg\mathcal{F}] + \Pr[\mathcal{F}] \\
&\leq n^{3\ell}W \cdot \left(\frac{3D}{\sigma}\right)^{kd} \cdot \left(\frac{32n^2dD^2\varepsilon}{\sigma^4}\right)^{\ell/4} + W^{-1} \\
&= n^{3\ell}W \cdot \left(\frac{1}{n^{3+2\kappa}}\right)^{kd} + W^{-1} \\
&\leq n^{3\ell+\kappa kd} \cdot \left(\frac{1}{n^{3+2\kappa}}\right)^{kd} + W^{-1} \\
&\leq n^{-\kappa kd} + W^{-1} \leq 2W^{-1}.
\end{aligned}$$

The equation is by our choice of ε , the inequalities are due to some simplifications and $W \leq n^{\kappa kd}$. \square

Lemma 6. *The probability of the event $\mathcal{E} \wedge \neg\mathcal{F}$ is bounded from above by*

$$\left(\frac{3D}{\sigma}\right)^{kd} \cdot \left(\frac{32n^2dD^2\varepsilon}{\sigma^4}\right)^{\ell/4}.$$

Proof. Let g_i be the random vector that equals the sum of the points in B_i , i.e.,

$$g_i := \sum_{b \in B_i} b.$$

Due to the union bound, the influence of A_i and A_j on the hyperplane $H_{i,j}$ is fixed. Since the union bound also fixes the number of points in B_i and B_j , it suffices to know the sums g_i and g_j to deduce the exact position of the hyperplane $H_{i,j}$. Hence, once all sums g_i are fixed, all hyperplanes are fixed as well. The drawback is, of course, that fixing the sums g_i has an impact on the distribution of the random positions of the points in B_i . We circumvent this problem as follows: We basically show that if B_i contains m points and the sum g_i is fixed, then we can still use the randomness of $m - 1$ of these points. For sets B_i that contain at least two points, this means that we can use the randomness of at least half of its points. Complications are only caused by sets B_i that consist of a single point. For such sets, fixing g_i is equivalent to fixing the position of the point, and we give a more direct analysis without fixing g_i .

For $y_i, y_j \in \mathbb{R}^d$, we denote by $H_{i,j}(y_i, y_j)$ the bisector of the clusters \mathcal{C}_i and \mathcal{C}_j that is obtained for $g_i = y_i$ and $g_j = y_j$. Let k^* be the number of clusters \mathcal{C}_i with $|B_i| > 0$. Without loss of generality, these are the clusters $\mathcal{C}_1, \dots, \mathcal{C}_{k^*}$. This convention allows us to rewrite the probability of $\mathcal{E} \wedge \neg\mathcal{F}$ as

$$\begin{aligned}
\Pr[\forall i, j: \forall b \in B_{i,j}: \text{dist}(b, H_{i,j}) \leq \varepsilon \wedge \neg\mathcal{F}] &\leq \int_{\mathcal{D}} \cdots \int_{\mathcal{D}} \left(\prod_{i=1}^{k^*} f_{g_i}(y_i) \right) \\
&\cdot \Pr[\forall i, j: \forall b \in B_{i,j}: \text{dist}(b, H_{i,j}(y_i, y_j)) \leq \varepsilon \mid \forall i: g_i = y_i] dy_{k^*} \cdots dy_1,
\end{aligned}$$

where f_{g_i} is the density of the random vector g_i . We admit that our notation is a bit sloppy: If $|B_{i,j}| > 0$ and $j \notin \{1, \dots, k^*\}$, then $H_{i,j}$ depends only on y_i . In this case, we should actually write $H_{i,j}(y_i)$ instead of $H_{i,j}(y_i, y_j)$ in the formula above. In order to keep the notation less

cumbersome, we ignore this subtlety and assume that $H_{i,j}(y_i, y_{j_i})$ is implicitly replaced by $H_{i,j}(y_i)$ whenever necessary. Points from different sets B_i and B_j are independent even under the assumption that the sums g_i and g_j are fixed. Hence, we can further rewrite the probability as

$$\int_{\mathcal{D}} \cdots \int_{\mathcal{D}} \left(\prod_{i=1}^{k^*} f_{g_i}(y_i) \right) \cdot \left(\prod_{i=1}^{k^*} \Pr[\forall j: \forall b \in B_{i,j}: \text{dist}(b, H_{i,j}(y_i, y_j)) \leq \varepsilon \mid g_i = y_i] \right) dy_{k^*} \dots dy_1. \quad (1)$$

Now let us consider the probability

$$\Pr[\forall j: \forall b \in B_{i,j}: \text{dist}(b, H_{i,j}(y_i, y_j)) \leq \varepsilon \mid g_i = y_i] \quad (2)$$

for a fixed i and for fixed values y_i and y_j . To simplify the notation, let $B_i = \{b_1, \dots, b_m\}$, and let the corresponding hyperplanes (which are fixed because y_i and the y_j 's are given) be H_1, \dots, H_m . (A hyperplane may occur several times in this list if more than one point goes from \mathcal{C}_i to some cluster \mathcal{C}_j .) Then the probability simplifies to

$$\Pr[\forall j: \text{dist}(b_j, H_j) \leq \varepsilon \mid g_i = y_i].$$

We distinguish between the cases $m = 1$ and $m > 1$.

Case 1: $m = 1$. The probability degenerates to

$$\Pr[\text{dist}(b_1, H_1) \leq \varepsilon \mid g_i = b_1 = y_i] = \begin{cases} 1 & \text{if } y_i \text{ is } \varepsilon\text{-close to } H_1, \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

So, given $g_i = y_i$, there is no randomness in the event that y_i is ε -close to H_1 . Choose j_i such that $b_1 \in B_{i,j_i}$ and denote by $\mathbb{I}_i(y_i, y_{j_i})$ the indicator variable defined in (3). We replace the corresponding probability in (1) by $\mathbb{I}_i(y_i, y_{j_i})$.

Case 2: $m > 1$. Let $H_j(\varepsilon)$ be the slab of width 2ε around H_j , i.e., $H_j(\varepsilon) = \{x \in \mathbb{R}^d \mid \text{dist}(x, H_j) \leq \varepsilon\}$. Let f be the joint density of the random vectors b_1, \dots, b_{m-1}, g_i . Then the probability (2) can be bounded from above by

$$\int_{z_1 \in H_1(\varepsilon)} \cdots \int_{z_{m-1} \in H_{m-1}(\varepsilon)} \frac{f(z_1, \dots, z_{m-1}, y_i)}{f_{g_i}(y_i)} dz_{m-1} \dots dz_1.$$

Now let f_i be the density of the random vector b_i . This allows us to rewrite the joint density, and we obtain the upper bound

$$\begin{aligned} & \int_{z_1 \in H_1(\varepsilon)} \cdots \int_{z_{m-1} \in H_{m-1}(\varepsilon)} \frac{f_1(z_1) \cdots f_{m-1}(z_{m-1}) \cdot f_m(y_i - \sum_{j=1}^{m-1} z_j)}{f_{g_i}(y_i)} dz_{m-1} \dots dz_1 \\ & \leq \frac{1}{f_{g_i}(y_i) \cdot \sigma^d} \int_{z_1 \in H_1(\varepsilon)} \cdots \int_{z_{m-1} \in H_{m-1}(\varepsilon)} f_1(z_1) \cdots f_{m-1}(z_{m-1}) dz_{m-1} \dots dz_1 \\ & = \frac{1}{f_{g_i}(y_i) \cdot \sigma^d} \left(\prod_{i=1}^{m-1} \int_{z_i \in H_i(\varepsilon)} f_i(z_i) dz_i \right) \\ & \leq \frac{1}{f_{g_i}(y_i) \cdot \sigma^d} \cdot \left(\frac{\varepsilon}{\sigma} \right)^{m-1}. \end{aligned}$$

The first inequality follows from $f_m(\cdot) \leq \sigma^{-d}$, and the last inequality follows because the probability that a Gaussian random vector assumes a position within distance ε of a given hyperplane is at most ε/σ .

Now we plug the bounds derived in Cases 1 and 2 into (1). Let $k_=-$ be the number of clusters \mathcal{C}_i with $|B_i| = 1$, and let us assume that these are the clusters $\mathcal{C}_1, \dots, \mathcal{C}_{k_-}$. Let

$$k_> = |\{i \mid |B_i| > 1\}| \quad \text{and} \quad m' = \sum_{i: |B_i| > 1} (|B_i| - 1),$$

that is, $k_>$ is the number of clusters that lose more than one point, and m' is the number of points leaving those clusters minus $k_>$, i.e., the number of points whose randomness we exploit. Note that $k^* = k_> + k_-$. Then (1) can be bounded from above by

$$\frac{1}{\sigma^{k_>d}} \cdot \left(\frac{\varepsilon}{\sigma}\right)^{m'} \int_{\mathcal{D}} \cdots \int_{\mathcal{D}} \left(\prod_{i=1}^{k_-} f_{g_i}(y_i) \cdot \mathbb{I}_i(y_i, y_{j_i}) \right) dy_{k^*} \dots dy_1, \quad (4)$$

where f_{g_i} cancels out for $i > k_-$. Observe that for fixed y_{j_i} the term

$$\int_{\mathcal{D}} f_{g_i}(y_i) \cdot \mathbb{I}_i(y_i, y_{j_i}) dy_i \quad (5)$$

describes the probability that the point b with $B_i = \{b\}$ lies in \mathcal{D} and is at distance at most ε from the bisector of \mathcal{C}_i and \mathcal{C}_{j_i} . For $y_{j_i} \in \mathcal{D}$, the point b can only lie in the hypercube \mathcal{D} if it has a distance of at most $\sqrt{d}(1 + 2D) \leq 3\sqrt{d}D$ from y_{j_i} . Hence, we can use Lemma 7 with $\delta = 3\sqrt{d}D$ to obtain that the probability in (5) is upper bounded by

$$\frac{2\sqrt{n3\sqrt{d}D}\varepsilon}{\sigma} \leq \frac{4\sqrt{n\sqrt{d}D}\varepsilon}{\sigma}.$$

Since there can be circular dependencies like, e.g., $j_i = i'$ and $j_{i'} = i$, it might not be possible to reorder the integrals in (4) such that all terms become isolated as in (5). We resolve these dependencies by only considering a subset of the clusters. To make this more precise, consider a graph whose nodes are the clusters and that has a directed edge from \mathcal{C}_i to \mathcal{C}_j if $|B_i| = |B_{i,j}| = 1$, i.e., for every i with $|B_i| = 1$, there is an edge from \mathcal{C}_i to \mathcal{C}_{j_i} . This graph contains exactly k_- edges and we can identify a subset of $k' \geq k_-/2$ edges that is cycle-free. The subset \mathcal{C}' of clusters that we consider consists of the tails of these edges. Since every node in the graph has an out-degree of at most one, \mathcal{C}' consists of exactly k' clusters. For each cluster \mathcal{C}_i not contained in \mathcal{C}' , we replace the corresponding $\mathbb{I}_i(y_i, y_{j_i})$ by the trivial upper bound of 1. Without loss of generality, the identified subset \mathcal{C}' consists of the clusters $\mathcal{C}_1, \dots, \mathcal{C}_{k'}$ and it is topologically sorted in the sense that $i < j_i$ for all $i \in \{1, \dots, k'\}$. Given this, (4) can be bounded from above by

$$\frac{1}{\sigma^{k_>d}} \cdot \left(\frac{\varepsilon}{\sigma}\right)^{m'} \underbrace{\int_{\mathcal{D}} \cdots \int_{\mathcal{D}}}_{k^* - k' \text{ integrals}} \underbrace{\int_{\mathcal{D}} f_{g_{k'}}(y_{k'}) \cdot \mathbb{I}_{k'}(y_{k'}, y_{j_{k'}}) \cdots \int_{\mathcal{D}} f_{g_1}(y_1) \cdot \mathbb{I}_1(y_1, y_{j_1})}_{k' \text{ integrals}} dy_1 \dots dy_{k^*}.$$

Evaluating this formula from right to left according to Lemma 7 yields

$$(3D)^{dk} \cdot \frac{1}{\sigma^{k_>d}} \cdot \left(\frac{\varepsilon}{\sigma}\right)^{m'} \left(\frac{4\sqrt{n\sqrt{d}D}\varepsilon}{\sigma}\right)^{k'} = (3D)^{dk} \cdot \left(4\sqrt{n\sqrt{d}D}\right)^{k'} \frac{\varepsilon^{m'+k'/2}}{\sigma^{k_>d+m'+k'}} ,$$

where the term $(3D)^{dk}$ comes from the $k^* - k' \leq k$ integrals over $y_{k'+1}, \dots, y_{k^*}$: Each of these integrals is over the hypercube \mathcal{D} , which has a volume of $(2D + 1)^d \leq (3D)^d$. The definitions directly yield that $k_{>} \leq k$ and $m' + k' \leq m' + k_{=} \leq \ell$. Furthermore,

$$m' \geq \frac{\sum_{i, |B_i| > 1} |B_i|}{2} \quad \text{and} \quad k' \geq \frac{k_{=}}{2} = \frac{\sum_{i, |B_i| = 1} |B_i|}{2}$$

implies $m' + k'/2 \geq \ell/4$. Altogether, this yields the desired upper bound of

$$(3D)^{dk} \left(4\sqrt{n\sqrt{d}D}\right)^\ell \cdot \frac{\varepsilon^{\ell/4}}{\sigma^{kd+\ell}} = \left(\frac{3D}{\sigma}\right)^{kd} \cdot \left(\frac{32n^2 d D^2 \varepsilon}{\sigma^4}\right)^{\ell/4}$$

for the probability of the event $\mathcal{E} \wedge \neg \mathcal{F}$. □

Lemma 7. *Let $o \in \mathbb{R}^d$ and $p \in \mathbb{R}^d$ be arbitrary points and let r denote a random point chosen according to a d -dimensional normal distribution with arbitrary mean and standard deviation σ . Moreover, let $\ell \in \{0, \dots, n-1\}$, and let*

$$q = \frac{\ell}{\ell+1} \cdot p + \frac{1}{\ell+1} \cdot r$$

be a convex combination of p and r . Then the probability that r is

(i) at a distance of at most $\delta > 0$ from o and

(ii) at a distance of at most $\varepsilon > 0$ from the bisector of o and q

is bounded from above by

$$\frac{2\sqrt{n\delta\varepsilon}}{\sigma}.$$

Proof. For ease of notation, we assume that o is the origin of the coordinate system, i.e., $o = (0, \dots, 0)$. Due to rotational symmetry, we can also assume that $p = (0, p_2, 0, \dots, 0)$ for some $p_2 \in \mathbb{R}$. Let $r = (r_1, \dots, r_d)$, and assume that the coordinates r_2, \dots, r_d are fixed arbitrarily. We only exploit the randomness of the coordinate r_1 , which is a one-dimensional Gaussian random variable with standard deviation σ . The condition that r has a distance of at most ε from the bisector of o and q can be expressed algebraically as

$$\frac{q - o}{\|q - o\|} \cdot \left(\frac{q + o}{2} - r\right) \in [-\varepsilon, \varepsilon].$$

Since $o = (0, \dots, 0)$, this simplifies to

$$\frac{q}{\|q\|} \cdot \left(\frac{q}{2} - r\right) \in [-\varepsilon, \varepsilon] \iff q \cdot \left(\frac{q}{2} - r\right) \in [-\|q\|\varepsilon, \|q\|\varepsilon].$$

Since r_i is fixed for $i \neq 1$, also the coordinates q_i of q are fixed for $i \neq 1$. Setting $2\lambda = 1/(\ell+1)$ and exploiting that the first coordinate of p is 0, we can further rewrite the previous expression as

$$\begin{pmatrix} 2\lambda r_1 \\ q_2 \\ \vdots \\ q_d \end{pmatrix} \cdot \begin{pmatrix} (\lambda - 1)r_1 \\ q_2/2 - r_2 \\ \vdots \\ q_d/2 - r_d \end{pmatrix} \in [-\|q\|\varepsilon, \|q\|\varepsilon],$$

which is equivalent to

$$2\lambda(\lambda - 1)r_1^2 + \sum_{i=2}^d q_i(q_i/2 - r_i) \in [-\|q\|\varepsilon, \|q\|\varepsilon].$$

Since the coordinates q_i and r_i are fixed for $i \neq 1$, this implies that r can only be at distance ε from the bisector of p and q if r_1^2 falls into a fixed interval of length

$$\frac{2\|q\|\varepsilon}{2\lambda(1 - \lambda)} = \frac{2(\ell + 1)\|q\|\varepsilon}{1 - \frac{1}{2^{\ell+1}}} \leq 4n\|q\|\varepsilon.$$

As we only consider the event that q has a distance of at most δ from o , we can replace $\|q\|$ by δ in the expression above, leaving us with the problem to find an upper bound for the probability that the random variable r_1^2 assumes a value in a fixed interval of length at most $4n\delta\varepsilon$. For this to happen, r_1 has to assume a value in one of two intervals, each of length at most $2\sqrt{n\delta\varepsilon}$. Now r_1 follows a Gaussian distribution with standard deviation σ . This means that the corresponding density is bounded from above by $(\sqrt{2\pi}\sigma)^{-1}$. Thus, the probability of this event is at most

$$\frac{4\sqrt{n\delta\varepsilon}}{\sqrt{2\pi}\sigma} < \frac{2\sqrt{n\delta\varepsilon}}{\sigma}. \quad \square$$

3 An Upper Bound

Lemma 5 yields an upper bound on the number of iterations that k -means needs: Since there are only few points close to hyperplanes, eventually a point switches from one cluster to another that initially was not close to a hyperplane. The results of this section lead to the proof of Theorem 2 in Section 7.3.

First, we bound the number of iterations in terms of the distance Δ of the closest cluster centers that occur during the run of k -means.

Lemma 8. *For every $a \in [k]$, with a probability of at least $1 - 3W^{-1}$, every sequence of $k^{kd/a} + 1$ consecutive steps of the k -means algorithm (not including the first one) reduces the potential by at least*

$$\frac{\varepsilon^2 \cdot \min\{\Delta^2, 1\}}{36dD^2k^{kd/a}},$$

where Δ denotes the smallest distance of two cluster centers that occurs during the sequence and ε is defined as in Lemma 5.

Proof. Consider the configuration directly before the sequence of steps is performed. Due to Lemma 5, the probability that more than kd/a points are within distance ε of one of the bisectors is at most $2W^{-1}$. Additionally, only with a probability of at most W^{-1} there exists a point from \mathcal{X} that does not lie in the hypercube \mathcal{D} . Let us assume in the following that none of these failure events occurs, which is the case with a probability of at least $1 - 3W^{-1}$.

These points can assume at most $k^{kd/a}$ different configurations. Thus, during the considered sequence, at least one point that is initially not within distance ε of one of the bisectors must change its assignment. Let us call this point x , and let us assume that it changes from cluster \mathcal{C}_1 to cluster \mathcal{C}_2 . Furthermore, let δ be the distance of the centers of \mathcal{C}_1 and \mathcal{C}_2 before the sequence, and let c_1 and c_2 be the positions of the centers before the sequence. We distinguish

two cases. First, if x is closer to c_2 than to c_1 already in the beginning of the sequence, then x will change its assignment in the first step. Let $v = \frac{c_2 - c_1}{\|c_2 - c_1\|}$ be the unit vector in $c_2 - c_1$ direction. We have $c_2 - c_1 = \delta v$ and $(2x - c_1 - c_2) \cdot v = \alpha v$ for some $\alpha \geq 2\varepsilon$. Then the switch of x from \mathcal{C}_1 to \mathcal{C}_2 reduces the potential by at least

$$\|x - c_1\|^2 - \|x - c_2\|^2 = (2x - c_1 - c_2) \cdot (c_2 - c_1) \geq 2\varepsilon\delta \geq 2\varepsilon\Delta \geq \frac{\varepsilon^2 \cdot \min\{\Delta^2, 1\}}{4D^2 d k^{kd/a}},$$

where the last inequality follows from $\varepsilon \leq 1$ and $D \geq 1$. This completes the first case. Second, if x is closer to c_1 than to c_2 , then

$$\|x - c_2\|^2 - \|x - c_1\|^2 \geq 2\varepsilon\delta,$$

and hence,

$$\|x - c_2\| - \|x - c_1\| \geq \frac{2\varepsilon\delta}{\|x - c_2\| + \|x - c_1\|} \geq \frac{\varepsilon\delta}{36D\sqrt{d}}.$$

In this case, x can only change to cluster \mathcal{C}_2 after at least one of the centers of \mathcal{C}_1 or \mathcal{C}_2 has moved. Consider the centers of \mathcal{C}_1 and \mathcal{C}_2 immediately before the reassignment of x from \mathcal{C}_1 to \mathcal{C}_2 . Let c'_1 and c'_2 denote these centers. Then

$$\|x - c'_1\| - \|x - c'_2\| > 0.$$

By combining these observations with the triangle inequality, we obtain

$$\begin{aligned} \|c_1 - c'_1\| + \|c_2 - c'_2\| &\geq (\|x - c'_1\| - \|x - c_1\|) + (\|x - c_2\| - \|x - c'_2\|) \\ &= (\|x - c'_1\| - \|x - c'_2\|) + (\|x - c_2\| - \|x - c_1\|) \geq \frac{\varepsilon\delta}{\sqrt{d}3D}. \end{aligned}$$

This implies that one of the centers must have moved by at least $\varepsilon\delta/(6\sqrt{d}D)$ during the considered sequence of steps. Each time the center moves by some amount ξ , the potential drops by at least ξ^2 (see Lemma 17). Since this function is concave, the smallest potential drop is obtained if the center moves by $\varepsilon\delta/(6\sqrt{d}Dk^{kd/a})$ in each iteration. Then the decrease of the potential due to the movement of the center is at least

$$k^{kd/a} \cdot \left(\frac{\varepsilon\delta}{6\sqrt{d}Dk^{kd/a}} \right)^2 \geq \frac{\varepsilon^2\Delta^2}{36dD^2k^{kd/a}},$$

which concludes the proof. \square

In order to obtain a bound on the number of iterations that k -means needs, we need to bound the distance Δ of the closest cluster centers. This is done in the following lemma, which exploits Lemma 8. The following lemma is the crucial ingredient of the proof of Theorem 2.

Lemma 9. *Let $a \in [k]$ be arbitrary. Then the expected number of steps until the potential drops by at least 1 is bounded from above by*

$$\gamma \cdot k^{2kd/a} \cdot nkd \left(\frac{d^2 n^4 D}{\sigma\varepsilon} \right)^2$$

for a sufficiently large absolute constant γ .

Proof. With a probability of at least $1 - 3W^{-1}$, the number of iterations until the potential drops by at least

$$\frac{\varepsilon^2 \cdot \min\{\Delta^2, 1\}}{36dD^2k^{kd/a}}$$

is at most $k^{kd/a} + 1$ due to Lemma 8. We estimate the contribution of the failure event, which occurs only with probability $3W^{-1}$, to the expected running time by 3 and ignore it in the following. Let T denote the random variable that equals the number of sequences of length $k^{kd/a} + 1$ until the potential has dropped by one.

The random variable T can only exceed t if

$$\min\{\Delta^2, 1\} \leq \frac{36dD^2k^{kd/a}}{\varepsilon^2 \cdot t},$$

leading to the following bound on the expected value of T :

$$\begin{aligned} \mathbb{E}[T] &= \sum_{t=1}^W \Pr[T \geq t] \leq \int_0^W \Pr\left[\min\{\Delta^2, 1\} \leq \frac{36dD^2k^{kd/a}}{\varepsilon^2 \cdot t}\right] dt \\ &\leq t' + \int_{t'}^W \Pr\left[\Delta \leq \frac{6\sqrt{d}Dk^{kd/(2a)}}{\varepsilon \cdot \sqrt{t}}\right] dt, \end{aligned}$$

for

$$t' = \left(\frac{(24d + 96)n^4\sqrt{d}Dk^{kd/(2a)}}{\sigma\varepsilon}\right)^2.$$

Let us consider a situation reached by k -means in which there are two clusters \mathcal{C}_1 and \mathcal{C}_2 whose centers are at a distance of δ from each other. We denote the positions of these centers by c_1 and c_2 . Let H be the bisector between c_1 and c_2 . The points c_1 and c_2 are the centers of mass of the points assigned to \mathcal{C}_1 and \mathcal{C}_2 , respectively. From this, we can conclude the following: for every point that is assigned to \mathcal{C}_1 or \mathcal{C}_2 and that has a distance of at least δ from the bisector H , as compensation another point must be assigned to \mathcal{C}_1 or \mathcal{C}_2 that has a distance of at most $\delta/2$ from H . Hence, the total number of points assigned to \mathcal{C}_1 or \mathcal{C}_2 can be at most twice as large as the total number of points assigned to \mathcal{C}_1 or \mathcal{C}_2 that are at a distance of at most δ from H . Hence, there can only exist two centers at a distance of at most δ if one of the following two properties is met:

1. There exists a hyperplane from which more than $2d$ points have a distance of at most δ .
2. There exist two subsets of points whose union has cardinality at most $4d$ and whose centers of mass are at a distance of at most δ .

The probability that one of these events occurs can be bounded as follows using a union bound and Lemma 13 (see also Arthur and Vassilvitskii [3, Proposition 5.6]):

$$\Pr[\Delta \leq \delta] \leq n^{2d} \left(\frac{4d\delta}{\sigma}\right)^{2d-d} + (2n)^{4d} \cdot \left(\frac{\delta}{\sigma}\right)^d \leq \left(\frac{(4d + 16)n^4\delta}{\sigma}\right)^d.$$

Hence,

$$\Pr\left[\Delta \leq \frac{6\sqrt{d}Dk^{kd/(2a)}}{\varepsilon \cdot \sqrt{t}}\right] \leq \left(\frac{\sqrt{t'}}{\sqrt{t}}\right)^d$$

and, for $d \geq 3$, we obtain

$$\begin{aligned} \mathbb{E}[T] &\leq t' + \int_{t'}^W \left(\frac{\sqrt{t'}}{\sqrt{t}} \right)^d dt \\ &\leq t' + t'^{d/2} \left[\frac{1}{(-d/2 + 1) \cdot t^{d/2-1}} \right]_{t'}^{\infty} = \frac{d}{d-2} \cdot t' \leq 2\kappa n k d \cdot t'. \end{aligned}$$

For $d = 2$, we obtain

$$\mathbb{E}[T] \leq t' + \int_{t'}^W \left(\frac{\sqrt{t'}}{\sqrt{t}} \right)^d dt \leq t' + t' \cdot [\ln(t)]_1^W = t' \cdot (1 + \ln(W)) \leq 2\kappa n k d \cdot t'.$$

Altogether, this shows that the expected number of steps until the potential drops by at least 1 can be bounded from above by

$$2 + \left(k^{kd/a} + 1 \right) \cdot 2\kappa n k d \cdot \left(\frac{(24d + 96)n^4 \sqrt{d} D k^{kd/(2a)}}{\sigma \varepsilon} \right)^2,$$

which can, for a sufficiently large absolute constant γ , be bounded from above by

$$\gamma \cdot k^{2kd/a} \cdot n k d \cdot \left(\frac{d^2 n^4 D}{\sigma \varepsilon} \right)^2. \quad \square$$

4 Iterations with at most \sqrt{k} Active Clusters

In this and the following section, we aim at proving the main lemmas that lead to Theorem 1, which we will prove in Section 7.2. To do this, we distinguish two cases: In this section, we deal with the case that at most \sqrt{k} are active. In this case, either few points change clusters, which yields a potential drop caused by the movement of the centers. Or many points change clusters. Then, in particular, many points switch between two clusters, and not all of them can be close to the hyperplane bisecting the corresponding centers, which yields the potential drop in this case.

We define an *epoch* to be a sequence of consecutive iterations in which no cluster center assumes more than two different positions. Equivalently, there are at most two different sets $\mathcal{C}'_i, \mathcal{C}''_i$ that every cluster \mathcal{C}_i assumes. The obvious upper bound for the length of an epoch is 2^k , which is stated also by Arthur and Vassilvitskii [3]: After that many iterations, at least one cluster must have assumed a third position. For our analysis, however, 2^k is too big, and we bring it down to a constant.

Lemma 10. *The length of any epoch is less than four.*

Proof. Let x be any data point that changes from one cluster to another during an epoch, and let i_1, i_2, \dots, i_ℓ be the indices of the different clusters to which x belongs in that order. (We have $i_j \neq i_{j+1}$, but x can change back to a cluster it has already visited. So, e.g., $i_j = i_{j+2}$ is allowed.) For every i_j , we then have two different sets \mathcal{C}'_{i_j} and \mathcal{C}''_{i_j} with centers c'_{i_j} and c''_{i_j} such that $x \in \mathcal{C}''_{i_j} \setminus \mathcal{C}'_{i_j}$. Since x belongs always to at exactly one cluster, we have $\mathcal{C}_{i_j} = \mathcal{C}'_{i_j}$ for all except for one j for which $\mathcal{C}_{i_j} = \mathcal{C}''_{i_j}$. Now assume that $\ell \geq 4$. Then, when changing from \mathcal{C}_{i_1}

to \mathcal{C}_{i_2} , we have $\|x - c'_{i_2}\| < \|x - c'_{i_4}\|$ since x prefers \mathcal{C}_{i_2} over \mathcal{C}_{i_4} and, when changing to \mathcal{C}_{i_4} , we have $\|x - c'_{i_4}\| < \|x - c'_{i_2}\|$. This contradicts the assumption that $\ell \geq 4$.

Now assume that x does not change from \mathcal{C}_{i_j} to $\mathcal{C}_{i_{j+1}}$ for a couple of steps, i.e., x waits until it eventually changes clusters. Then the reason for eventually changing to $\mathcal{C}_{i_{j+1}}$ can only be that either \mathcal{C}_{i_j} has changed to some $\tilde{\mathcal{C}}_{i_j}$, which makes x prefer $\mathcal{C}_{i_{j+1}}$. But, since $\tilde{\mathcal{C}}_{i_j} \neq \mathcal{C}_{i_j}''$ and $x \in \tilde{\mathcal{C}}_{i_j}$, we have a third cluster for \mathcal{C}_{i_j} . Or $\mathcal{C}_{i_{j+1}}$ has changed to $\tilde{\mathcal{C}}_{i_{j+1}}$, and x prefers $\tilde{\mathcal{C}}_{i_{j+1}}$. But then $\tilde{\mathcal{C}}_{i_{j+1}} \neq \mathcal{C}_{i_{j+1}}'$ and $x \notin \tilde{\mathcal{C}}_{i_{j+1}}$, and we have a third cluster for $\mathcal{C}_{i_{j+1}}$.

We can conclude that x visits at most three different clusters, and changes its cluster in every iteration of the epoch. Furthermore, the order in which x visits its clusters is periodic with a period length of at most three. Finally, even a period length of three is impossible: Suppose x visits \mathcal{C}_{i_1} , \mathcal{C}_{i_2} , and \mathcal{C}_{i_3} . Then, to go from \mathcal{C}_{i_j} to $\mathcal{C}_{i_{j+1}}$ (arithmetic is modulo 3), we have $\|x - c'_{i_{j+1}}\| < \|x - c'_{i_{j-1}}\|$. Since this holds for $j = 1, 2, 3$, we have a contradiction.

This holds for every data point. Thus, after at most four iterations either k -means terminates, which is fine, or some cluster assumes a third configuration, which ends the epoch, or some clustering repeats, which is impossible. \square

Similar to Arthur and Vassilvitskii [3], we define a *key-value* to be an expression of the form $K = \frac{s}{t} \cdot \text{cm}(S)$, where $s, t \in \mathbb{N}$, $s \leq n^2$, $t < n$, and $S \subseteq \mathcal{X}$ is a set of at most $4d\sqrt{k}$ points. (Arthur and Vassilvitskii allow up to $4dk$ points.) For two key-values K_1, K_2 , we write $K_1 \equiv K_2$ if and only if they have identical coefficients for every data point.

We say that \mathcal{X} is δ -sparse if, for every key-values K_1, K_2, K_3, K_4 with $\|K_1 + K_2 - K_3 - K_4\| \leq \delta$, we have $K_1 + K_2 \equiv K_3 + K_4$.

Lemma 11. *The probability that the point set \mathcal{X} is not δ -sparse is at most*

$$n^{16d\sqrt{k}+12} \cdot \left(\frac{n^4\delta}{\sigma}\right)^d.$$

Proof. Let us first bound the number of possible key-values: There are at most n^3 possibilities for choosing s and t and $n^{4d\sqrt{k}}$ possibilities for choosing the set S . Thus, there are at most $n^{16d\sqrt{k}+12}$ possibilities for choosing four key-values K_1, \dots, K_4 . We fix K_1, \dots, K_4 arbitrarily. The rest follows from a union bound and the proof of Proposition 5.3 of Arthur and Vassilvitskii [3]. \square

After four iterations, one cluster has assumed a third center or k -means terminates. This yields the following lemma (see also Arthur and Vassilvitskii [3, Corollary 5.2]).

Lemma 12. *Assume that \mathcal{X} is δ -sparse. Then, in every sequence of four consecutive iterations that do not lead to termination and such that in every of these iterations*

- *at most \sqrt{k} clusters are active and*
- *each cluster gains or loses at most $2d\sqrt{k}$ points,*

the potential decreases by at least $\frac{\delta^2}{4n^4}$.

We say that \mathcal{X} is ε -separated if, for every hyperplane H , there are at most $2d$ points in \mathcal{X} that are within distance ε of H . The following lemma, due to Arthur and Vassilvitskii [3, Proposition 5.6], shows that \mathcal{X} is likely to be ε -separated.

Lemma 13 (Arthur and Vassilvitskii [3]). *The point set \mathcal{X} is not ε -separated with a probability of at most*

$$n^{2d} \cdot \left(\frac{4d\varepsilon}{\sigma} \right)^d.$$

Given that \mathcal{X} is ε -separated, every iteration with at most \sqrt{k} active clusters in which one cluster gains or loses at least $2d\sqrt{k}$ points yields a significant decrease of the potential.

Lemma 14. *Assume that \mathcal{X} is ε -separated. For every iteration with at most \sqrt{k} active clusters, the following holds: If a cluster gains or loses more than $2d\sqrt{k}$ points, then the potential drops by at least $2\varepsilon^2/n$.*

This lemma is similar to Proposition 5.4 of Arthur and Vassilvitskii [3]. We present here a corrected proof based on private communication with Vassilvitskii.

Proof. If a cluster \mathcal{C}_i gains or loses more than $2d\sqrt{k}$ points in a single iteration with at most \sqrt{k} active clusters, then there exists another cluster \mathcal{C}_j with which \mathcal{C}_i exchanges at least $2d+1$ points. Since \mathcal{X} is ε -separated, one of these points, say, x , must be at a distance of at least ε from the hyperplane bisecting the cluster centers c_i and c_j . Assume that x switches from \mathcal{C}_i to \mathcal{C}_j .

Then the potential decreases by at least $\|c_i - x\|^2 - \|c_j - x\|^2 = (2x - c_i - c_j) \cdot (c_j - c_i)$. Let v be the unit vector in $c_j - c_i$ direction. Then $(2x - c_i - c_j) \cdot v \geq 2\varepsilon$. We have $c_j - c_i = \alpha v$ for $\alpha = \|c_j - c_i\|$, and hence, it remains to bound $\|c_j - c_i\|$ from below. If we can prove $\alpha \geq \varepsilon/n$, then we have a potential drop of at least $(2x - c_i - c_j) \cdot \alpha v \geq \alpha 2\varepsilon \geq 2\varepsilon^2/n$ as claimed.

Let H be the hyperplane bisecting the centers of \mathcal{C}_i and \mathcal{C}_j in the previous iteration. While H does not necessarily bisect c_i and c_j , it divides the data points belonging to \mathcal{C}_i and \mathcal{C}_j correctly. In particular, this implies that $\|c_i - c_j\| \geq \text{dist}(c_i, H) + \text{dist}(c_j, H)$.

Consider the at least $2d+1$ data points switching between \mathcal{C}_i and \mathcal{C}_j . One of them must be at a distance of at least ε of H since \mathcal{X} is ε -separated. Let us assume w.l.o.g. that this point switches to \mathcal{C}_i . This yields $\text{dist}(c_i, H) \geq \varepsilon/n$ since \mathcal{C}_i contains at most n points. Thus, $\|c_i - c_j\| \geq \varepsilon/n$, which yields $\alpha \geq \varepsilon/n$ as desired. \square

Now set $\delta_i = n^{-16-(16+i)\sqrt{k}} \cdot \sigma$ and $\varepsilon_i = \sigma \cdot n^{-4-i\sqrt{k}}$. Then the probability that the instance is not δ_i -sparse is bounded from above by

$$n^{16d\sqrt{k}+12+4d-16d-(16+i)d\sqrt{k}} \leq n^{-id\sqrt{k}}.$$

The probability that the instance is not ε_i -separated is bounded from above by (we use $d \leq n$ and $4 \leq n$)

$$n^{4d-4d-id\sqrt{k}} = n^{-id\sqrt{k}}.$$

We abbreviate the fact that an instance is δ_i -sparse and ε_i -separated by i -nice. Now Lemmas 12 and 14 immediately yield the following lemma.

Lemma 15. *Assume that \mathcal{X} is i -nice. Then the number of sequences of at most four consecutive iterations, each of which with at most \sqrt{k} active clusters, until the potential has dropped by at least 1 is bounded from above by*

$$\left(\min \left\{ \frac{1}{4} \cdot n^{-36-(32+2i)\sqrt{k}} \cdot \sigma^2, 2\sigma^2 \cdot n^{-9-i2\sqrt{k}} \right\} \right)^{-1} \leq \frac{n^{(c+2i)\sqrt{k}}}{\sigma^2} =: S_i$$

for a suitable constant c .

The first term comes from δ_i , which yields a potential drop of at least $\delta_i^2/(4n^4)$. The second term comes from ε_i , which yields a drop of at least $2\varepsilon_i^2/n$.

Putting the pieces together yields the main lemma of this section.

Lemma 16. *The expected number of sequences of at most four consecutive iterations, each of which with at most \sqrt{k} active clusters, until the potential has dropped by at least 1 is bounded from above by*

$$\text{poly}\left(n^{\sqrt{k}}, \frac{1}{\sigma}\right).$$

Proof. The probability that it takes more than S_i such sequences is bounded from above by the probability that the instance is not i -nice, which is bounded from above by $2n^{-id\sqrt{k}}$. Let T be the random variable of the number of sequences of at most four consecutive iterations, each with at most \sqrt{k} active clusters, that it takes until we have a potential drop of at least 1.

We observe that k -means runs always at most $W \leq n^{\kappa kd}$ iterations. This yields that we have to consider i only up to $\kappa\sqrt{k}$. We assume further that $d \geq 2$. Putting all observations together yields the lemma:

$$\begin{aligned} \mathbb{E}[T] &\leq S_0 + \sum_{i=0}^{\kappa\sqrt{k}} \Pr[\text{not } i\text{-nice, but } (i+1)\text{-nice}] \cdot S_{i+1} + W \cdot \Pr[\text{not } \kappa\sqrt{k}\text{-nice}] \\ &\leq S_0 + \sum_{i=0}^{\kappa\sqrt{k}} \Pr[\text{not } i\text{-nice}] \cdot S_{i+1} + n^{\kappa kd} \cdot 2n^{-\kappa dk} \\ &\leq \frac{n^{c\sqrt{k}}}{\sigma^2} + \sum_{i=0}^{\kappa\sqrt{k}} 2n^{-id\sqrt{k}} \cdot \frac{n^{(c+2+2i)\sqrt{k}}}{\sigma^2} + 2 \\ &\leq \frac{n^{c\sqrt{k}}}{\sigma^2} + \sum_{i=0}^{\kappa\sqrt{k}} 2 \cdot \frac{n^{(c+2)\sqrt{k}}}{\sigma^2} + 2 \leq \text{poly}\left(n^{\sqrt{k}}, \frac{1}{\sigma}\right). \quad \square \end{aligned}$$

5 Iterations with at least \sqrt{k} Active Clusters

In this section, we consider steps of the k -means algorithm in which at least \sqrt{k} different clusters gain or lose points. The improvement yielded by such a step can only be small if none of the cluster centers changes its position significantly due to the reassignment of points, which, intuitively, becomes increasingly unlikely the more clusters are active. We show that, indeed, if at least \sqrt{k} clusters are active, then with high probability one of them changes its position by $n^{-O(\sqrt{k})}$, yielding a potential drop in the same order of magnitude.

The following observation, which has also been used by Arthur and Vassilvitskii [3], relates the movement of a cluster center to the potential drop.

Lemma 17. *If in an iteration of the k -means algorithm a cluster center changes its position from c to c' , then the potential drops by at least $\|c - c'\|^2$.*

Proof. The potential is defined as

$$\sum_{x \in \mathcal{X}} \|x - c_x\|^2,$$

where c_x denotes the center that is closest to x . We can rewrite this as

$$\sum_{c \in \mathcal{C}} \sum_{x \in X_c} \|x - c\|^2 = \sum_{c \in \mathcal{C}} \left(\sum_{x \in X_c} \|x - \text{cm}(X_c)\|^2 + |X_c| \cdot \|\text{cm}(X_c) - c\|^2 \right),$$

where $X_c \subseteq \mathcal{X}$ denotes all points from \mathcal{X} whose closest center is c and where $\text{cm}(X_c)$ denotes the center of mass of X_c .

Let us consider the case that one center changes its position from c to c' . Then c' must be the center of mass of X_c . Furthermore, $|X_c| \geq 1$. Hence, the potential drops by at least

$$\|\text{cm}(X_c) - c\|^2 - \|\text{cm}(X_c) - c'\|^2 = \|c' - c\|^2 - \|c' - c'\|^2 = \|c' - c\|^2. \quad \square$$

Now we are ready to prove the main lemma of this section.

Lemma 18. *The expected number of steps with at least \sqrt{k} active clusters until the potential drops by at least 1 is bounded from above by*

$$\text{poly} \left(n^{\sqrt{k}}, \frac{1}{\sigma} \right).$$

Proof. We consider one step of the k -means algorithm with at least \sqrt{k} active clusters. Let ε be defined as in Lemma 5 for $a = 1$. We distinguish two cases: Either one point that is reassigned during the considered iteration has a distance of at least ε from the bisector that it crosses, or all points are at a distance of at most ε from their respective bisectors. In the former case, we immediately get a potential drop of at least $2\varepsilon\Delta$, where Δ denotes the minimal distance of two cluster centers. In the latter case, Lemma 5 implies that with high probability less than kd points are reassigned during the considered step. We apply a union bound over the choices for these points. In the union bound, we fix not only these points but also the clusters they are assigned to before and after the step. We denote by A_i the set of points that are assigned to cluster \mathcal{C}_i in both configurations and we denote by B_i and B'_i the sets of points assigned to cluster \mathcal{C}_i before and after the step, respectively, except for the points in A_i . Analogously to Lemma 5, we assume that the positions of the points in $A_1 \cup \dots \cup A_k$ are fixed adversarially, and we apply a union bound on the different partitions A_1, \dots, A_k realizable. Altogether, we have a union bound over less than

$$n^{\kappa kd} \cdot n^{3kd} \leq n^{(\kappa+3) \cdot kd}$$

events. Let c_i be the position of the cluster center of \mathcal{C}_i before the reassignment, and let c'_i be the position after the reassignment. Then

$$c_i = \frac{|A_i| \cdot \text{cm}(A_i) + |B_i| \cdot \text{cm}(B_i)}{|A_i| + |B_i|},$$

where $\text{cm}(\cdot)$ denotes the center of mass of a point set. Since c'_i can be expressed analogously, we can write the change of position of the cluster center of \mathcal{C}_i as

$$c_i - c'_i = |A_i| \cdot \text{cm}(A_i) \left(\frac{1}{|A_i| + |B_i|} - \frac{1}{|A_i| + |B'_i|} \right) + \frac{|B_i| \cdot \text{cm}(B_i)}{|A_i| + |B_i|} - \frac{|B'_i| \cdot \text{cm}(B'_i)}{|A_i| + |B'_i|}.$$

Due to the union bound, $\text{cm}(A_i)$ and $|A_i|$ are fixed. Additionally, also the sets B_i and B'_i are fixed but not the positions of the points in these two sets. If we considered only a single center,

then we could easily estimate the probability that $\|c_i - c'_i\| \leq \beta$. For this, we additionally fix all positions of the points in $B_i \cup B'_i$ except for one of them, say b_i . Given this, we can express the event $\|c_i - c'_i\| \leq \beta$ as the event that b_i assumes a position in a ball whose position depends on the fixed values and whose radius, which depends on the number of points in $|A_i|$, $|B_i|$, and $|B'_i|$, is not larger than $n\beta$. Hence, the probability is bounded from above by

$$\left(\frac{n\beta}{\sigma}\right)^d.$$

However, we are interested in the probability that this is true for all centers simultaneously. Unfortunately, the events are not independent for different clusters. We estimate this probability by identifying a set of $\ell/2$ clusters whose randomness is independent enough, where $\ell \geq \sqrt{k}$ is the number of active clusters. To be more precise, we do the following: Consider a graph whose nodes are the active clusters and that contains an edge between two nodes if and only if the corresponding clusters exchange at least one point. We identify a dominating set in this graph, i.e., a subset of nodes that covers the graph in the sense that every node not belonging to this subset has at least one edge into the subset. We can assume that the dominating set, which we identify, contains at most half of the active clusters. (In order to find such a dominating set, start with the graph and throw out edges until the remaining graph is a tree. Then put the nodes on odd layers to the left side and the nodes on even layers to the right side, and take the smaller side as the dominating set.)

For every active center C that is not in the dominating set, we do the following: We assume that all the positions of the points in $B_i \cup B'_i$ are already fixed except for one of them. Given this, we can use the aforementioned estimate for the probability of $\|c_i - c'_i\| \leq \beta$. If we iterate this over all points not in the dominating set, we can always use the same estimate; the reason is that the choice of the subset guarantees that, for every node not in the subset, we have a point whose position is not fixed yet. This yields an upper bound of

$$\left(\frac{n\beta}{\sigma}\right)^{d\ell/2}.$$

Combining this probability with the number of choices in the union bound yields a bound of

$$n^{(\kappa+3)\cdot kd} \cdot \left(\frac{n\beta}{\sigma}\right)^{d\ell/2} \leq n^{(\kappa+3)\cdot kd} \cdot \left(\frac{n\beta}{\sigma}\right)^{d\sqrt{k}/2}.$$

For

$$\beta = \frac{\sigma}{n^{(4\kappa+6)\cdot\sqrt{k}+1}}$$

the probability can be bounded from above by $n^{-\kappa kd} \leq W^{-1}$.

Now we also take into account the failure probability of $2W^{-1}$ from Lemma 5. This yields that, with a probability of at least $1 - 3W^{-1}$, the potential drops in every iteration, in which at least \sqrt{k} clusters are active, by at least

$$\begin{aligned} \Gamma &:= \min\{2\varepsilon\Delta, \beta^2\} \geq \min\left\{\frac{\sigma^8\Delta}{1296n^{14+8\kappa}D^6d}, \frac{\sigma^2}{n^{(8\kappa+12)\cdot\sqrt{k}+2}}\right\} \\ &\geq \min\left\{\Delta \cdot \text{poly}(n^{-1}, \sigma), \text{poly}(n^{-\sqrt{k}}, \sigma)\right\} \end{aligned}$$

since $d \leq n$ and D is polynomially bounded in σ and n . The number T of steps with at least \sqrt{k} active clusters until the potential has dropped by one can only exceed t if $\Gamma \leq 1/t$. Hence,

$$\begin{aligned}
\mathbb{E}[T] &\leq \sum_{t=1}^{\infty} \Pr[T \geq t] + 3W^{-1} \cdot W \leq 3 + \int_{t=0}^{\infty} \Pr[T \geq t] dt \\
&\leq 4 + \int_{t=1}^{\infty} \Pr\left[\Gamma \leq \frac{1}{t}\right] dt \leq 4 + \beta^{-2} + \int_{t=\beta^{-2}}^{\infty} \Pr\left[\Gamma \leq \frac{1}{t}\right] dt \\
&\leq 4 + \beta^{-2} + \int_{t=\beta^{-2}}^{\infty} \Pr\left[\Delta \cdot \text{poly}\left(\frac{1}{n}, \sigma\right) \leq \frac{1}{t}\right] dt \\
&\leq 4 + \beta^{-2} + \int_{t=\beta^{-2}}^{\infty} \Pr\left[\Delta \leq \frac{1}{t} \cdot \text{poly}\left(n, \frac{1}{\sigma}\right)\right] dt \\
&\leq 4 + \beta^{-2} + \int_{t=\beta^{-2}}^{\infty} \min\left\{1, \left(\frac{(4d+16) \cdot n^4 \cdot \text{poly}(n, \sigma^{-1})}{t \cdot \sigma}\right)^d\right\} dt = \text{poly}\left(n^{\sqrt{k}}, \frac{1}{\sigma}\right),
\end{aligned}$$

where the integral is upper bounded as in the proof of Lemma 9. \square

6 A Polynomial Bound in One Dimension

In this section, we consider a one-dimensional set $\mathcal{X} \subseteq \mathbb{R}$ of points. The aim of this section is to prove that the expected number of steps until the potential has dropped by at least 1 is bounded by a polynomial in n and $1/\sigma$.

We say that the point set \mathcal{X} is ε -spreaded if the following conditions are fulfilled:

- There is no interval of length ε that contains three or more points of \mathcal{X} .
- For any four points x_1, x_2, x_3, x_4 , where x_2 and x_3 may denote the same point, we have $|x_1 - x_2| > \varepsilon$ or $|x_3 - x_4| > \varepsilon$.

The following lemma justifies the notion of ε -spreadedness.

Lemma 19. *Assume that \mathcal{X} is ε -spreaded. Then the potential drops by at least $\frac{\varepsilon^2}{4n^2}$ in every iteration.*

Proof. Let \mathcal{C}_i be the left-most active cluster, and let \mathcal{C}_j be the right-most active cluster.

We consider \mathcal{C}_i first. \mathcal{C}_i exchanges only points with the clusters to its right, for otherwise it would not be the leftmost active cluster. Thus, it cannot gain and lose points simultaneously. Assume that it gains points. Let A_i be the set of points of \mathcal{C}_i before the iteration, and let B_i be the set of points that it gains. Obviously, $\min B_i > \max A_i$. If $B_i \cup A_i$ contains at least three points, then we are done: If $|A_i| \geq 2$, then we consider the two rightmost points $x_1 \leq x_2$ of A_i and the leftmost point x_3 of B_i . These points are not within a common interval of size ε . Hence, x_3 has a distance of at least $\varepsilon/2$ from the center of mass $\text{cm}(A_i)$ because $\text{dist}(x_1, x_3) \geq \varepsilon$, $x_1 \leq x_2 \leq x_3$, and $\text{cm}(A_i) \leq (x_1 + x_2)/2$. Hence,

$$\text{cm}(B_i) \geq \text{cm}(A_i) + \frac{\varepsilon}{2}.$$

Thus, the cluster center moves to the right from $\text{cm}(A_i)$ to

$$\text{cm}(A_i \cup B_i) = \frac{|A_i| \cdot \text{cm}(A_i) + |B_i| \cdot \text{cm}(B_i)}{|A_i \cup B_i|} \geq \frac{|A_i \cup B_i| \cdot \text{cm}(A_i) + |B_i| \cdot \frac{\varepsilon}{2}}{|A_i \cup B_i|} \geq \text{cm}(A_i) + \frac{\varepsilon}{2n}.$$

The case $|A_i| = 1$ and $|B_i| \geq 2$ is analogous. The same holds if cluster \mathcal{C}_j switches from A_j to $A_j \cup B_j$ with $|A_j \cup B_j| \geq 3$, or if \mathcal{C}_i or \mathcal{C}_j lose points but initially have at least three points. Thus, in these cases, a cluster moves by at least $\varepsilon/(2n)$, which causes a potential drop by at least $\varepsilon^2/(4n^2)$.

It remains to consider the case that $|A_i \cup B_i| = 2 = |A_j \cup B_j|$. Thus, $A_i = \{a_i\}$, $B_i = \{b_i\}$, and also $A_j = \{a_j\}$, $B_j = \{b_j\}$. We restrict ourselves to the case that \mathcal{C}_i consists only of a_i and gains b_i and that \mathcal{C}_j has a_j and b_j and loses b_j . If only two clusters are active, we have $b_i = b_j$, and we have only three different points. Otherwise, all four points are distinct. This allows us to bring ε -spreadedness into play. We have either $|a_i - b_i| \geq \varepsilon$ or $|a_j - b_j| \geq \varepsilon$. But then either the center of \mathcal{C}_i or the center of \mathcal{C}_j moves by at least $\varepsilon/2$, which implies that the potential decreases by at least $\varepsilon^2/4 \geq \varepsilon^2/(4n^2)$. \square

Assume that \mathcal{X} is ε -spreaded. Then the number of iterations until the potential has dropped by at least 1 is at most $4n^2/\varepsilon^2$ by the lemma above. Let us estimate the probability that \mathcal{X} is ε -spreaded.

Lemma 20. *The probability that \mathcal{X} is not ε -spreaded is bounded from above by $\frac{2n^4\varepsilon^2}{\sigma^2}$.*

Proof. For the first property, let us consider any three points x_1, x_2, x_3 , and assume that x_1 is fixed arbitrarily. Then, in order to share an interval of size ε , we must have $|x_i - x_1| \leq \varepsilon$ for $i = 2, 3$. Since x_2 and x_3 are independent, this happens with a probability of at most $(\frac{2\varepsilon}{\sqrt{2\pi}\sigma})^2 \leq (\frac{\varepsilon}{\sigma})^2$. There are at most $n^3 \leq n^4$ choices for x_1, x_2, x_3 .

For the second property, consider any x_1, \dots, x_4 , and assume that x_2 and x_3 are fixed. Then the probability that $|x_1 - x_2| \leq \varepsilon$ and $|x_3 - x_4| \leq \varepsilon$ is at most $(\frac{\varepsilon}{\sigma})^2$. There are at most n^4 choices for x_1, \dots, x_4 .

Overall, by a union bound, the probability that \mathcal{X} is not ε -spreaded is at most $\frac{2n^4\varepsilon^2}{\sigma^2}$. \square

Now we have all ingredients for the proof of the main lemma of this section.

Lemma 21. *The number of iterations of k -means until the potential has dropped by at least 1 is bounded by a polynomial in n and $1/\sigma$.*

Proof. Let T be the random variable of the number of iterations until the potential has dropped by at least 1. If $T \geq t$, then \mathcal{X} cannot be ε -spreaded with $4n^2/\varepsilon^2 \leq t$. Thus, in this case, \mathcal{X} is not ε -spreaded with $\varepsilon = \frac{2n}{\sqrt{t}}$. In the worst case, k -means runs for at most $n^{\kappa k}$ iterations. Hence,

$$\begin{aligned} \mathbb{E}[T] &= \sum_{t=1}^{n^{\kappa k}} \Pr[T \geq t] \leq \sum_{t=1}^{n^{\kappa k}} \Pr\left[\mathcal{X} \text{ is not } \frac{2n}{\sqrt{t}}\text{-spreaded}\right] \\ &\leq \sum_{t=1}^{n^{\kappa k}} \frac{8n^4 n^2}{t\sigma^2} \in O\left(\frac{n^6}{\sigma^2} \cdot \log n^{\kappa k}\right) \subseteq O\left(\frac{n^7}{\sigma^2} \cdot \log n\right). \quad \square \end{aligned}$$

Finally, we remark that, by choosing $\varepsilon = \frac{\sigma}{n^{2+c}}$, we obtain that the probability that the number of iterations until the potential has dropped by at least exceeds a polynomial in n and $1/\sigma$ is bounded from above by $O(n^{-2c})$. This yields a bound on the running-time of k -means for $d = 1$ that holds with high probability.

7 Putting the Pieces Together

In the previous sections, we have only analyzed the expected number of iterations until the potential drops by at least 1. To bound the expected number of iterations that k -means needs to terminate, we need an upper on the potential in the beginning. To get this, we use the following lemma.

Lemma 22. *Let x be a one-dimensional Gaussian random variable with standard deviation σ and mean $\mu \in [0, 1]$. Then, for all $t \geq 1$,*

$$\Pr[x \notin [-t, 1 + t]] < \sigma \cdot \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

For $D = \sqrt{2\sigma^2 \ln(n^{1+\kappa kd} d \sigma)} \leq \text{poly}(n, \sigma)$, the probability that any component of any of the n data points is not contained in the hypercube $\mathcal{D} = [-D, 1 + D]^d$ is bounded from above by $n^{-\kappa kd} \leq W^{-1}$. This implies that $\mathcal{X} \subseteq \mathcal{D}$ with a probability of at least $1 - W^{-1}$.

In the beginning, we made the assumption that $\sigma \leq 1$. While this covers the small values of σ , which we consider as more relevant, the assumption is only a technical requirement, and we can get rid of it: The number of iterations that k -means needs is invariant under scaling of the point set \mathcal{X} . Now assume that $\sigma > 1$. Then we consider \mathcal{X} scaled down by $1/\sigma$, which corresponds to the following model: The adversary chooses points from the hypercube $[0, 1/\sigma]^d \subseteq [0, 1]^d$, and then we add d -dimensional Gaussian vectors with standard deviation 1 to every data point. The expected running-time that k -means needs on this instance is bounded from above by the running-time needed for adversarial points chosen from $[0, 1]^d$ and $\sigma = 1$, which is $\text{poly}(n) \leq \text{poly}(n, 1/\sigma)$.

7.1 Proof of Theorem 4

We obtain a bound that is polynomial in n and $1/\sigma$ from Lemmas 21 and 22: First, after one iteration, the potential is bounded from above by $\text{poly}(n, D) = \text{poly}(n)$. If this is not the case, we bound the number of iterations by W , which adds $W \cdot W^{-1}$ to the expected number of iterations. Second, the expected number of iterations until the potential has dropped by at least 1 is bounded by $\text{poly}(n, 1/\sigma)$, which yields a bound of $\text{poly}(n, 1/\sigma)$ until the algorithm terminates. This proves Theorem 4.

The result that the probability that the number of iterations exceeds a polynomial in n and $1/\sigma$ is at most $O(1/\text{poly}(n))$ follows immediately.

7.2 Proof of Theorem 1

In the remainder of this section, we restrict ourselves to $d \geq 2$. For $d = 1$, we already have a polynomial bound according to Theorem 4 and Section 6.

After $\text{poly}(n^{\sqrt{k}}, 1/\sigma)$ iterations, we have

- at least $\text{poly}(n^{\sqrt{k}}, 1/\sigma)$ sequences of four consecutive iterations, each of which with at most \sqrt{k} active clusters, or
- at least $\text{poly}(n^{\sqrt{k}}, 1/\sigma)$ iterations with at least \sqrt{k} active clusters.

Thus, by Lemmas 16 and 18, the expected number of steps until the potential has dropped by one is at most $\text{poly}(n^{\sqrt{k}}, 1/\sigma)$.

After the first iteration, the potential is at most $nd \cdot (2D + 1)^2$. As argued above, we can restrict ourselves to $\sigma \leq 1$, which implies $D \leq \text{poly}(n)$. This yields that the expected number of steps until termination is at most

$$nd \cdot (2D + 1)^2 \cdot \text{poly}\left(n^{\sqrt{k}}, \frac{1}{\sigma}\right) = \text{poly}\left(n^{\sqrt{k}}, \frac{1}{\sigma}\right),$$

provided that $\mathcal{X} \subseteq \mathcal{D}$. If $\mathcal{X} \not\subseteq \mathcal{D}$, we bound the number of iterations by the worst-case bound of W , which contributes only $W^{-1} \cdot W = 1$ to the expected number of iterations. This proves Theorem 1.

7.3 Proofs of Theorem 2 and Corollary 3

We exploit Lemma 9 with $a = 2$. Then the expected number of iterations until the potential has dropped by at least 1 is bounded from above by

$$k^{kd} \cdot \text{poly}\left(n, \frac{1}{\sigma}\right).$$

Again, after the first iteration, the potential is at most $nd \cdot (2D + 1)^2 \leq \text{poly}(n)$ with a probability of at least $1 - W^{-1}$. This shows that the expected number of iterations, provided $\mathcal{X} \subseteq \mathcal{D}$, is bounded from above by

$$k^{kd} \cdot \text{poly}\left(n, \frac{1}{\sigma}\right).$$

The event $\mathcal{X} \not\subseteq \mathcal{D}$ contributes only 1 to the expected number of iterations as argued in the previous section. This completes the proof of Theorem 2.

If $k, d \in O(\sqrt{\log n / \log \log n})$, then $k^{kd} \leq \text{poly}(n)$, which proves Corollary 3.

8 Conclusions

We have proved two upper bounds for the smoothed running-time of the k -means method: The first bound is $\text{poly}(n^{\sqrt{k}}, 1/\sigma)$. The second bound is $k^{kd} \cdot \text{poly}(n, 1/\sigma)$, which decouples the exponential growth in k and d from the number of points and the standard deviation. In particular, this yields a smoothed running-time that is polynomial in n and $1/\sigma$ for $k, d \in O(\sqrt{\log n / \log \log n})$.

The obvious question now is whether a bound exists that is polynomial in n and $1/\sigma$, without exponential dependence on k or d . We believe that such a bound exists. However, we suspect that new techniques are required to prove it; bounding the smallest possible improvement from below might not be sufficient. The reason for this is that the number of possible partitions, and thus the number of possible k -means steps, grows exponentially in k , which makes it more likely for small improvements to exist as k grows.

Finally, we are curious if our techniques carry over to other heuristics. In particular Lemma 5 is quite general, as it bounds the number of points from above that are close to the boundaries of the Voronoi partitions that arise during the execution of k -means. In fact, we believe that a slightly weaker version of Lemma 5 is also true for arbitrary Voronoi partitions and not only for those arising during the execution of k -means. This insight might turn out to be helpful in other contexts as well.

Acknowledgement

We thank David Arthur, Dan Spielman, Shang-Hua Teng, and Sergei Vassilvitskii for fruitful discussions and comments.

References

- [1] David Arthur. Smoothed analysis of the k -means method. Manuscript, 2008.
- [2] David Arthur and Sergei Vassilvitskii. How slow is the k -means method? In Nina Amenta and Otfried Cheong, editors, *Proc. of the 22nd ACM Symposium on Computational Geometry (SOCG)*, pages 144–153. ACM Press, 2006.
- [3] David Arthur and Sergei Vassilvitskii. Worst-case and smoothed analysis of the ICP algorithm, with an application to the k -means method. In *Proc. of the 47th Ann. IEEE Symp. on Foundations of Computer Science (FOCS)*, pages 153–164. IEEE Computer Society, 2006.
- [4] Mihai Bădoiu, Sariel Har-Peled, and Piotr Indyk. Approximate clustering via core-sets. In *Proc. of the 34th Ann. ACM Symposium on Theory of Computing (STOC)*, pages 250–257. ACM Press, 2002.
- [5] Pavel Berkhin. Survey of clustering data mining techniques. Technical report, Accrue Software, San Jose, CA, USA, 2002.
- [6] Richard O. Duda, Peter E. Hart, and David G. Stork. *Pattern Classification*. John Wiley & Sons, 2000.
- [7] Sariel Har-Peled and Bardia Sadri. How fast is the k -means method? *Algorithmica*, 41(3):185–202, 2005.
- [8] Mary Inaba, Naoki Katoh, and Hiroshi Imai. Variance-based k -clustering algorithms by Voronoi diagrams and randomization. *IEICE Transactions on Information and Systems*, E83-D(6):1199–1206, 2000.
- [9] Amit Kumar, Yogish Sabharwal, and Sandeep Sen. A simple linear time $(1 + \varepsilon)$ -approximation algorithm for k -means clustering in any dimensions. In *Proc. of the 45th Ann. IEEE Symp. on Foundations of Computer Science (FOCS)*, pages 454–462, 2004.
- [10] Stuart P. Lloyd. Least squares quantization in PCM. *IEEE Transactions on Information Theory*, 28(2):129–137, 1982.
- [11] Jiří Matoušek. On approximate geometric k -clustering. *Discrete and Computational Geometry*, 24(1):61–84, 2000.
- [12] Daniel A. Spielman and Shang-Hua Teng. Smoothed analysis of algorithms: Why the simplex algorithm usually takes polynomial time. *Journal of the ACM*, 51(3):385–463, 2004.

- [13] Daniel A. Spielman and Shang-Hua Teng. Smoothed analysis of algorithms and heuristics: Progress and open questions. In Luis M. Pardo, Allan Pinkus, Endre Süli, and Michael J. Todd, editors, *Foundations of Computational Mathematics, Santander 2005*, pages 274–342. Cambridge University Press, 2006.