

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Problem Reformulation and Mathematical Tools</b>	<b>5</b>
2.1	Reformulation of the Clustering Problem . . . . .	5
2.2	Introduction to the PALM Theory . . . . .	5
<b>3</b>	<b>Clustering: The Squared Euclidean Norm Case</b>	<b>8</b>
3.1	Clustering with PALM . . . . .	8
<b>4</b>	<b>Clustering: The Euclidean Norm Case</b>	<b>14</b>
4.1	A Smoothed Clustering Problem . . . . .	14
4.2	Different Approach Towards Solving the Smoothed $H\varepsilon$ . . . . .	23
<b>5</b>	<b>Returning to KMEANS</b>	<b>27</b>
5.1	Similarity to KMEANS . . . . .	27
5.2	KMEANS Local Minima Convergence Proof . . . . .	29
<b>6</b>	<b>Numeric Results</b>	<b>31</b>
6.1	Iris Dataset . . . . .	31
6.2	Synthetic Dataset . . . . .	32

# 1 Introduction

The clustering problem is the task of grouping objects which are similar. It consists of partitioning a dataset into subsets, called clusters, such that the data points in each cluster are similar with respect to a specific criteria.

The clustering problem is a fundamental problem in machine learning field, and it arises in wide scope of applications, such as data mining, pattern recognition, information retrieval and many others. For example, in image segmentation, one is interested in partitioning the pixels of an image into objects, where each pixel can be described via its location in the image and its color given in RGB format. Another example is learning the probability density of some data, where the data is assumed to be drawn from a mixtures of distributions. Each partition of the data is represented by a unimodal probability density model, and summing of all the cluster models gives a multimodal density for the entire dataset. Vector quantization is yet another example, where large sets of points are represented by their centroid point. This method can be used for data compression, data correction and pattern recognition.

There are several categories of clustering methods, each has a direct impact of the final clustering structure.

- (i) Hierarchical versus partitioning clustering: In partitioning clustering the dataset is divided into clusters, whereas in hierarchical clustering each cluster may have sub-clusters, thus forming a tree which leaves are the single points of the dataset.
- (ii) Hard versus soft and fuzzy clustering: In hard clustering each data point is assigned to single cluster, versus a soft clustering where each point may be assigned to more than one cluster, hence clusters may overlap. In fuzzy clustering for each point there is a distribution that describes the probability of a point to be part of a cluster.
- (iii) Complete versus partial clustering: In complete clustering all points in the dataset are assigned to clusters, whereas in partial clustering some points may be intentionally skipped and are not being assigned to a cluster.

Finding the optimal partition of a fixed number of clusters for some given dataset is known to be a NP-hard problem, and hence cannot be solved efficiently. Most algorithms seek to minimize some mathematical criteria, and usually achieve local rather than global minimum solution. In this work we focus on partitioning clustering, where the number of clusters is known in advance. Most partitioning clustering methods iteratively update the cluster centers, and hence they are often referred as center-based clustering methods. We introduce few notations for the upcoming discussion. Let  $\mathcal{A} = \{a^1, a^2, \dots, a^m\}$  be a given set of points in  $\mathbb{R}^n$ , and let  $1 < k < m$  be a fixed given number of clusters. The clustering problem consists of partitioning the dataset  $\mathcal{A}$  into  $k$  subsets  $\{C^1, C^2, \dots, C^k\}$ , called clusters. For each  $l = 1, 2, \dots, k$ , the cluster  $C^l$  is represented by its center  $x^l \in \mathbb{R}^n$ . We describe several well-known center-based clustering algorithms.

- (i) K-means algorithm: This algorithm is probably the most famous within the clustering scope, and dates back to MacQueen (1967). The k-means algorithm partitions the

data into  $k$  sets. The solution is then a set of  $k$  centers, each of which is located at the centroid of the data for which it is the closest center. The k-means algorithm performs hard clustering, and each point is labeled according to its closest center. The objective function that the k-means algorithm minimizes is

$$f_{KM}(x) = \sum_{i=1}^m \min_{1 \leq l \leq k} \|a^i - x^l\|^2.$$

The simplicity of the algorithm both in the theoretical and implementation aspects made it very popular.

- (ii) Fuzzy k-means (FKM) algorithm: The FKM algorithm is a soft clustering method. For each data point the result of the FKM algorithm is a distribution of membership over the clusters. The objective function that the FKM algorithm minimizes is

$$f_{FKM}(x) = \sum_{i=1}^m \sum_{l=1}^k (w_l^i)^\beta \|a^i - x^l\|^2.$$

The parameter  $w_l^i$  denotes the probability that data point  $a^i$  is assigned to cluster  $x^l$ , thus it is under the constraints  $\sum_{l=1}^k w_l^i = 1$  for all  $1 \leq i \leq m$  and  $w_l^i \geq 0$ . The parameter  $\beta > 1$  governs the "fuzzy partition". Setting  $\beta = 1$  results in the standard k-means algorithm.

- (iii) Expectation-Maximization (EM) algorithm: The EM algorithm is used extensively in statistical estimation problems for learning mixtures of distributions. It is a soft clustering algorithm. The objective function that EM maximizes is

$$f_{EM} = \sum_{i=1}^m \log \left( \sum_{l=1}^k p(a^i | x^l) p(x^l) \right),$$

where  $p(a^i | x^l)$  is the probability of  $a^i$  given that it is generated by the Gaussian distribution with center  $x^l$  and  $p(x^l)$  is the prior probability of center  $x^l$ .

An interesting paper of Teboulle [8] shows that these center-based clustering algorithms can be recovered from the proposed continuous optimization framework. The smoothing methodologies for the clustering problem are based on nonlinear means and on approximation of appropriate asymptotic functions.

Most of the existing clustering methods are sensitive to the starting point, namely choosing different starting point result in significant changes in the final clustering. There are plethora of heuristic initialization method. One such initialization method is choosing random  $k$  data points as a starting centers for clustering, assuming uniform distribution or some other prior distribution on the data. Another popular method is k-means++, where the first center is chosen at random from the dataset, and for each  $2 \leq l \leq k$ , the center  $x^l$  is the furthest point from the points chosen so far.

The starting point of our work is a formulation of the clustering problem which consists of minimizing the sum of finite collection of min-functions, which is a nonsmooth and nonconvex optimization problem, in its most the general case. The clustering problem is given by

$$\min_{x \in \mathbb{R}^{nk}} \left\{ F(x) := \sum_{i=1}^m \min_{1 \leq l \leq k} d(x^l, a^i) \right\}, \quad (1.1)$$

with  $d(\cdot, \cdot)$  being a distance-like function.

The lack of smoothness in this formulation can be overcome, yet the nonconvex nature of the clustering problem shall accompany the discussions throughout this work. Significant amount of studies have been made on convex models, even though in many cases the original optimization problem is nonconvex. To overcome the lack of convexity the common approach is usually achieved by relaxation of the original problem. Motivated by papers of Attouch et al. [1, 2] that established convergence of the sequences generated by the proximal Gauss-Seidel scheme in the general nonconvex and nonsmooth settings, and similar result for the proximal-forward-backward algorithm applied to the nonconvex and nonsmooth minimization of the sum of a nonsmooth function with a smooth one. This approach assumes that the objective function to be minimized satisfies the Kurdyka-Łojasiewicz (KL) property. The convergence results were further extended in the recent work by Bolte et al. [5], where the objective function is a function of finite blocks of variables.

We focus on two cases of distance-like functions. The first is the squared Euclidean norm, which is the standard proximity measure used by k-means. For this case, we derive an equivalent smooth optimization problem for the clustering problem presented in (1.1) and prove our convergence result for the suggested algorithm via the methodology which is discussed in [5]. The second distance-like function that we study is the Euclidean norm. In this case we present two approximations, in order to overcome the lack of smoothness in the problem, and then proceed with the same methodology as in the first case. We present numeric experiments, that show the superiority of the Euclidean norm distance function for datasets in which the data points are spread relatively sparsely from their centers.

**Outline of the thesis.** This work is organized as follows. In the following section we transform the initial formulation of the clustering problem into a smooth one. In addition, we introduce the KL theory and the general methodology that will be used in our analysis of the proposed algorithms. In section 3 and 4 we suggest two algorithms, KPLAM and  $\varepsilon$ -KPALM, respectively. KPALM addresses the clustering problem with squared Euclidean norm distance-like function and  $\varepsilon$ -KPALM the standard Euclidean norm distance function. Employing the methodology of Section 2 we establish our convergence results. In section 5 we prove the convergence of k-means algorithm to a critical point and under certain assumption extend the convergence to a local minimum. Finally, in Section 6 we compare the performance of these algorithm according to some common criteria.

## 2 Problem Reformulation and Mathematical Tools

### 2.1 Reformulation of the Clustering Problem

We begin with a reformulation of the clustering problem which will be the basis for our developments in this work. The reformulation is based on the following fact:

$$\min_{1 \leq l \leq k} u_l = \min \{ \langle u, v \rangle : v \in \Delta \},$$

where  $\Delta$  denotes the well-known simplex defined by

$$\Delta = \left\{ u \in \mathbb{R}^k : \sum_{l=1}^k u_l = 1, u \geq 0 \right\}.$$

Using this fact in Problem (1.1) and introducing new variables  $w^i \in \mathbb{R}^k$ ,  $i = 1, 2, \dots, m$ , gives a smooth reformulation of the clustering problem

$$\min_{x \in \mathbb{R}^{nk}} \sum_{i=1}^m \min_{w^i \in \Delta} \langle w^i, d^i(x) \rangle, \quad (2.1)$$

where

$$d^i(x) = (d(x^1, a^i), d(x^2, a^i), \dots, d(x^k, a^i)) \in \mathbb{R}^k, \quad i = 1, 2, \dots, m.$$

Replacing further the constraint  $w^i \in \Delta$  by adding the indicator function  $\delta_\Delta(\cdot)$ , which is defined to be 0 in  $\Delta$  and  $\infty$  otherwise, to the objective function, results in a equivalent formulation

$$\min_{x \in \mathbb{R}^{nk}, w \in \mathbb{R}^{km}} \left\{ \sum_{i=1}^m (\langle w^i, d^i(x) \rangle + \delta_\Delta(w^i)) \right\}, \quad (2.2)$$

where  $w = (w^1, w^2, \dots, w^m) \in \mathbb{R}^{km}$ . Finally, for the simplicity of the yet to come expositions, we define the following functions

$$H(w, x) := \sum_{i=1}^m H^i(w, x) = \sum_{i=1}^m \langle w^i, d^i(x) \rangle \quad \text{and} \quad G(w) = \sum_{i=1}^m G^i(w^i) := \sum_{i=1}^m \delta_\Delta(w^i).$$

Replacing the terms in Problem (2.2) with the functions defined above gives a compact equivalent form of the original clustering problem

$$\min \{ \Psi(z) := H(w, x) + G(w) \mid z := (w, x) \in \mathbb{R}^{km} \times \mathbb{R}^{nk} \}. \quad (2.3)$$

### 2.2 Introduction to the PALM Theory

In this subsection we give a brief review of the main developments established in [5]. These developments which include the proximal alternating linearized minimization (PALM) algorithm and general procedure for proving global convergence of generic algorithm play a central rule in this work. First, let us recall several definitions which are needed for the upcoming discussion.

**Definition 1** (Subdifferentials). Let  $\sigma : \mathbb{R}^d \rightarrow (-\infty, +\infty]$  be a proper and lower semicontinuous function.

- (i) For a given  $x \in \text{dom } \sigma := \{x \in \mathbb{R}^d : \sigma(x) < \infty\}$ , the Fréchet subdifferential of  $\sigma$  at  $x$ , written  $\widehat{\partial}\sigma(x)$ , is the set of all vectors  $u \in \mathbb{R}^d$  which satisfy

$$\liminf_{y \neq x, y \rightarrow x} \frac{\sigma(y) - \sigma(x) - \langle u, y - x \rangle}{\|y - x\|} \geq 0.$$

When  $x \notin \text{dom } \sigma$ , we set  $\widehat{\partial}\sigma(x) = \emptyset$ .

- (ii) The limiting-subdifferential, or subdifferential in short, of  $\sigma$  at  $x \in \mathbb{R}^n$ , written  $\partial\sigma(x)$ , is defined through the following closure process

$$\partial\sigma(x) := \left\{ u \in \mathbb{R}^d : \exists x^k \rightarrow x, \sigma(x^k) \rightarrow \sigma(x) \text{ and } u^k \in \widehat{\partial}\sigma(x^k) \text{ as } k \rightarrow \infty \right\}.$$

In the nonsmooth context, as in the smooth case, the well-known Fermat's rule remains unchanged, that is, if  $x \in \mathbb{R}^d$  is a local minimizer of  $\sigma$  then  $0 \in \partial\sigma(x)$ . Points whose subdifferential contains 0 are called *critical points*, and the set of all critical points of  $\sigma$  is denoted by  $\text{crit}\sigma$ .

Now we present the Kurdyka-Łojasiewicz property, which plays a central role in PALM's analysis. Let  $\eta \in (0, +\infty]$ . Denote the following class of functions

$$\Phi_\eta = \left\{ \varphi \in C([0, \eta], \mathbb{R}_+), \text{concave} : \varphi \in C^1((0, \eta)), \varphi' > 0, \varphi(0) = 0 \right\}.$$

**Definition 2** (Kurdyka-Łojasiewicz property). Let  $\sigma : \mathbb{R}^d \rightarrow (-\infty, +\infty]$  be proper and lower semicontinuous.

- (i) The function  $\sigma$  is said to have the Kurdyka-Łojasiewicz (KL) property at  $\bar{u} \in \text{dom } \partial\sigma := \{u \in \mathbb{R}^d : \partial\sigma \neq \emptyset\}$  if there exist  $\eta \in (0, +\infty]$ , a neighborhood  $U$  of  $\bar{u}$  and a function  $\varphi \in \Phi_\eta$ , such that for all

$$u \in U \cap \{x \in \mathbb{R}^d : \sigma(\bar{u}) < \sigma(x) < \sigma(\bar{u}) + \eta\},$$

the following inequality holds

$$\varphi'(\sigma(u) - \sigma(\bar{u})) \text{dist}(0, \partial\sigma(u)) \geq 1,$$

where  $\text{dist}(x, S) := \inf \{\|y - x\| : y \in S\}$  denotes the distance from  $x \in \mathbb{R}^d$  to  $S \subset \mathbb{R}^d$ .

- (ii) If  $\sigma$  satisfy the KL property at each point of  $\text{dom } \sigma$  then  $\sigma$  is called a KL function.

**Definition 3** (Semi-algebraic sets and functions). (i) A subset  $S \subset \mathbb{R}^d$  is a real semi-algebraic set if there exists a finite number of real polynomial functions  $g_{ij}, h_{ij} : \mathbb{R}^d \rightarrow \mathbb{R}$  such that

$$S = \bigcup_{j=1}^p \bigcap_{i=1}^q \{u \in \mathbb{R}^d : g_{ij} = 0 \text{ and } h_{ij}(u) < 0\}$$

(ii) A function  $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$  is called *semi-algebraic* if its graph

$$\{(u, t) \in \mathbb{R}^{d+1} : h(u) = t\}$$

is a semi-algebraic subset of  $\mathbb{R}^{d+1}$ .

**Theorem 1.** Let  $\sigma : \mathbb{R}^d \rightarrow (-\infty, +\infty]$  be a proper and lower semicontinuous function. If  $\sigma$  is semi-algebraic then it satisfies the KL property at any point of  $\text{dom } \sigma$ .

The class of semi-algebraic function is very broad, it includes real polynomial functions; indicator functions of semi-algebraic sets; finite sums and products of semi-algebraic functions; composition of semi-algebraic functions, and many more.

Equipped with these definitions, we present the methodology that was used for PALM algorithm, and will be used several times throughout this work. Let  $\Psi : \mathbb{R}^d \rightarrow (-\infty, +\infty]$  be a proper and lower semicontinuous function which is bounded from below and consider the problem

$$(P) \quad \{\Psi(z) : z \in \mathbb{R}^d\}.$$

Suppose that we are given a generic algorithm  $\mathcal{A}$  which generates a sequence  $\{z^k\}_{k \in \mathbb{N}}$  via the following scheme:

$$z^0 \in \mathbb{R}^d, z^{k+1} \in \mathcal{A}(z^k), \quad k = 0, 1, \dots$$

Building upon [1, 2], the following three requirements are sufficient to assure the convergence of the whole sequence  $\{z^k\}_{k \in \mathbb{N}}$  to a critical point of  $\Psi$ .

(C1) *Sufficient decrease property:* There exists a positive constant  $\rho_1$ , such that

$$\rho_1 \|z^{k+1} - z^k\|^2 \leq \Psi(z^k) - \Psi(z^{k+1}), \quad \forall k = 0, 1, \dots$$

(C2) *A subgradient lower bound for iterates gap:* Assuming that the sequence generated by the algorithm  $\mathcal{A}$  is bounded. There exists a positive constant  $\rho_2$ , such that

$$\|w^{k+1}\| \leq \rho_2 \|z^{k+1} - z^k\|, \quad w^k \in \partial \Psi(z^k) \quad \forall k = 0, 1, \dots$$

(C3) *KL property:* The function  $\Psi$  is a KL function.

Due to Theorem 1, property (C3) follows whenever  $\Psi$  is a semi-algebraic function. Finally, we present the proximal alternating linearized minimization (PALM) algorithm which solves the nonconvex and nonsmooth minimization problem of the following form

$$(M) \quad \text{minimize } \Psi(x, y) := f(x) + g(y) + H(x, y) \text{ over all } (x, y) \in \mathbb{R}^n \times \mathbb{R}^m,$$

where  $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$  and  $g : \mathbb{R}^m \rightarrow (-\infty, +\infty]$  are proper and lower semicontinuous functions while  $H : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$  is a  $C^1$  function. In addition, partial gradients of  $H$  are Lipschitz continuous, namely,  $H(\cdot, y) \in C_{L_1(y)}^{1,1}$  and  $H(x, \cdot) \in C_{L_2(x)}^{1,1}$ . PALM is alternating the steps of the proximal forward-backward (PFB) scheme. PFB scheme tackles the problem of minimizing the sum of a smooth function  $h$  with a nonsmooth one  $\sigma$ , and can be viewed as the proximal regularization of  $h$  linearized at a given point  $x^k$ , that is,

$$x^{k+1} \in \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \langle x - x^k, \nabla h(x^k) \rangle + \frac{t}{2} \|x - x^k\|^2 + \sigma(x) \right\}, \quad (t > 0).$$

Adopting this scheme on the two blocks  $(x, y)$  leads towards the following approximations

$$\widehat{\Psi}(x, y^k) = \langle x - x^k, \nabla_x H(x^k, y^k) \rangle + \frac{c^k}{2} \|x - x^k\|^2 + f(x), \quad (c_k > 0),$$

and

$$\widetilde{\Psi}(x^{k+1}, y) = \langle y - y^k, \nabla_y H(x^{k+1}, y^k) \rangle + \frac{d^k}{2} \|y - y^k\|^2 + g(y), \quad (d_k > 0).$$

Thus, PALM is alternating between the following two subproblems

$$x^{k+1} \in \operatorname{argmin} \left\{ \widehat{\Psi}(x, y^k) : x \in \mathbb{R}^n \right\} \quad \text{and} \quad y^{k+1} \in \operatorname{argmin} \left\{ \widetilde{\Psi}(x^{k+1}, y) : y \in \mathbb{R}^m \right\}.$$

Assuming  $\Psi$  is KL function and the generated sequence by PALM,  $\{z^k := (x^k, y^k)\}_{k \in \mathbb{N}}$ , is bounded, Bolte et al. [5] proved that the sequence obeys properties (C1) and (C2), and it converges to a critical point of  $\Psi$ .

### 3 Clustering: The Squared Euclidean Norm Case

#### 3.1 Clustering with PALM

In this section we tackle the clustering problem, given in (2.3), for which the proximity function  $d(\cdot, \cdot)$  is taken to be the classical distance function defined by  $d(u, v) = \|u - v\|^2$ . We devise a PALM-like algorithm, based on the discussion in the previous subsection. Since the clustering problem has a specific structure, we are ought to exploit it in the following manner.

- (1) The function  $w \mapsto H(w, x)$ , for fixed  $x$ , is linear and therefore there is no need to linearize it as suggested in the framework which was discussed in Section 2.2.
- (2) The function  $x \mapsto H(w, x)$ , for fixed  $w$ , is quadratic and convex. Hence, there is no need to add a proximal term as suggested in the framework which was discussed in Section 2.2.



As in the PALM algorithm, our algorithm is based on the old approach of alternating minimization, with the following adaptations which are motivated by the observations mentioned above. More precisely, with respect to  $w$  we suggest to regularize the first subproblem with proximal term as follows

$$w^i(t+1) = \arg \min_{w^i \in \Delta} \left\{ \langle w^i, d^i(x(t)) \rangle + \frac{\alpha_i(t)}{2} \|w^i - w^i(t)\|^2 \right\}, \quad i = 1, 2, \dots, m, \quad (3.1)$$

where  $\alpha_i(t) > 0$  for all  $i = 1, 2, \dots, m$ . On the other hand, with respect to  $x$  we perform exact minimization, that is,

$$x(t+1) = \operatorname{argmin} \{ H(w(t+1), x) \mid x \in \mathbb{R}^{nk} \}. \quad (3.2)$$

It is easy to check that all subproblems, with respect to  $w^i$ ,  $i = 1, 2, \dots, m$ , and  $x$ , can be rewritten explicitly (where we use  $P_\Delta$  for the orthogonal projection onto the set  $\Delta$ ). Thus, we can present now the KPALM algorithm.

### KPALM

(1) Initialization:  $(w(0), x(0)) \in \Delta^m \times \mathbb{R}^{nk}$ .

(2) General step ( $t = 0, 1, \dots$ ):

(2.1) Cluster assignment: choose certain  $\alpha_i(t) > 0$ ,  $i = 1, 2, \dots, m$ , and compute

$$w^i(t+1) = P_\Delta \left( w^i(t) - \frac{d^i(x(t))}{\alpha_i(t)} \right). \quad (3.3)$$

(2.2) Center update: for each  $l = 1, 2, \dots, k$  compute

$$x^l(t+1) = \frac{\sum_{i=1}^m w_l^i(t+1) a^i}{\sum_{i=1}^m w_l^i(t+1)}. \quad (3.4)$$

We begin our analysis of the KPALM algorithm with the following boundedness property of the generated sequence. For simplicity, from now on, we denote  $z(t) := (w(t), x(t))$ ,  $t \in \mathbb{N}$ .

**Proposition 3.1** (Boundedness of KPALM sequence). *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by KPALM. Then, the following statements hold true.*

(i) *For all  $l = 1, 2, \dots, k$ , the sequence  $\{x^l(t)\}_{t \in \mathbb{N}}$  is contained in  $\operatorname{Conv}(\mathcal{A})$ , the convex hull of  $\mathcal{A}$ , and therefore bounded by  $M = \max_{1 \leq i \leq m} \|a^i\|$ .*

(ii) *The sequence  $\{z(t)\}_{t \in \mathbb{N}}$  is bounded in  $\mathbb{R}^{km} \times \mathbb{R}^{nk}$ .*

*Proof.* (i) Let  $1 \leq l \leq k$ . We set  $\lambda_i = w_l^i(t) / \sum_{j=1}^m w_l^j(t)$ ,  $i = 1, 2, \dots, m$ , then  $\lambda_i \geq 0$  and

$\sum_{i=1}^m \lambda_i = 1$ . From (3.4) we have

$$x^l(t) = \frac{\sum_{i=1}^m w_l^i(t) a^i}{\sum_{i=1}^m w_l^i(t)} = \sum_{i=1}^m \left( \frac{w_l^i(t)}{\sum_{j=1}^m w_l^j(t)} \right) a^i = \sum_{i=1}^m \lambda_i a^i \in \text{Conv}(\mathcal{A}), \quad (3.5)$$

which proves that  $x^l(t)$  is in the convex hull of  $\mathcal{A}$ , for all  $l = 1, 2, \dots, k$  and  $t \in \mathbb{N}$ . Taking the norm of  $x^l(t)$  and using (3.5) yields that

$$\|x^l(t)\| = \left\| \sum_{i=1}^m \lambda_i a^i \right\| \leq \sum_{i=1}^m \lambda_i \|a^i\| \leq \sum_{i=1}^m \lambda_i \max_{1 \leq i \leq m} \|a^i\| = M.$$

(ii) The sequence  $\{w(t)\}_{t \in \mathbb{N}}$  is bounded, since  $w^i(t) \in \Delta$  for all  $i = 1, 2, \dots, m$  and  $t \in \mathbb{N}$ . Combined with the previous item, the result follows.  $\square$

The following assumption will be crucial for the coming analysis.

**Assumption 1.** (i) The chosen sequences of parameters  $\{\alpha_i(t)\}_{t \in \mathbb{N}}$ ,  $i = 1, 2, \dots, m$ , are bounded, that is, there exist  $\underline{\alpha}_i > 0$  and  $\overline{\alpha}_i < \infty$  for all  $i = 1, 2, \dots, m$ , such that

$$\underline{\alpha}_i \leq \alpha_i(t) \leq \overline{\alpha}_i, \quad \forall t \in \mathbb{N}. \quad (3.6)$$

(ii) For all  $t \in \mathbb{N}$  there exists  $\underline{\beta} > 0$  such that

$$2 \min_{1 \leq l \leq k} \sum_{i=1}^m w_l^i(t) := \beta(w(t)) \geq \underline{\beta}. \quad (3.7)$$

It should be noted that Assumption 1(i) is very mild since the parameters  $\alpha_i(t)$ ,  $1 \leq i \leq m$  and  $t \in \mathbb{N}$ , can be chosen arbitrarily by the user and therefore it can be controlled such that the boundedness property holds true. Assumption 1(ii) is essential since if it is not true then  $w_l^i(t) = 0$  for all  $1 \leq i \leq m$ , which means that the center  $x^l$  does not play any role in the solution process which is, of course, meaningless situation.

**Lemma 3.1.1** (Strong convexity of  $H(w, x)$  in  $x$ ). *The function  $x \mapsto H(w, x)$  is strongly convex with parameter  $\beta(w)$  which defined in (3.7), whenever  $\beta(w) > 0$ .*

*Proof.* Since the function  $x \mapsto H(w, x) = \sum_{l=1}^k \sum_{i=1}^m w_l^i \|x^l - a^i\|^2$  is  $C^2$ , it is strongly convex if and only if the smallest eigenvalue of the corresponding Hessian matrix is positive. Indeed, the Hessian is given by

$$\nabla_{x^j} \nabla_{x^l} H(w, x) = \begin{cases} 0 & \text{if } j \neq l, \quad 1 \leq j, l \leq k, \\ 2 \sum_{i=1}^m w_l^i & \text{if } j = l, \quad 1 \leq j, l \leq k. \end{cases}$$

Since the Hessian is a diagonal matrix, the smallest eigenvalue is  $\beta(w) = 2 \min_{1 \leq l \leq k} \sum_{i=1}^m w_l^i$ , and the result follows.  $\square$

Now we are ready to prove global convergence of the sequence  $\{z(t)\}_{t \in \mathbb{N}}$  generated by KPALM to a critical point of  $\Psi$  given in (2.3). We will follow here the general procedure which was discussed in Section 2.2. Therefore we need to prove that  $\{z(t)\}_{t \in \mathbb{N}}$  satisfies the conditions (C1) and (C2). We begin by proving condition (C1).

**Proposition 3.2** (Sufficient decrease property). *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by KPALM. Suppose that Assumption 1 holds true, then there exists  $\rho_1 > 0$  such that*

$$\rho_1 \|z(t+1) - z(t)\|^2 \leq \Psi(z(t)) - \Psi(z(t+1)), \quad \forall t \in \mathbb{N}.$$

*Proof.* From step (3.3), see also (3.1), we derive, for each  $i = 1, 2, \dots, m$ , the following inequality

$$\begin{aligned} H^i(w(t+1), x(t)) + \frac{\alpha_i(t)}{2} \|w^i(t+1) - w^i(t)\|^2 &= \langle w^i(t+1), d^i(x(t)) \rangle + \frac{\alpha_i(t)}{2} \|w^i(t+1) - w^i(t)\|^2 \\ &\leq \langle w^i(t), d^i(x(t)) \rangle + \frac{\alpha_i(t)}{2} \|w^i(t) - w^i(t)\|^2 \\ &= \langle w^i(t), d^i(x(t)) \rangle \\ &= H^i(w(t), x(t)). \end{aligned}$$

Hence, we obtain for all  $t \in \mathbb{N}$ , that

$$\frac{\alpha_i(t)}{2} \|w^i(t+1) - w^i(t)\|^2 \leq H^i(w(t), x(t)) - H^i(w(t+1), x(t)). \quad (3.8)$$

Denote  $\underline{\alpha} = \min_{1 \leq i \leq m} \underline{\alpha}_i$ . Summing inequality (3.8) over  $i = 1, 2, \dots, m$  yields

$$\begin{aligned} \frac{\underline{\alpha}}{2} \|w(t+1) - w(t)\|^2 &= \frac{\underline{\alpha}}{2} \sum_{i=1}^m \|w^i(t+1) - w^i(t)\|^2 \\ &\leq \sum_{i=1}^m \frac{\alpha_i(t)}{2} \|w^i(t+1) - w^i(t)\|^2 \\ &\leq \sum_{i=1}^m [H^i(w(t), x(t)) - H^i(w(t+1), x(t))] \\ &= H(w(t), x(t)) - H(w(t+1), x(t)), \end{aligned} \quad (3.9)$$

where the first inequality follows from Assumption 1(i) and the definition of  $\underline{\alpha}$ .

From Assumption 1(ii) we have that  $\beta(w(t)) \geq \underline{\beta}$ , for all  $t \in \mathbb{N}$ , and from Lemma 3.1.1 it follows that the function  $x \mapsto H(w(t), x)$  is strongly convex with parameter  $\beta(w(t))$ . Using the strong convexity yields that

$$\begin{aligned}
H(w(t+1), x(t)) - H(w(t+1), x(t+1)) &\geq \\
&\geq \langle \nabla_x H(w(t+1), x(t+1)), x(t) - x(t+1) \rangle + \frac{\beta(w(t))}{2} \|x(t) - x(t+1)\|^2 \\
&= \frac{\beta(w(t))}{2} \|x(t+1) - x(t)\|^2 \\
&\geq \frac{\underline{\beta}}{2} \|x(t+1) - x(t)\|^2,
\end{aligned} \tag{3.10}$$

where the equality follows from (3.2), since  $\nabla_x H(w(t+1), x(t+1)) = 0$ . Set  $\rho_1 = \frac{1}{2} \min \{\underline{\alpha}, \underline{\beta}\}$ , by combining (3.9) and (3.10), we get

$$\begin{aligned}
\rho_1 \|z(t+1) - z(t)\|^2 &= \rho_1 (\|w(t+1) - w(t)\|^2 + \|x(t+1) - x(t)\|^2) \leq \\
&\leq [H(w(t), x(t)) - H(w(t+1), x(t))] + [H(w(t+1), x(t)) - H(w(t+1), x(t+1))] \\
&= H(z(t)) - H(z(t+1)) \\
&= \Psi(z(t)) - \Psi(z(t+1)),
\end{aligned}$$

where the last equality follows from the fact that  $G(w(t)) = 0$ , since  $w(t) \in \Delta^m$  for all  $t \in \mathbb{N}$ , and therefore  $H(z(t)) = \Psi(z(t))$ ,  $t \in \mathbb{N}$ .  $\square$

Now, we are focusing on proving condition (C2) and we will need the following technical result.

**Lemma 3.2.1.** *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by KPALM. Then we have a Lipschitz property of  $d^i(\cdot)$  for all  $t \in \mathbb{N}$ , given by*

$$\|d^i(x(t+1)) - d^i(x(t))\| \leq 4M \|x(t+1) - x(t)\|, \quad \forall i = 1, 2, \dots, m, \quad t \in \mathbb{N},$$

where  $M = \max_{1 \leq i \leq m} \|a^i\|$ .

*Proof.* Since  $d(u, v) = \|u - v\|^2$ , we get that

$$\begin{aligned}
\|d^i(x(t+1)) - d^i(x(t))\| &= \left[ \sum_{l=1}^k \left| \|x^l(t+1) - a^i\|^2 - \|x^l(t) - a^i\|^2 \right|^2 \right]^{\frac{1}{2}} \\
&= \left[ \sum_{l=1}^k \left| \|x^l(t+1)\|^2 - 2\langle x^l(t+1), a^i \rangle + \|a^i\|^2 - \|x^l(t)\|^2 + 2\langle x^l(t), a^i \rangle - \|a^i\|^2 \right|^2 \right]^{\frac{1}{2}} \\
&\leq \left[ \sum_{l=1}^k \left( \left| \|x^l(t+1)\|^2 - \|x^l(t)\|^2 \right| + |2\langle x^l(t) - x^l(t+1), a^i \rangle| \right)^2 \right]^{\frac{1}{2}} \\
&\leq \left[ \sum_{l=1}^k \left( \left| \|x^l(t+1)\| - \|x^l(t)\| \right| \cdot \left| \|x^l(t+1)\| + \|x^l(t)\| \right| + 2\|x^l(t) - x^l(t+1)\| \cdot \|a^i\| \right)^2 \right]^{\frac{1}{2}} \\
&\leq \left[ \sum_{l=1}^k \left( \|x^l(t+1) - x^l(t)\| \cdot 2M + 2\|x^l(t+1) - x^l(t)\|M \right)^2 \right]^{\frac{1}{2}} \\
&= \left[ \sum_{l=1}^k (4M)^2 \|x^l(t+1) - x^l(t)\|^2 \right]^{\frac{1}{2}} = 4M \|x(t+1) - x(t)\|,
\end{aligned}$$

where the last inequality follows from the fact that  $x^l(t) \in \text{Conv}(\mathcal{A})$  and hence  $\|x^l(t)\| \leq M$  for all  $t \in \mathbb{N}$  and  $l = 1, 2, \dots, k$  (see Proposition 3.1(i)). This proves the desired result.  $\square$

Now, using this result we can show that  $\{z(t)\}_{t \in \mathbb{N}}$  satisfies condition (C2).

**Proposition 3.3** (Subgradient lower bound for the iterates gap). *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by KPALM. Then, there exists  $\rho_2 > 0$  and  $\gamma(t+1) \in \partial\Psi(z(t+1))$  such that*

$$\|\gamma(t+1)\| \leq \rho_2 \|z(t+1) - z(t)\|, \quad \forall t \in \mathbb{N}.$$

*Proof.* By the definition of  $\Psi$  (see (2.3)) we get

$$\partial\Psi = \nabla H + \partial G = \left( (\nabla_{w^i} H^i + \partial_{w^i} \delta_\Delta)_{i=1,2,\dots,m}, \nabla_x H \right).$$

Evaluating the last relation at  $z(t+1)$  yields

$$\begin{aligned}
\partial\Psi(z(t+1)) &= \\
&= \left( (\nabla_{w^i} H^i(w(t+1), x(t+1)) + \partial_{w^i} \delta_\Delta(w^i(t+1)))_{i=1,2,\dots,m}, \nabla_x H(w(t+1), x(t+1)) \right) \\
&= \left( (d^i(x(t+1)) + \partial_{w^i} \delta_\Delta(w^i(t+1)))_{i=1,2,\dots,m}, \nabla_x H(w(t+1), x(t+1)) \right) \\
&= \left( (d^i(x(t+1)) + \partial_{w^i} \delta_\Delta(w^i(t+1)))_{i=1,2,\dots,m}, \mathbf{0} \right), \tag{3.11}
\end{aligned}$$

where the last equality follows from (3.2), that is, the optimality condition of  $x(t+1)$ .

The optimality condition of  $w^i(t+1)$  which derived from (3.1), yields that for all  $i = 1, 2, \dots, m$  there exists  $u^i(t+1) \in \partial\delta_\Delta(w^i(t+1))$  such that

$$d^i(x(t)) + \alpha_i(t) (w^i(t+1) - w^i(t)) + u^i(t+1) = \mathbf{0}. \quad (3.12)$$

Setting  $\gamma(t+1) := \left( (d^i(x(t+1)) + u^i(t+1))_{i=1,2,\dots,m}, \mathbf{0} \right)$ , it follows from (3.11) that  $\gamma(t+1) \in \partial\Psi(z(t+1))$ . Using (3.12) we obtain

$$\gamma(t+1) = \left( (d^i(x(t+1)) - d^i(x(t)) - \alpha_i(t)(w^i(t+1) - w^i(t)))_{i=1,2,\dots,m}, \mathbf{0} \right).$$

Hence, by defining  $\bar{\alpha} = \max_{1 \leq i \leq m} \bar{\alpha}_i$ , we obtain

$$\begin{aligned} \|\gamma(t+1)\| &\leq \sum_{i=1}^m \|d^i(x(t+1)) - d^i(x(t)) - \alpha_i(t)(w^i(t+1) - w^i(t))\| \\ &\leq \sum_{i=1}^m \|d^i(x(t+1)) - d^i(x(t))\| + \sum_{i=1}^m \alpha_i(t) \|w^i(t+1) - w^i(t)\| \\ &\leq \sum_{i=1}^m 4M \|x(t+1) - x(t)\| + \bar{\alpha} \sqrt{m} \|w(t+1) - w(t)\| \\ &\leq (4Mm + \bar{\alpha} \sqrt{m}) \|z(t+1) - z(t)\|, \end{aligned}$$

where the third inequality follows from Lemma 3.2.1. Define  $\rho_2 = 4Mm + \bar{\alpha} \sqrt{m}$ , and the result follows.  $\square$

## 4 Clustering: The Euclidean Norm Case

### 4.1 A Smoothed Clustering Problem

In the previous section we have formulated the clustering problem in the following equivalent form

$$\min \left\{ \Psi(z) := H(w, x) + G(w) \mid z := (w, x) \in \mathbb{R}^{km} \times \mathbb{R}^{nk} \right\},$$

where, in this setting, the involved functions are

$$H(w, x) = \sum_{i=1}^m \langle w^i, d^i(x) \rangle = \sum_{i=1}^m \sum_{l=1}^k w_l^i \|x^l - a^i\| \quad \text{and} \quad G(w) = \sum_{i=1}^m \delta_\Delta(w^i).$$

In order to be able to use the theory mentioned in Section 2.2, we have used, in Section 3, the fact that the coupled function  $H(w, x)$  is smooth, which is not the case now. Therefore, for any  $\varepsilon > 0$ , it leads us to the following smoothed form of the clustering problem

$$\min \left\{ \Psi_\varepsilon(z) := H_\varepsilon(w, x) + G(w) \mid z := (w, x) \in \mathbb{R}^{km} \times \mathbb{R}^{nk} \right\}, \quad (4.1)$$

where

$$H_\varepsilon(w, x) = \sum_{l=1}^k H_\varepsilon^l(w, x) = \sum_{l=1}^k \sum_{i=1}^m w_l^i (\|x^l - a^i\|^2 + \varepsilon^2)^{1/2}, \quad (4.2)$$

and for all  $i = 1, 2, \dots, m$ ,

$$d_\varepsilon^i(x) = \left( (\|x^1 - a^i\|^2 + \varepsilon^2)^{1/2}, (\|x^2 - a^i\|^2 + \varepsilon^2)^{1/2}, \dots, (\|x^k - a^i\|^2 + \varepsilon^2)^{1/2} \right) \in \mathbb{R}^k. \quad (4.3)$$

Note that  $\Psi_\varepsilon(z)$  is a perturbed form of  $\Psi(z)$  for a small  $\varepsilon > 0$ , and obviously  $\Psi_0(z) = \Psi(z)$ . The following lemma shows that the smoothed function  $H_\varepsilon(w, x)$  indeed approximates  $H(w, x)$ .

**Lemma 4.0.1** (Closeness of smooth). *For any  $(w, x) \in \Delta^m \times \mathbb{R}^{nk}$  and  $\varepsilon > 0$  the following relations hold true*

$$H(w, x) \leq H_\varepsilon(w, x) \leq H(w, x) + m\varepsilon.$$

*Proof.* It is clear that for all  $\lambda \geq 0$  we have

$$\forall \lambda \geq 0, \quad \lambda \leq \sqrt{\lambda^2 + \varepsilon^2} \leq \lambda + \varepsilon.$$

Applying this inequality with  $\lambda = \|x^l - a^i\|$ , yields

$$\|x^l - a^i\| \leq (\|x^l - a^i\|^2 + \varepsilon^2)^{1/2} \leq \|x^l - a^i\| + \varepsilon,$$

for all  $l = 1, 2, \dots, k$  and  $i = 1, 2, \dots, m$ . By multiplying each inequality by  $w_l^i$  and summing over  $l = 1, 2, \dots, k$  and  $i = 1, 2, \dots, m$  we obtain

$$H(w, x) \leq H_\varepsilon(w, x) \leq H(w, x) + \sum_{i=1}^m \sum_{l=1}^k w_l^i \varepsilon.$$

Since for all  $i = 1, 2, \dots, m$ ,  $w^i \in \Delta$ , the result follows.  $\square$

Now we would like to develop an algorithm which is based on the methodology of PALM (see Section 2.2) to solve Problem (4.1). It is easy to see that with respect to  $w$ , the objective function  $\Psi_\varepsilon$  keeps on the same structure as  $\Psi$  and therefore we apply the same step as in KPALM. More precisely, for all  $i = 1, 2, \dots, m$ , we have

$$\begin{aligned} w^i(t+1) &= \arg \min_{w^i \in \Delta} \left\{ \langle w^i, d_\varepsilon^i(x(t)) \rangle + \frac{\alpha_i(t)}{2} \|w^i - w^i(t)\|^2 \right\} \\ &= P_\Delta \left( w^i(t) - \frac{d_\varepsilon^i(x(t))}{\alpha_i(t)} \right), \quad \forall t \in \mathbb{N}, \end{aligned}$$

where  $\alpha_i(t)$ ,  $i = 1, 2, \dots, m$ , is arbitrarily chosen. On the other hand, with respect to  $x$  we tackle the subproblem differently than in KPALM. Here we follow exactly the idea of PALM, that is, linearizing the function  $x \rightarrow H(w, \cdot)$ , for fixed  $w$ , and adding a regularizing term

$$x^l(t+1) = \arg \min_{x^l} \left\{ \langle x^l - x^l(t), \nabla_{x^l} H_\varepsilon(w(t+1), x(t)) \rangle + \frac{L_\varepsilon^l(w(t+1), x(t))}{2} \|x^l - x^l(t)\|^2 \right\},$$

where

$$L_\varepsilon^l(w(t+1), x(t)) := \sum_{i=1}^m \frac{w_l^i(t+1)}{(\|x^l(t) - a^i\|^2 + \varepsilon^2)^{1/2}}, \quad \forall l = 1, 2, \dots, k. \quad (4.4)$$

The motivation to use this specific regularizing parameter (see (4.4)) will be discussed later.

Now we present our algorithm for solving Problem (4.1), we call it  $\varepsilon$ -KPALM. The algorithm alternates between cluster assignment step, similar to KPALM, and centers update step that is based on certain gradient step.

#### $\varepsilon$ -KPALM

(1) Initialization:  $(w(0), x(0)) \in \Delta^m \times \mathbb{R}^{nk}$ .

(2) General step ( $t = 0, 1, \dots$ ):

(2.1) Cluster assignment: choose certain  $\alpha_i(t) > 0$ ,  $i = 1, 2, \dots, m$ , and compute

$$w^i(t+1) = P_\Delta \left( w^i(t) - \frac{d_\varepsilon^i(x(t))}{\alpha_i(t)} \right). \quad (4.5)$$

(2.2) Center update: for each  $l = 1, 2, \dots, k$  compute

$$x^l(t+1) = x^l(t) - \frac{1}{L_\varepsilon^l(w(t+1), x(t))} \nabla_{x^l} H_\varepsilon(w(t+1), x(t)). \quad (4.6)$$

**Remark 1.** Similarly to the KPALM algorithm, the sequence generated by  $\varepsilon$ -KPALM is also bounded, since here we also have that

$$\begin{aligned} x^l(t+1) &= x^l(t) - \frac{1}{L_\varepsilon^l(w(t+1), x(t))} \nabla_{x^l} H(w(t+1), x(t)) \\ &= x^l(t) - \frac{1}{L_\varepsilon^l(w(t+1), x(t))} \sum_{i=1}^m w_l^i(t+1) \cdot \frac{x^l(t) - a^i}{(\|x^l(t) - a^i\|^2 + \varepsilon^2)^{1/2}} \\ &= \frac{1}{L_\varepsilon^l(w(t+1), x(t))} \sum_{i=1}^m \left( \frac{w_l^i(t+1)}{(\|x^l(t) - a^i\|^2 + \varepsilon^2)^{1/2}} \right) a^i \in \text{Conv}(\mathcal{A}). \end{aligned}$$

Before we will be able to prove the two properties (see Section 2.2) needed for global convergence of the sequence  $\{z(t)\}_{t \in \mathbb{N}}$  generated by  $\varepsilon$ -KPALM, we will need several auxiliary results. For the simplicity of the expositions we define the function  $f_\varepsilon : \mathbb{R}^n \rightarrow \mathbb{R}$  by

$$f_\varepsilon(x) = \sum_{i=1}^m v_i (\|x - a^i\|^2 + \varepsilon^2)^{1/2},$$



for fixed non-negative numbers (not all zero)  $v_1, v_2, \dots, v_m \in \mathbb{R}$  and  $a^i \in \mathbb{R}^n$ ,  $i = 1, 2, \dots, m$ . We also need the following auxiliary function  $h_\varepsilon : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$h_\varepsilon(x, y) = \sum_{i=1}^m \frac{v_i (\|x - a^i\|^2 + \varepsilon^2)}{(\|y - a^i\|^2 + \varepsilon^2)^{1/2}}.$$

Finally we introduce the following modulus,  $L_\varepsilon : \mathbb{R}^n \rightarrow \mathbb{R}$  defined by

$$L_\varepsilon(x) = \sum_{i=1}^m \frac{v_i}{(\|x - a^i\|^2 + \varepsilon^2)^{1/2}}.$$

**Lemma 4.0.2** (Properties of the auxiliary function  $h_\varepsilon$ ). *The following properties of  $h_\varepsilon$  hold.*

(i) For any  $y \in \mathbb{R}^n$ ,

$$h_\varepsilon(y, y) = f_\varepsilon(y).$$

(ii) For any  $x, y \in \mathbb{R}^n$ ,

$$h_\varepsilon(x, y) \geq 2f_\varepsilon(x) - f_\varepsilon(y).$$

(iii) For any  $x, y \in \mathbb{R}^n$ ,

$$f_\varepsilon(x) \leq f_\varepsilon(y) + \langle \nabla f_\varepsilon(y), x - y \rangle + \frac{L_\varepsilon(y)}{2} \|x - y\|^2.$$

*Proof.* (i) Follows by substituting  $x = y$  in  $h_\varepsilon(x, y)$ .

(ii) For any two numbers  $a \in \mathbb{R}$  and  $b > 0$  the inequality

$$\frac{a^2}{b} \geq 2a - b,$$

holds true. Thus, for every  $i = 1, 2, \dots, m$ , we have that

$$\frac{\|x - a^i\|^2 + \varepsilon^2}{(\|y - a^i\|^2 + \varepsilon^2)^{1/2}} \geq 2(\|x - a^i\|^2 + \varepsilon^2)^{1/2} - (\|y - a^i\|^2 + \varepsilon^2)^{1/2}.$$

Multiplying the last inequality by  $v_i$  and summing over  $i = 1, 2, \dots, m$ , the results follows.

(iii) The function  $x \mapsto h_\varepsilon(x, y)$  is quadratic with associated matrix  $L_\varepsilon(y)\mathbf{I}$ . Therefore, its second-order taylor expansion around  $y$  leads to the following identity

$$h_\varepsilon(x, y) = h_\varepsilon(y, y) + \langle \nabla_x h_\varepsilon(y, y), x - y \rangle + L_\varepsilon(y) \|x - y\|^2.$$

Using the first two items and the fact that  $\nabla_x h_\varepsilon(y, y) = 2\nabla f_\varepsilon(y)$  yields the desired result. □

Now we get back to the  $\varepsilon$ -KPALM algorithm and prove few technical results about the involved functions which are based on the properties obtained above.

**Proposition 4.1** (Bounds for  $L_\varepsilon^l$ ). *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by  $\varepsilon$ -KPALM. Then, the following two statements hold true.*

(i) *For all  $t \in \mathbb{N}$  and  $l = 1, 2, \dots, k$  we have*

$$L_\varepsilon^l(w(t+1), x(t)) \geq \frac{\underline{\beta}}{(d_{\mathcal{A}}^2 + \varepsilon^2)^{1/2}},$$

where  $d_{\mathcal{A}} = \text{diam}(\text{Conv}(\mathcal{A}))$  is the diameter of  $\text{Conv}(\mathcal{A})$  and  $\underline{\beta}$  is given in (3.7).

(ii) *For all  $t \in \mathbb{N}$  and  $l = 1, 2, \dots, k$  we have*

$$L_\varepsilon^l(w(t+1), x(t)) \leq \frac{m}{\varepsilon}.$$

*Proof.* (i) From Assumption 1(ii) and the fact that  $x^l(t) \in \text{Conv}(\mathcal{A})$  for all  $1 \leq l \leq k$ , it follows that

$$L_\varepsilon^l(w(t+1), x(t)) = \sum_{i=1}^m \frac{w_l^i(t+1)}{(\|x^l(t) - a^i\|^2 + \varepsilon^2)^{1/2}} \geq \frac{\sum_{i=1}^m w_l^i(t+1)}{(d_{\mathcal{A}}^2 + \varepsilon^2)^{1/2}} \geq \frac{\underline{\beta}}{(d_{\mathcal{A}}^2 + \varepsilon^2)^{1/2}},$$

where the first inequality follows from the fact that  $\|x^l(t) - a^i\| \leq d_{\mathcal{A}}$ , for all  $1 \leq l \leq k$ .

(ii) Since  $w(t+1) \in \Delta^m$  we have

$$L_\varepsilon^l(w(t+1), x(t)) = \sum_{i=1}^m \frac{w_l^i(t+1)}{(\|x^l(t) - a^i\|^2 + \varepsilon^2)^{1/2}} \leq \sum_{i=1}^m \frac{1}{\varepsilon} = \frac{m}{\varepsilon},$$

as asserted. □

Now we prove the following result.

**Proposition 4.2.** *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by  $\varepsilon$ -KPALM. Then, for all  $t \in \mathbb{N}$ , we have*

$$\begin{aligned} H_\varepsilon(w(t+1), x(t+1)) &\leq H_\varepsilon(w(t+1), x(t)) + \langle \nabla_x H_\varepsilon(w(t+1), x(t)), x(t+1) - x(t) \rangle \\ &\quad + \sum_{l=1}^k \frac{L_\varepsilon^l(w(t+1), x(t))}{2} \|x^l(t+1) - x^l(t)\|^2. \end{aligned}$$

*Proof.* By definition (see (4.2)) we have, for  $i = 1, 2, \dots, m$ , that

$$H_\varepsilon^l(w(t+1), x(t)) = f_\varepsilon(x^l(t)),$$

where  $v_i = w_i^l(t+1)$ ,  $i = 1, 2, \dots, m$ . Therefore, by applying Lemma 4.0.2(iii) with  $x = x^l(t+1)$  and  $y = x^l(t)$ , we get

$$\begin{aligned} H_\varepsilon^l(w(t+1), x(t+1)) &\leq H_\varepsilon^l(w(t+1), x(t)) + \langle \nabla_{x^l} H_\varepsilon^l(w(t+1), x(t)), x(t+1) - x(t) \rangle \\ &\quad + \frac{L_\varepsilon^l(w(t+1), x(t))}{2} \|x^l(t+1) - x^l(t)\|^2. \end{aligned}$$

Summing the last inequality over  $l = 1, 2, \dots, k$ , yields

$$\begin{aligned} H_\varepsilon(w(t+1), x(t+1)) &\leq H_\varepsilon(w(t+1), x(t)) + \sum_{l=1}^k \frac{L_\varepsilon^l(w(t+1), x(t))}{2} \|x^l(t+1) - x^l(t)\|^2 \\ &\quad + \sum_{l=1}^k \langle \nabla_{x^l} H_\varepsilon(w(t+1), x(t)), x^l(t+1) - x^l(t) \rangle. \end{aligned}$$

Replacing the last term with the following compact form

$$\sum_{l=1}^k \langle \nabla_{x^l} H_\varepsilon(w(t+1), x(t)), x^l(t+1) - x^l(t) \rangle = \langle \nabla_x H_\varepsilon(w(t+1), x(t)), x(t+1) - x(t) \rangle,$$

and the result follows.  $\square$

Now we are finally ready to prove the two properties needed for guaranteeing that the sequence  $\{z(t)\}_{t \in \mathbb{N}}$  which is generated by  $\varepsilon$ -KPALM converges to a critical point of  $\Psi_\varepsilon$ .

**Proposition 4.3** (Sufficient decrease property). *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by  $\varepsilon$ -KPALM. Then, there exists  $\rho_1 > 0$  such that*

$$\rho_1 \|z(t+1) - z(t)\|^2 \leq \Psi_\varepsilon(z(t)) - \Psi_\varepsilon(z(t+1)), \quad \forall t \in \mathbb{N}.$$

*Proof.* As we already mentioned, the step with respect to  $w$  of KPALM and  $\varepsilon$ -KPALM are similar in nature and therefore following the same arguments given at the beginning of the proof of Proposition 3.2 we have that

$$\frac{\underline{\alpha}}{2} \|w(t+1) - w(t)\|^2 \leq H_\varepsilon(w(t), x(t)) - H_\varepsilon(w(t+1), x(t)), \quad (4.7)$$

where  $\underline{\alpha} = \min_{1 \leq i \leq m} \alpha_i$ . Applying Proposition 4.2 and using (4.6) we get for all  $t \in \mathbb{N}$  that

$$\begin{aligned} H_\varepsilon(w(t+1), x(t)) - H_\varepsilon(w(t+1), x(t+1)) &\geq \sum_{l=1}^k \frac{L_\varepsilon^l(w(t+1), x(t))}{2} \|x^l(t+1) - x^l(t)\|^2 \\ &\geq \frac{\underline{\beta}}{(d_{\mathcal{A}}^2 + \varepsilon^2)^{1/2}} \sum_{l=1}^k \|x^l(t+1) - x^l(t)\|^2 \\ &\geq \frac{\underline{\beta}}{(d_{\mathcal{A}}^2 + \varepsilon^2)^{1/2}} \|x(t+1) - x(t)\|^2, \end{aligned} \quad (4.8)$$

where the second inequality follows from Proposition 4.1(i). Set  $\rho_1 = \frac{1}{2} \min \left\{ \underline{\alpha}, \underline{\beta} / (d_{\mathcal{A}}^2 + \varepsilon^2)^{1/2} \right\}$ . Summing (4.7) and (4.8) yields

$$\begin{aligned} \rho_1 \|z(t+1) - z(t)\|^2 &= \rho_1 (\|w(t+1) - w(t)\|^2 + \|x(t+1) - x(t)\|^2) \leq \\ &\leq [H_\varepsilon(w(t), x(t)) - H_\varepsilon(w(t+1), x(t))] + [H_\varepsilon(w(t+1), x(t)) - H_\varepsilon(w(t+1), x(t+1))] \\ &= H_\varepsilon(z(t)) - H_\varepsilon(z(t+1)) \\ &= \Psi_\varepsilon(z(t)) - \Psi_\varepsilon(z(t+1)), \end{aligned}$$

where the last equality follows from the fact that  $G(w(t)) = 0$ , since  $w(t) \in \Delta^m$  for all  $t \in \mathbb{N}$ . This proves the desired result.  $\square$

The next two lemmas will be useful in proving the subgradient lower bounds for the iterates gap property of the sequence generated by  $\varepsilon$ -KPALM.

**Lemma 4.3.1.** *For all  $y, z \in \mathbb{R}^n$  the following statement holds true*

$$\|\nabla f_\varepsilon(y) - \nabla f_\varepsilon(z)\| \leq \frac{2L_\varepsilon(z)L_\varepsilon(y)}{L_\varepsilon(z) + L_\varepsilon(y)} \|z - y\|.$$

*Proof.* Let  $z \in \mathbb{R}^n$  be a fixed vector. Define the following function

$$\tilde{f}_\varepsilon(y) = f_\varepsilon(y) - \langle \nabla f_\varepsilon(z), y \rangle,$$

hence,

$$f_\varepsilon(y) = \tilde{f}_\varepsilon(y) + \langle \nabla f_\varepsilon(z), y \rangle. \quad (4.9)$$

Substituting (4.9) into Lemma 4.0.2(iii) yields

$$\tilde{f}_\varepsilon(x) \leq \tilde{f}_\varepsilon(y) + \left\langle \nabla \tilde{f}_\varepsilon(y), x - y \right\rangle + \frac{L_\varepsilon(y)}{2} \|x - y\|^2. \quad (4.10)$$

It is clear that the optimal point of  $\tilde{f}_\varepsilon$  is  $z$  since  $\nabla \tilde{f}_\varepsilon(z) = 0$ , therefore using (4.10) with  $x = y - (1/L_\varepsilon(y)) \nabla \tilde{f}_\varepsilon(y)$  yields

$$\begin{aligned} \tilde{f}_\varepsilon(z) &\leq \tilde{f}_\varepsilon\left(y - \frac{1}{L_\varepsilon(y)} \nabla \tilde{f}_\varepsilon(y)\right) \leq \tilde{f}_\varepsilon(y) + \left\langle \nabla \tilde{f}_\varepsilon(y), -\frac{1}{L_\varepsilon(y)} \nabla \tilde{f}_\varepsilon(y) \right\rangle + \frac{L_\varepsilon(y)}{2} \left\| \frac{1}{L_\varepsilon(y)} \nabla \tilde{f}_\varepsilon(y) \right\|^2 \\ &= \tilde{f}_\varepsilon(y) - \frac{1}{2L_\varepsilon(y)} \left\| \nabla \tilde{f}_\varepsilon(y) \right\|^2. \end{aligned}$$

Thus, using the definition of  $\tilde{f}_\varepsilon$  and the fact that  $\nabla \tilde{f}_\varepsilon(y) = \nabla f_\varepsilon(y) - \nabla f_\varepsilon(z)$ , yields that

$$f_\varepsilon(z) \leq f_\varepsilon(y) + \langle \nabla f_\varepsilon(z), z - y \rangle - \frac{1}{2L_\varepsilon(y)} \|\nabla f_\varepsilon(y) - \nabla f_\varepsilon(z)\|^2.$$

Now, following the same arguments we can show that

$$f_\varepsilon(y) \leq f_\varepsilon(z) + \langle \nabla f_\varepsilon(y), y - z \rangle - \frac{1}{2L_\varepsilon(z)} \|\nabla f_\varepsilon(z) - \nabla f_\varepsilon(y)\|^2.$$

Combining the last two inequalities yields that

$$\left( \frac{1}{2L_\varepsilon(z)} + \frac{1}{2L_\varepsilon(y)} \right) \|\nabla f_\varepsilon(y) - \nabla f_\varepsilon(z)\|^2 \leq \langle \nabla f_\varepsilon(z) - \nabla f_\varepsilon(y), z - y \rangle,$$

that is,

$$\|\nabla f_\varepsilon(y) - \nabla f_\varepsilon(z)\| \leq \frac{2L_\varepsilon(z)L_\varepsilon(y)}{L_\varepsilon(z) + L_\varepsilon(y)} \|z - y\|,$$

for all  $z, y \in \mathbb{R}^n$ . This proves the desired result.  $\square$

**Lemma 4.3.2.** *For any  $x, y \in \mathbb{R}^{nk}$  such that  $x^l, y^l \in \text{Conv}(\mathcal{A})$  for all  $1 \leq l \leq k$  the following inequality holds*

$$\|d_\varepsilon^i(x) - d_\varepsilon^i(y)\| \leq \frac{d_{\mathcal{A}}}{\varepsilon} \|x - y\|, \quad \forall i = 1, 2, \dots, m,$$

with  $d_{\mathcal{A}} = \text{diam}(\text{Conv}(\mathcal{A}))$  and  $d_\varepsilon^i(\cdot)$  is defined in (4.3).

*Proof.* Define  $\psi(t) = \sqrt{t + \varepsilon^2}$ , for  $t \geq 0$ . Using the Lagrange mean value theorem over  $a > b \geq 0$  yields

$$\frac{\psi(a) - \psi(b)}{a - b} = \psi'(c) = \frac{1}{2\sqrt{c + \varepsilon^2}} \leq \frac{1}{2\varepsilon},$$

where  $c \in (b, a)$ . Therefore, for all  $i = 1, 2, \dots, m$  and  $l = 1, 2, \dots, k$  we have

$$\begin{aligned} \left| (\|x^l - a^i\|^2 + \varepsilon^2)^{1/2} - (\|y^l - a^i\|^2 + \varepsilon^2)^{1/2} \right| &\leq \frac{1}{2\varepsilon} \left| \|x^l - a^i\|^2 + \varepsilon^2 - (\|y^l - a^i\|^2 + \varepsilon^2) \right| \\ &= \frac{1}{2\varepsilon} \left| \|x^l - a^i\|^2 - \|y^l - a^i\|^2 \right| \\ &= \frac{1}{2\varepsilon} \left| \|x^l - a^i\| + \|y^l - a^i\| \right| \cdot \left| \|x^l - a^i\| - \|y^l - a^i\| \right| \\ &\leq \frac{1}{\varepsilon} d_{\mathcal{A}} \|x^l - y^l\|, \end{aligned}$$

where the last inequality follows from  $\|x^l - a^i\|, \|y^l - a^i\| \leq d_{\mathcal{A}}$  and the reverse triangle inequality. Therefore,

$$\begin{aligned} \|d_\varepsilon^i(x) - d_\varepsilon^i(y)\| &= \left[ \sum_{l=1}^k \left| (\|x - a^i\|^2 + \varepsilon^2)^{1/2} - (\|y - a^i\|^2 + \varepsilon^2)^{1/2} \right|^2 \right]^{\frac{1}{2}} \\ &\leq \left[ \sum_{l=1}^k \left( \frac{1}{\varepsilon} d_{\mathcal{A}} \|x^l - y^l\| \right)^2 \right]^{\frac{1}{2}} \\ &= \frac{d_{\mathcal{A}}}{\varepsilon} \|x - y\|, \end{aligned}$$

as asserted.  $\square$

**Proposition 4.4** (Subgradient lower bound for the iterates gap). *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by  $\varepsilon$ -KPALM. Then, there exists  $\rho_2 > 0$  and  $\gamma(t+1) \in \partial \Psi_\varepsilon(z(t+1))$  such that*

$$\|\gamma(t+1)\| \leq \rho_2 \|z(t+1) - z(t)\|, \quad \forall t \in \mathbb{N}.$$

*Proof.* Repeating the steps of the proof in the case of KPALM (see Proposition 3.3) yields that

$$\gamma(t+1) := \left( (d_\varepsilon^i(x(t+1)) + u^i(t+1))_{i=1,\dots,m}, \nabla_x H_\varepsilon(w(t+1), x(t+1)) \right) \in \partial \Psi_\varepsilon(z(t+1)), \quad (4.11)$$

where  $u^i(t+1) \in \partial \delta_\Delta(w^i(t+1))$ ,  $i = 1, 2, \dots, m$ . Now, writing the optimality condition of step (4.5), yields that

$$d_\varepsilon^i(x(t)) + \alpha_i(t) (w^i(t+1) - w^i(t)) + u^i(t+1) = \mathbf{0}. \quad (4.12)$$

Plugging (4.12) into (4.11), and taking the norm yields

$$\begin{aligned} \|\gamma(t+1)\| &\leq \sum_{i=1}^m \|d_\varepsilon^i(x(t+1)) - d_\varepsilon^i(x(t)) - \alpha_i(t) (w^i(t+1) - w^i(t))\| \\ &\quad + \|\nabla_x H_\varepsilon(w(t+1), x(t+1))\| \\ &\leq \sum_{i=1}^m \|d_\varepsilon^i(x(t+1)) - d_\varepsilon^i(x(t))\| + \sum_{i=1}^m \alpha_i(t) \|w^i(t+1) - w^i(t)\| \\ &\quad + \|\nabla_x H_\varepsilon(w(t+1), x(t+1))\| \\ &\leq \frac{md_A}{\varepsilon} \|x(t+1) - x(t)\| + \bar{\alpha} \sqrt{m} \|w(t+1) - w(t)\| + \|\nabla_x H_\varepsilon(w(t+1), x(t+1))\|, \end{aligned}$$

where the last inequality follows from Lemma 4.3.2 and the fact that  $\bar{\alpha} = \max_{1 \leq i \leq m} \bar{\alpha}_i$ .

Next we will show that  $\|\nabla_x H_\varepsilon(w(t+1), x(t+1))\| \leq c \|x(t+1) - x(t)\|$ , for some constant  $c > 0$ . Indeed, for all  $l = 1, 2, \dots, k$ , we have

$$\begin{aligned} \nabla_{x^l} H_\varepsilon(w(t+1), x(t+1)) &= \nabla_{x^l} H_\varepsilon(w(t+1), x(t+1)) - \nabla_{x^l} H_\varepsilon(w(t+1), x(t)) \\ &\quad + \nabla_{x^l} H_\varepsilon(w(t+1), x(t)) \\ &= \nabla_{x^l} H_\varepsilon(w(t+1), x(t+1)) - \nabla_{x^l} H_\varepsilon(w(t+1), x(t)) \\ &\quad + L_\varepsilon^l(w(t+1), x(t)) (x^l(t) - x^l(t+1)), \end{aligned} \quad (4.13)$$

where the last equality follows from (4.6). Therefore,

$$\begin{aligned} \|\nabla_x H_\varepsilon(w(t+1), x(t+1))\| &\leq \sum_{l=1}^k \|\nabla_{x^l} H_\varepsilon(w(t+1), x(t+1))\| \\ &\leq \sum_{l=1}^k L_\varepsilon^l(w(t+1), x(t)) \|x^l(t+1) - x^l(t)\| \\ &\quad + \sum_{l=1}^k \|\nabla_{x^l} H_\varepsilon(w(t+1), x(t+1)) - \nabla_{x^l} H_\varepsilon(w(t+1), x(t))\| \\ &\leq \frac{m}{\varepsilon} \sum_{l=1}^k \|x^l(t+1) - x^l(t)\| + \sum_{l=1}^k \gamma^l(t) \|x^l(t+1) - x^l(t)\|, \end{aligned} \quad (4.14)$$

where the last inequality follows from Proposition 4.1(ii) and Lemma 4.3.1 using

$$\gamma^l(t) = \frac{2L_\varepsilon^l(w(t+1), x(t))L_\varepsilon^l(w(t+1), x(t+1))}{L_\varepsilon^l(w(t+1), x(t)) + L_\varepsilon^l(w(t+1), x(t+1))}, \quad l = 1, 2, \dots, k.$$

From Proposition 4.1(ii) we obtain that

$$\gamma^l(t) = \frac{2}{\frac{1}{L_\varepsilon^l(w(t+1), x(t))} + \frac{1}{L_\varepsilon^l(w(t+1), x(t+1))}} \leq \frac{2}{\frac{\varepsilon}{m} + \frac{\varepsilon}{m}} = \frac{m}{\varepsilon}.$$

Hence, from (4.14), we have

$$\|\nabla_x H_\varepsilon(w(t+1), x(t+1))\| \leq \frac{2m}{\varepsilon} \sum_{l=1}^k \|x^l(t+1) - x^l(t)\| \leq \frac{2m\sqrt{k}}{\varepsilon} \|x(t+1) - x(t)\|. \quad (4.15)$$

Therefore, setting  $\rho_2 = \frac{md_A}{\varepsilon} + \bar{\alpha}\sqrt{m} + \frac{2m\sqrt{k}}{\varepsilon}$ , yields the result.  $\square$

## 4.2 Different Approach Towards Solving the Smoothed $H_\varepsilon$

In this section we describe a different approach towards solving the smoothed clustering problem described in (4.1). Using the Arithmetic-Geometric inequality we derive the following simple observation

$$\frac{1}{2} \min_{s \geq 0} \left\{ s\lambda + \frac{1}{\lambda} \right\} \geq \min_{s \geq 0} \left\{ \sqrt{s\lambda \cdot \frac{1}{s}} \right\} = \sqrt{\lambda}, \quad \forall \lambda \geq 0,$$

and the unique minimizer is given by  $s^* = 1/\sqrt{\lambda}$ . Using this fact we can write

$$\sqrt{\|u\|^2 + \varepsilon^2} = \frac{1}{2} \min_{v \geq 0} \left\{ v (\|u\|^2 + \varepsilon^2) + \frac{1}{v} \right\}, \quad (4.16)$$

with  $v^* = 1/\sqrt{\|u\|^2 + \varepsilon^2}$ . Thus, instead of solving problem (4.1) with  $H_\varepsilon(\cdot, \cdot)$ , defined in (4.2), we replace it with the following function

$$B_\varepsilon(v, w, x) = \frac{1}{2} \sum_{i=1}^m \sum_{l=1}^k \left\{ v_l^i w_l^i (\|x^l - a^i\|^2 + \varepsilon^2) + \frac{w_l^i}{v_l^i} \right\}, \quad (4.17)$$

where  $v = (v^1, v^2, \dots, v^m)$ , and then problem (4.1) can be written equivalently as

$$\min_{x, v, w} \{ B_\varepsilon(v, w, x) + G(w) : v \geq 0 \}.$$

For all  $i = 1, 2, \dots, m$  we define  $b_\varepsilon^i : \mathbb{R}^{mk} \times \mathbb{R}^{nk} \rightarrow \mathbb{R}^k$  by

$$b_\varepsilon^i(v, x) = \left( \frac{1}{2} v_l^i (\|x^l - a^i\|^2 + \varepsilon^2) + \frac{1}{2v_l^i} \right)_{l=1,2,\dots,k} \in \mathbb{R}^k,$$

and we have that

$$B_\varepsilon(v, w, x) = \sum_{i=1}^m \langle w^i, b_\varepsilon^i(v, x) \rangle. \quad (4.18)$$

Now the situation is similar to that of Section 3, namely

- (1) The function  $w \mapsto B_\varepsilon(v, w, x)$ , for fixed  $v$  and  $x$ , is linear;
- (2) The function  $x \mapsto B_\varepsilon(v, w, x)$ , for fixed  $v$  and  $w$ , is quadratic and convex.

Hence we can tackle these two steps as in KPALM.

Equipped with these observation we proceed to a PALM-like algorithm, which is based on three steps alternating minimization. More precisely, with respect to  $v$  we perform exact minimization

$$v(t+1) = \operatorname{argmin} \{B_\varepsilon(v, w(t), x(t)) : v \geq 0\}$$

It should be noted that this problem can be written equivalently by

$$v(t+1) = \operatorname{argmin} \{B_\varepsilon(v, w(t), x(t)) : v \in I^{mk}\},$$

where  $I := [1/\kappa, 1/\varepsilon]$  and  $\kappa = \sqrt{d_{\mathcal{A}}^2 + \varepsilon^2}$ . With respect to  $w$ , as in KPALM case, we need to solve the subproblem given by

$$w^i(t+1) = \operatorname{argmin}_{w^i \in \Delta} \left\{ \langle w^i, b_\varepsilon^i(v(t+1), x(t)) \rangle + \frac{\alpha_i(t)}{2} \|w^i - w^i(t)\|^2 \right\}, \quad i = 1, 2, \dots, m, \quad (4.19)$$

where  $\alpha_i(t) > 0$ ,  $i = 1, 2, \dots, m$ . With respect to  $x$ , again as in KPALM case, we perform exact minimization

$$x(t+1) = \operatorname{argmin} \{B_\varepsilon(v(t+1), w(t+1), x) \mid x \in \mathbb{R}^{nk}\}. \quad (4.20)$$

It is easy to check that explicit solutions to all three subproblems are given by

$$v_l^i(t+1) = \frac{1}{\left(\|x^l(t) - a^i\|^2 + \varepsilon^2\right)^{1/2}}, \quad i = 1, 2, \dots, m, \quad l = 1, 2, \dots, k, \quad (4.21)$$

$$w^i(t+1) = P_\Delta \left( w^i(t) - \frac{b_\varepsilon^i(v(t+1), x(t))}{\alpha_i(t)} \right), \quad i = 1, 2, \dots, m, \quad (4.22)$$

and

$$x^l(t+1) = \frac{\sum_{i=1}^m w_l^i(t+1) v_l^i(t+1) a^i}{\sum_{i=1}^m w_l^i(t+1) v_l^i(t+1)}, \quad l = 1, 2, \dots, k. \quad (4.23)$$

From the subproblem for  $v$  and the observation given in (4.16), we derive the following three relations

$$B_\varepsilon(v(t+1), w, x(t)) = H_\varepsilon(w, x(t)), \quad \forall t \in \mathbb{N}, \quad \forall w \in \Delta^m, \quad (4.24)$$



$$b_\varepsilon^i(v(t+1), x(t)) = d_\varepsilon^i(x(t)), \quad \forall t \in \mathbb{N}, i = 1, 2, \dots, m, \quad (4.25)$$

where  $d_\varepsilon^i$  is defined in (4.3), and

$$B_\varepsilon(v, w, x) \geq H_\varepsilon(w, x), \quad \forall (v, w, x) \in I^{mk} \times \Delta^m \times \mathbb{R}^{nk}. \quad (4.26)$$

Substituting (4.25) into (4.22) yields

$$w^i(t+1) = P_\Delta \left( w^i(t) - \frac{d_\varepsilon^i(x(t))}{\alpha_i(t)} \right), \quad i = 1, 2, \dots, m.$$

Moreover, substituting (4.21) into (4.23) yields

$$x^l(t+1) = \frac{1}{L_\varepsilon^l(w(t+1), x(t))} \sum_{i=1}^m \left( \frac{w_l^i(t+1)}{(\|x^l(t) - a^i\|^2 + \varepsilon^2)^{1/2}} \right) a^i, \quad l = 1, 2, \dots, k,$$

where  $L_\varepsilon^l$  is defined in (4.4). Thus, we recover the  $\varepsilon$ -KPALM algorithm, which means these two different approaches lead to the same iterative algorithm. However, with the current approach we can swiftly prove the sufficient decrease and the subgradient lower bound for the iterates gap properties, which are needed to obtain global convergence of  $\{z(t)\}_{t \in \mathbb{N}}$  that generated by  $\varepsilon$ -KPALM.

**Proposition 4.5** (Sufficient decrease property). *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by  $\varepsilon$ -KPALM. Then, there exists  $\rho_1 > 0$  such that*

$$\rho_1 \|z(t+1) - z(t)\|^2 \leq \Psi_\varepsilon(z(t)) - \Psi_\varepsilon(z(t+1)), \quad \forall t \in \mathbb{N}.$$

*Proof.* From (4.19) we have

$$\langle w^i(t+1), b_\varepsilon^i(v(t+1), x(t)) \rangle + \frac{\alpha_i(t)}{2} \|w^i(t+1) - w^i(t)\|^2 \leq \langle w^i(t), b_\varepsilon^i(v(t+1), x(t)) \rangle.$$

Summing the last inequality over  $i = 1, 2, \dots, m$  and applying (4.18) yields

$$B_\varepsilon(v(t+1), w(t+1), x(t)) + \sum_{i=1}^m \frac{\alpha_i(t)}{2} \|w^i(t+1) - w^i(t)\|^2 \leq B_\varepsilon(v(t+1), w(t), x(t)).$$

Using Assumption 1(i) we derive

$$\begin{aligned} \frac{\alpha}{2} \|w(t+1) - w(t)\|^2 &\leq B_\varepsilon(v(t+1), w(t), x(t)) - B_\varepsilon(v(t+1), w(t+1), x(t)) \\ &\leq H_\varepsilon(w(t), x(t)) - H_\varepsilon(w(t+1), x(t)), \end{aligned} \quad (4.27)$$

where the last inequality follows from (4.24) and (4.26).

Since the function  $x \mapsto B_\varepsilon(v, w, x)$  is  $C^2$ , and

$$\nabla_{x^j} \nabla_{x^l} B_\varepsilon(v, w, x) = \begin{cases} 0 & \text{if } j \neq l, \quad 1 \leq j, l \leq k, \\ \sum_{i=1}^m w_l^i v_l^i & \text{if } j = l, \quad 1 \leq j, l \leq k, \end{cases}$$

it follows that, the function  $x \mapsto B_\varepsilon(v(t+1), w(t), x)$  is strongly convex with parameter  $\underline{\beta}/2\kappa$ , for all  $t \in \mathbb{N}$ . Indeed,

$$\nabla_{x^i}^2 B_\varepsilon(v(t+1), w(t), x) = \sum_{i=1}^m w_i^i(t) v_i^i(t+1) \geq \frac{1}{\kappa} \sum_{i=1}^m w_i^i(t) \geq \frac{\beta(w^i(t))}{2\kappa} > \frac{\underline{\beta}}{2\kappa} > 0,$$

where the first inequality follows from the fact that  $v_i^i(t) \in I$  for all  $t \in \mathbb{N}$ ,  $\beta(\cdot)$  is defined in (3.7), and the second inequality is due to Assumption 1(ii). Using the strong convexity property we deduce the sufficient decrease in  $x$ , as follows,

$$\begin{aligned} \frac{\underline{\beta}}{4\kappa} \|x(t+1) - x(t)\|^2 &= \langle \nabla_x B_\varepsilon(z(t+1)), x(t) - x(t+1) \rangle + \frac{\underline{\beta}}{4\kappa} \|x(t+1) - x(t)\|^2 \\ &\leq B_\varepsilon(v(t+1), w(t+1), x(t)) - B_\varepsilon(z(t+1)) \\ &\leq H_\varepsilon(w(t+1), x(t)) - H_\varepsilon(w(t+1), x(t+1)), \end{aligned} \quad (4.28)$$

where the first equality follows from (4.20), the second inequality follows from the strong convexity, and the last inequality is due to (4.24) and (4.26). Set  $\rho_1 = \min\{\underline{\alpha}/2, \underline{\beta}/4\kappa\}$ . Summing (4.27) and (4.28), we get

$$\begin{aligned} \rho_1 \|z(t+1) - z(t)\|^2 &= \rho_1 (\|w(t+1) - w(t)\|^2 + \|x(t+1) - x(t)\|^2) \leq \\ &\leq [H_\varepsilon(z(t)) - H_\varepsilon(w(t+1), x(t))] + [H_\varepsilon(w(t+1), x(t)) - H_\varepsilon(z(t+1))] \\ &= H_\varepsilon(z(t)) - H_\varepsilon(z(t+1)) \\ &= \Psi_\varepsilon(z(t)) - \Psi_\varepsilon(z(t+1)), \end{aligned}$$

where the last equality follows from the fact that  $G(w(t)) = 0$ , since  $w(t) \in \Delta^m$  for all  $t \in \mathbb{N}$ . This proves the desired result.  $\square$

**Proposition 4.6** (Subgradient lower bound for the iterates gap). *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by  $\varepsilon$ -KPALM. Then, there exists  $\rho_2 > 0$  and  $\gamma(t+1) \in \partial \Psi_\varepsilon(z(t+1))$  such that*

$$\|\gamma(t+1)\| \leq \rho_2 \|z(t+1) - z(t)\|, \quad \forall t \in \mathbb{N}.$$

*Proof.* By the definition of  $\Psi_\varepsilon$  (see (4.1)) we get

$$\Psi_\varepsilon(w, x) = H_\varepsilon(w, x) + \sum_{i=1}^m \delta_\Delta(w^i) \quad (4.29)$$

Differentiating (4.29) with respect to  $x$  and evaluating it in  $z(t+1)$  yields

$$\partial_x \Psi_\varepsilon(z(t+1)) = \nabla_x H_\varepsilon(z(t+1)). \quad (4.30)$$

Similarly, differentiating (4.29) with respect to  $w^i$  and evaluating it in  $z(t+1)$  yields

$$\partial_{w^i} \Psi_\varepsilon(z(t+1)) = d_\varepsilon^i(x(t+1)) + \partial_{w^i} \delta_\Delta(w^i(t+1)). \quad (4.31)$$

The optimality condition of  $w^i(t+1)$  which derived from (4.19), yields that for all  $i = 1, 2, \dots, m$  there exists  $u^i(t+1) \in \partial\delta_\Delta(w^i(t+1))$  such that

$$\begin{aligned} \mathbf{0} &= b_\varepsilon^i(v(t+1), x(t)) + \alpha_i(t) (w^i(t+1) - w^i(t)) + u^i(t+1) \\ &= d_\varepsilon^i(x(t)) + \alpha_i(t) (w^i(t+1) - w^i(t)) + u^i(t+1), \end{aligned} \quad (4.32)$$

where the last equality follows from (4.25). Substituting (4.32) into (4.31) and combining with (4.30) we deduce that

$$\gamma(t+1) := \left( (d^i(x(t+1)) - d^i(x(t)) - \alpha_i(t)(w^i(t+1) - w^i(t)))_{i=1,2,\dots,m}, \nabla_x H_\varepsilon(z(t+1)) \right) \in \partial\Psi_\varepsilon(z(t+1)).$$

Therefore,

$$\begin{aligned} \|\gamma(t+1)\| &\leq \sum_{i=1}^m \|d^i(x(t+1)) - d^i(x(t)) - \alpha_i(t)(w^i(t+1) - w^i(t))\| + \|\nabla_x H_\varepsilon(z(t+1))\| \\ &\leq \sum_{i=1}^m \|d^i(x(t+1)) - d^i(x(t))\| + \bar{\alpha} \sum_{i=1}^m \|w^i(t+1) - w^i(t)\| + \frac{2m\sqrt{k}}{\varepsilon} \|x(t+1) - x(t)\| \\ &\leq \sum_{i=1}^m \frac{d_A}{\varepsilon} \|x(t+1) - x(t)\| + \bar{\alpha}\sqrt{m} \|w(t+1) - w(t)\| + \frac{2m\sqrt{k}}{\varepsilon} \|x(t+1) - x(t)\| \\ &\leq \left( \frac{md_A}{\varepsilon} + \bar{\alpha}\sqrt{m} + \frac{2m\sqrt{k}}{\varepsilon} \right) \|z(t+1) - z(t)\|, \end{aligned}$$

where the second inequality was established in Proposition 4.4 (see (4.15)) and the third inequality follows from Lemma 4.3.2. Define  $\rho_2 = \frac{md_A}{\varepsilon} + \bar{\alpha}\sqrt{m} + \frac{2m\sqrt{k}}{\varepsilon}$ , and the result follows.  $\square$

## 5 Returning to KMEANS

### 5.1 Similarity to KMEANS

The famous KMEANS algorithm has close relation to KPALM algorithm. KMEANS alternates between cluster assignment and centers update steps as well. In detail, we can write its steps in the following manner

#### KMEANS

- (1) Initialization:  $x(0) \in \mathbb{R}^{nk}$ .
- (2) General step ( $t = 0, 1, \dots$ ):

(2.1) Cluster assignment: for  $i = 1, 2, \dots, m$  compute

$$w^i(t+1) = \arg \min_{w^i \in \Delta} \{ \langle w^i, d^i(x(t)) \rangle \}. \quad (5.1)$$

(2.2) Center update: for  $l = 1, 2, \dots, k$  compute

$$x^l(t+1) = \frac{\sum_{i=1}^m w_l^i(t+1) a^i}{\sum_{i=1}^m w_l^i(t+1)}. \quad (5.2)$$

It is easy to see that if we take  $\alpha_i(t) = 0$  for all  $1 \leq i \leq m$  and  $t \in \mathbb{N}$ , then KPALM becomes KMEANS. We aim to use the theory described in Section 2.2 once again and show that the sequence generated by KMEANS converges to a critical point of  $\Psi(\cdot)$ , as defined in (2.3). The sufficient decrease proof of Section 3 collapses in this case, since it is based on Assumption 1(i), that is,  $\alpha_i(t) > \underline{\alpha}_i > 0$ , for all  $t \in \mathbb{N}$  and  $i = 1, 2, \dots, m$ . However, the proof of the subgradient lower bound for the iterates gap property follows through as is. In the following discussion we present the means to treat the case that  $\alpha_i(t) = 0$ , and prove the sufficient decrease property.

**Lemma 5.0.1.** *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by KMEANS. Then, there exists  $c > 0$  such that*

$$\|w^i(t+1) - w^i(t)\| \leq c \|x(t+1) - x(t)\|, \quad \forall i = 1, 2, \dots, m, t \in \mathbb{N}.$$

*Proof.* At each iteration KMEANS partitions the set  $\mathcal{A}$  into  $k$  clusters, and the center of each cluster is its mean. Since the number of these partitions is finite, there exists a finite set  $\mathcal{C} = \{x^1, x^2, \dots, x^N\} \subset \mathbb{R}^{nk}$  such that for all  $t \in \mathbb{N}$ ,  $x(t) \in \mathcal{C}$ . We denote

$$r = \min_{1 \leq j < l \leq N} \|x^j - x^l\|,$$

and set  $c = \sqrt{2}/r$ . At each iteration, the point  $a^i$  can move from one cluster to another, hence

$$\|w^i(t+1) - w^i(t)\| \leq \sqrt{2}.$$

Therefore, combining these arguments yields

$$\frac{\|w^i(t+1) - w^i(t)\|}{\|x(t+1) - x(t)\|} \leq \frac{\sqrt{2}}{r}.$$

In case that  $x(t+1) = x(t)$ , this implies that none of the clusters has changed, hence we proved the statement in both cases.  $\square$

Equipped with the last lemma we briefly prove the sufficient decrease property of KMEANS.

**Proposition 5.1** (Sufficient decrease property for KMEANS sequence). *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by KMEANS. Then, there exists  $\rho_1 > 0$  such that*

$$\rho_1 \|z(t+1) - z(t)\|^2 \leq \Psi_\varepsilon(z(t)) - \Psi_\varepsilon(z(t+1)) \quad \forall t \in \mathbb{N}.$$

*Proof.* The function  $x \mapsto H(w(t), x)$  remains strongly convex with parameter  $\beta(w(t))$  (see (3.10)), hence we have a sufficient decrease in the  $x$  variable, namely,

$$\frac{\beta}{2} \|x(t+1) - x(t)\|^2 \leq H(w(t), x(t)) - H(w(t+1), x(t+1)). \quad (5.3)$$

Setting  $\rho_1 = \underline{\beta}/2(1 + mc^2)$ , we can write

$$\begin{aligned} \rho_1 \|z(t+1) - z(t)\|^2 &= \rho_1 \sum_{i=1}^m \|w^i(t+1) - w^i(t)\|^2 + \rho_1 \|x(t+1) - x(t)\|^2 \\ &\leq \rho_1 (1 + mc^2) \|x(t+1) - x(t)\|^2 \\ &\leq H(w(t), x(t)) - H(w(t+1), x(t+1)) \\ &= \Psi(z(t)) - \Psi(z(t+1)) \end{aligned}$$

where the first inequality follows from Lemma 5.0.1, the second follows from (5.3), and the last equality follows from the fact that  $G(w(t)) = 0$ , for all  $t \in \mathbb{N}$ .  $\square$

## 5.2 KMEANS Local Minima Convergence Proof

In this section we present a simple and direct proof that KMEANS converges to local minima. We start with rewriting the KMEANS algorithm, in its most familiar form

### KMEANS

(1) Initialization:  $x(0) \in \mathbb{R}^{nk}$ .

(2) General step ( $t = 0, 1, \dots$ ):

(2.1) Cluster assignment: for  $i = 1, 2, \dots, m$  compute

$$C^l(t+1) = \{a \in \mathcal{A} \mid \|a - x^l(t)\| \leq \|a - x^j(t)\|, \quad \forall 1 \leq l \leq k\}. \quad (5.4)$$

(2.2) Center update: for  $l = 1, 2, \dots, k$  compute

$$x^l(t+1) = \text{mean}(C^l(t+1)) := \frac{1}{|C^l(t+1)|} \sum_{a \in C^l(t+1)} a. \quad (5.5)$$

(2.3) Stopping criteria: halt if

$$\forall 1 \leq l \leq k \quad C^l(t+1) = C^l(t) \quad (5.6)$$

As in KPALM, KMEANS needs Assumption 1(ii) for step (5.5) to be well defined. In order to prove the convergence of KMEANS to local minimum, we will need to following assumption.

**Assumption 2.** *Let  $t \in \mathbb{N}$  be the final iteration of KMEANS run, then we assume that each  $a \in \mathcal{A}$  belongs exclusively to single cluster  $C^l(t)$ .*

For any  $x \in \mathbb{R}^{nk}$  we denote the super-partition of  $\mathcal{A}$  with respect to  $x$  by  $\overline{C^l}(x) = \{a \in \mathcal{A} \mid \|a - x^l\| \leq \|a - x^j\|, \forall j \neq l\}$ , for all  $1 \leq l \leq k$ , and the sub-partition of  $\mathcal{A}$  by  $\underline{C^l}(x) = \{a \in \mathcal{A} \mid \|a - x^l\| < \|a - x^j\|, \forall j \neq l\}$ . Moreover, denote  $R_{lj}(t) = \min_{a \in C^l(t)} \{\|a - x^j(t)\| - \|a - x^l(t)\|\}$  for all  $1 \leq l, j \leq k$ , and  $r(t) = \min_{l \neq j} R_{lj}$ .

Due to Assumption 2 we have that  $\overline{C^l}(x(t)) = \underline{C^l}(x(t)) = C^l(t+1)$ , for all  $1 \leq l \leq k$ ,  $t \in \mathbb{N}$ , we also have that  $r(t) > 0$  for all  $t \in \mathbb{N}$ .

**Proposition 5.2.** *Let  $(C(t), x(t))$  be the clusters and centers KMEANS returns. Denote by  $U = B\left(x^1(t), \frac{r(t)}{2}\right) \times B\left(x^2(t), \frac{r(t)}{2}\right) \times \cdots \times B\left(x^l(t), \frac{r(t)}{2}\right)$  an open neighbourhood of  $x(t)$ , then for any  $x \in U$  we have  $C^l(t) = \underline{C^l}(x)$  for all  $1 \leq l \leq k$ .*

*Proof.* Pick some  $a \in C^l(t)$ , then  $x^l(t-1)$  is the closest center among the centers of  $x(t-1)$ . Since KMEANS halts at step  $t$ , then from (5.6) we have  $x(t) = x(t-1)$ , thus  $x^l(t)$  is the closest center to  $a$  among the centers of  $x(t)$ . Further we have

$$r(t) \leq \|x^j(t) - a\| - \|x^l(t) - a\| \quad \forall j \neq l. \quad (5.7)$$

Next, we show that  $a \in \underline{C^l}(x)$ , indeed

$$\begin{aligned} \|a - x^l\| - \|a - x^j\| &\leq \|a - x^l(t)\| + \|x^l(t) - x^l\| - (\|a - x^j(t)\| - \|x^j(t) - x^j\|) \\ &= \|a - x^l\| - \|a - x^j(t)\| + \|x^l(t) - x^l\| + \|x^j(t) - x^j\| \\ &< \|a - x^l\| - \|a - x^j(t)\| + r(t) \\ &\leq -r(t) + r(t) = 0, \end{aligned}$$

where the second inequality holds since  $x^l \in B\left(x^l(t), \frac{r(t)}{2}\right)$  and  $x^j \in B\left(x^j(t), \frac{r(t)}{2}\right)$ , and the third inequality follows from (5.7), and we get that  $C^l(t) \subseteq \underline{C^l}(x)$ . By definition of  $\underline{C^l}(x)$  we have that for any  $l \neq j$ ,  $\underline{C^l}(x) \cap \underline{C^j}(x) = \emptyset$ , and for all  $1 \leq l \leq k$ ,  $\underline{C^l}(x) \subseteq \mathcal{A}$ . Now, since  $C(t)$  is a partition of  $\mathcal{A}$ , then  $C^l(t) = \underline{C^l}(x)$  for all  $1 \leq l \leq k$ .  $\square$

**Proposition 5.3** (KMEANS converges to local minimum). *Let  $(C(t), x(t))$  be the clusters and centers KMEANS returns, then  $x(t)$  is local minimum of  $F$  in  $U = B\left(x^1(t), \frac{r(t)}{2}\right) \times B\left(x^2(t), \frac{r(t)}{2}\right) \times \cdots \times B\left(x^l(t), \frac{r(t)}{2}\right) \subset \mathbb{R}^{nk}$ .*

*Proof.* The minimum of  $F$  in  $U$  is

$$\min_{x \in U} F(x) = \min_{x \in U} \sum_{l=1}^k \sum_{a \in C^l(x)} \|a - x^l\|^2 = \min_{x \in U} \sum_{l=1}^k \sum_{a \in C^l(t)} \|a - x^l\|^2,$$

where the last equality follows from Proposition 5.2.

The function  $x \mapsto \sum_{l=1}^k \sum_{a \in C^l(t)} \|a - x^l\|^2$  is strictly convex, separable in  $x^l$  for all  $1 \leq l \leq k$ , and reaches its minimum at  $\frac{1}{|C^l(t)|} \sum_{a \in C^l(t)} a = \text{mean}(C^l(t)) = x^l(t)$ , and the result follows.  $\square$

## 6 Numeric Results

In this section we show the numeric results and compare the algorithms presented in this work with other algorithms that are commonly used to address the clustering problem.

The initialization points used within the implementation of the compared algorithms are as follows. KMEANS starting point is constructed by randomly choosing  $k$  different points from the dataset. The same technique is employed in the cases of KPALM and  $\varepsilon$ -KPALM, for the  $x(0)$  variable. Whereas for the  $w(0)$  variable, it is chosen at random from  $\Delta^m$ . KMEANS++ takes also part in our comparison, and it is basically the same as KMEANS, with the exception of its starting point that is constructed in the following manner. The first center  $x^1(0)$  is chosen randomly from the dataset  $\mathcal{A}$ . Suppose that  $1 \leq l < k$  centers have already been chosen, set  $x^{l+1}(0)$  to be the point in the dataset that is the furthest from its closest center.

Since it is impractical to compare the function values achieved with the algorithms which solve the squared Euclidean clustering problem with that of the algorithms which solve the Euclidean clustering problem, we used some criteria devised to compare clustering partitions. Criteria such as *variation of information (VI)*, *Mirkin metric*, and *Van Dongen metric*, are few examples for metrics that measure the difference between two clustering partitions (see [6]). With these metrics we compared the similarity of the partition achieved with each algorithm to the desired partition of each dataset. The goal is to decrease the value of the metrics.

### 6.1 Iris Dataset

We used the famous Iris dataset to test the performance of the KPALM algorithm. It is important to note that choosing the parameter  $\alpha$  is left to the user, and as presented below, has a significant effect on the convergence rate and the quality of the achieved clustering, namely the value of the objective function over the generated series. All the plots in this section are made by averaging over 100 trials, each trial with random starting point.

Figure 1 shows that dynamic values of the parameter  $\alpha$  which decreases fast, such as  $\alpha_i(t) = \frac{\text{diam}(\mathcal{A})}{2^{i-1}}$ , achieve smaller function values. In Figure 2 we made a comparison between KPALM with dynamic rule for choosing the parameter  $\alpha$ , that is  $\alpha_i(t) = \frac{\text{diam}(\mathcal{A})}{2^{i-1}}$ , with KMEANS and KMEANS++. It demonstrates that KPALM can reach lower objective function values than KMEANS, and these are similar to the values achieved with KMEANS++. In addition, the KPALM++ are the objective function values achieved with KPALM when

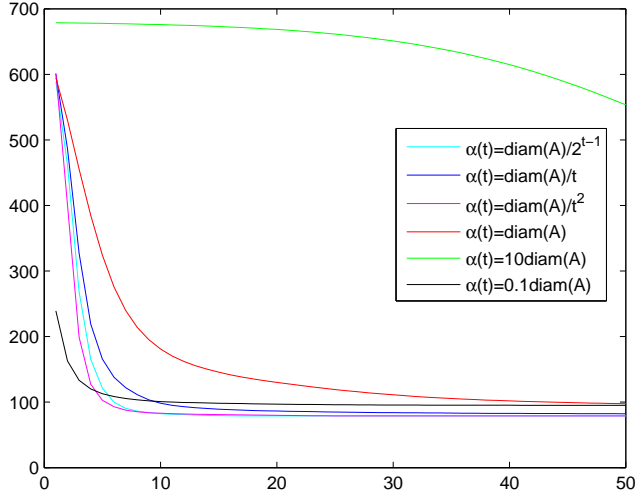


Figure 1: Comparison of the objective values for different values of  $\alpha$ .

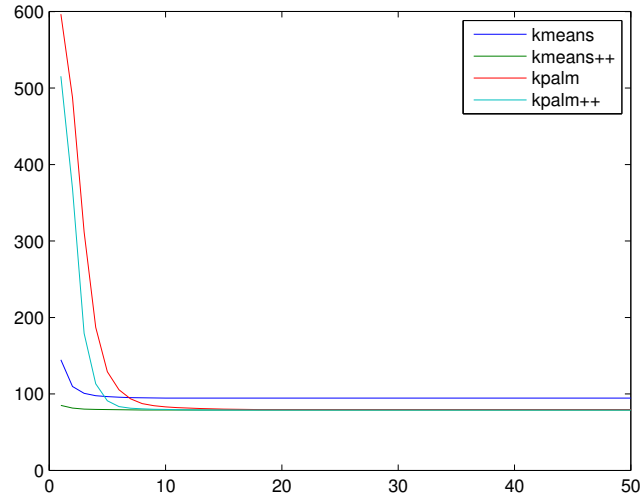
the  $x$  variable is initialized as in KMEANS++. Unlike KMEANS, the objective function values KPALM converge to are more stable and less sensitive to its starting point.

Figure 3 shows the number of iteration needed to reach precision of  $1e-3$  between consecutive objective function values. Similarly to Figure 1, in Figure 4 we can see a comparison of the objective values of  $\Psi_\varepsilon$  for different function values. The value of  $\varepsilon$  is set to be  $1e-5$ .

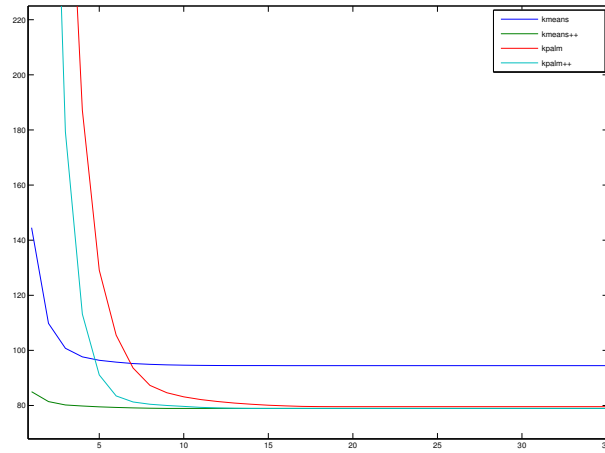
## 6.2 Synthetic Dataset

In this section we show that  $\varepsilon$ -KPALM is less sensitive to outliers in the data verses algorithms that suit the squared Euclidean norm (e.g., KMEANS, KMEANS++ and KPALM). We generated two synthetic datasets, each contains 300 points in the plane, by sampling three two-dimensional Gaussian, 100 samples each. In Figure 5(5a) the clusters are denser than in Figure 5(5b). Then we run the clustering algorithms and compared their clustering results, namely, how many points were clustered correctly. From Figure 6(6a) it is evident that KMEANS is superior to other algorithms in the dense case and  $\varepsilon$ -KPALM is quite sensitive. Whereas, in the sparse case in Figure 6(6b),  $\varepsilon$ -KPALM is superior, and less sensitive to outliers. In Figure 7 we compare the distance of clusterings achieved with different algorithms to the desired clustering, where kpalm1, kpalm2 and kpalm3 match using  $\alpha(t) = \text{diam}(\mathcal{A})/2^{t-1}$ ,  $\alpha(t) = \text{diam}(\mathcal{A})/t^2$  and  $\alpha(t) = \text{diam}(\mathcal{A})$  respectively, and similarly for  $\varepsilon$ -kpalmi,  $i \in \{1, 2, 3\}$ . In Figure 7(7a) we witness that for dense dataset, the resulting clusterings of squared Euclidean algorithms, namely, KMEANS, KMEANS++ and KPALM, are superior to the clustering  $\varepsilon$ -KPALM, where KPALM with  $\alpha(t) = \text{diam}(\mathcal{A})/2^{t-1}$  gives the best result, that is, the clustering in this setting is the closest to the desired clustering. Whereas in the sparse dataset, the clustering achieved with  $\varepsilon$ -KPALM with  $\alpha(t) = \text{diam}(\mathcal{A})/t^2$  is the closest to the desired clustering, as reflected from Figure 7(7b).



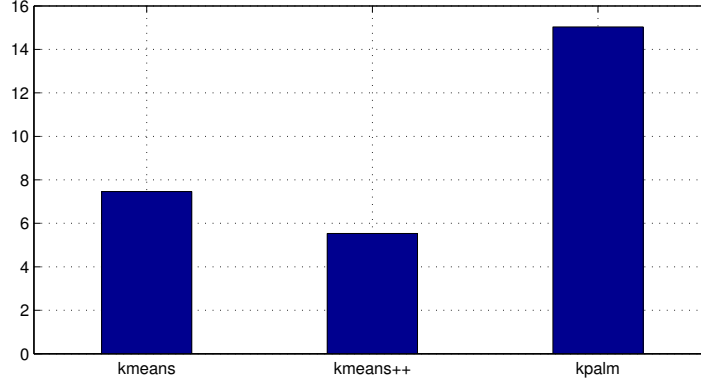


(a) Comparison of objective function values.

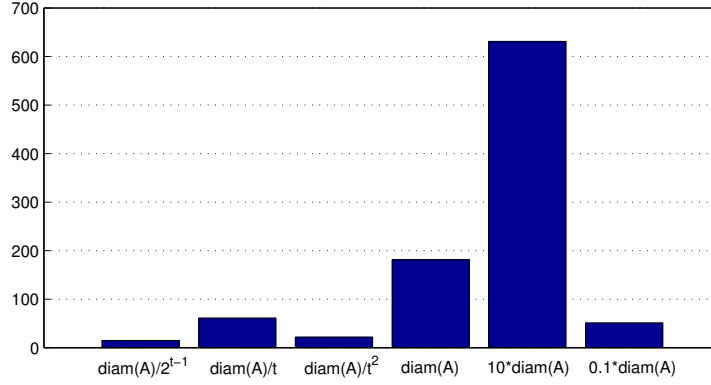


(b) Zoom of Figure 2a.

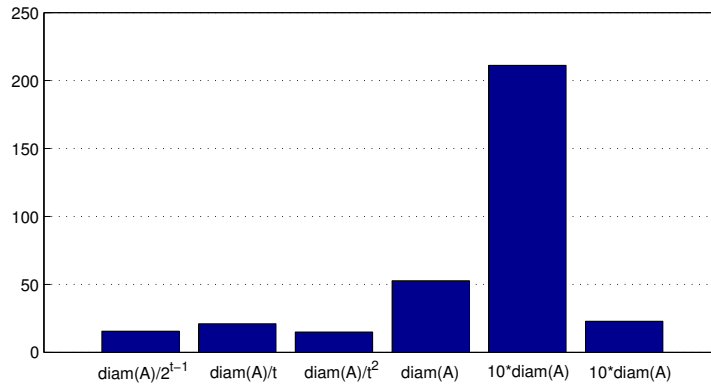
Figure 2: Comparison of objective function values for KMEANS, KMEANS++, KPALM and KMEANS++.



(a) Number of iterations of KMEANS, KMEANS++ and KPALM with  $\alpha(t) = \text{diam}(\mathcal{A})/2^{t-1}$ .



(b) Number of iterations of KPALM with different updates of  $\alpha(t)$ .



(c) Number of iterations of  $\varepsilon$ -KPALM with different updates of  $\alpha(t)$ .

Figure 3: Comparison of number of iterations needed to reach  $1e-3$  precision of  $\Psi$ .

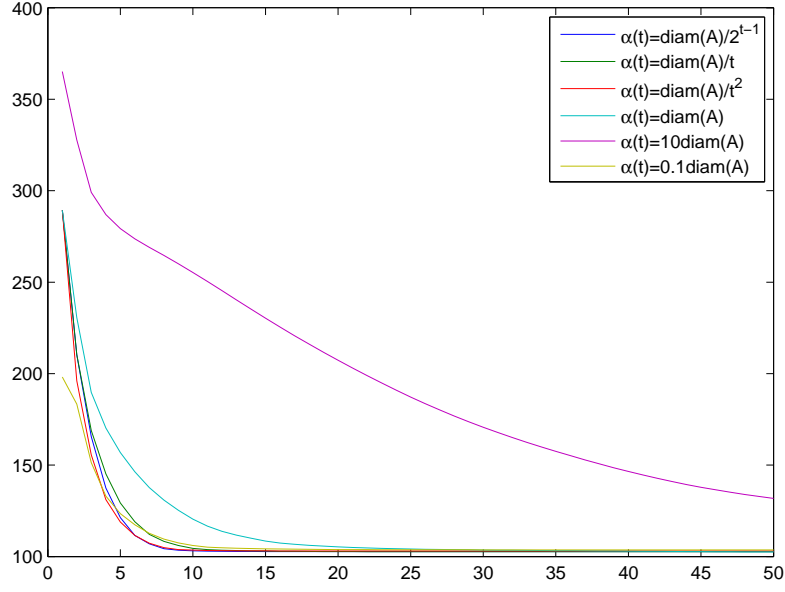
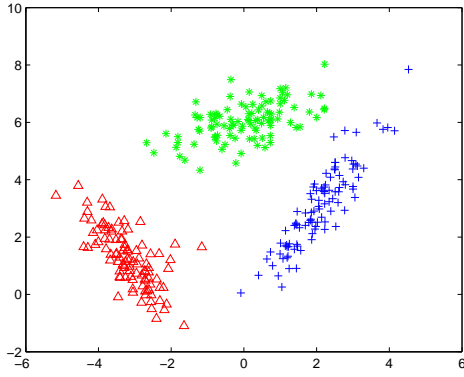
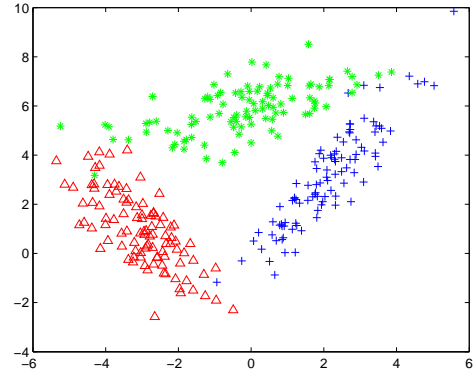


Figure 4: Comparison of the objective values for different values of  $\alpha$ .

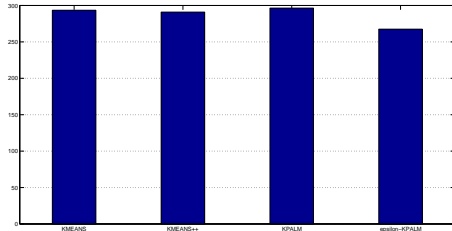


(a) Dense Gaussians.

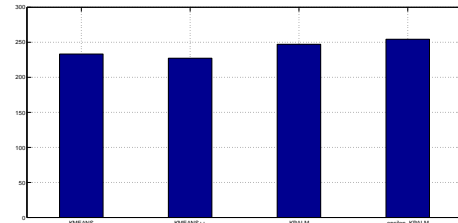


(b) Sparse Gaussians.

Figure 5: Two datasets, each 300 points.

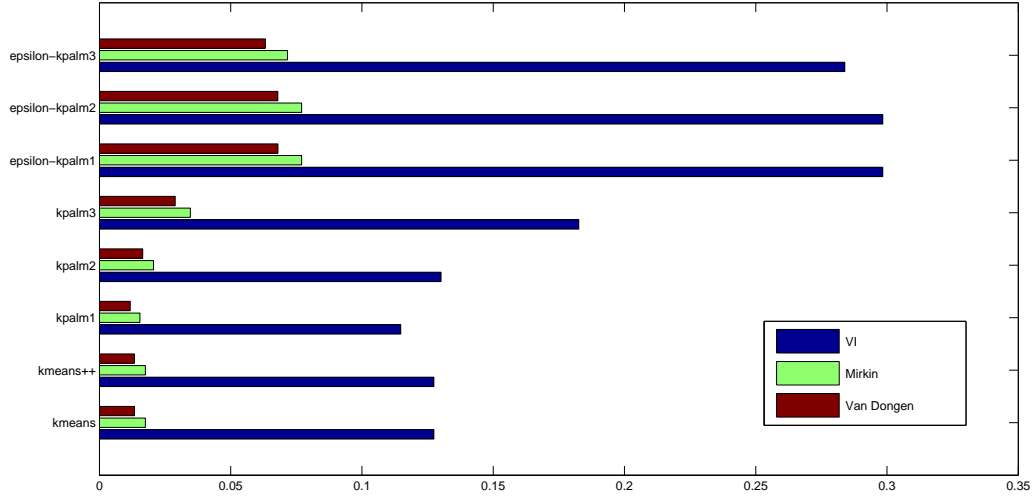


(a) Dense Gaussians clustering.

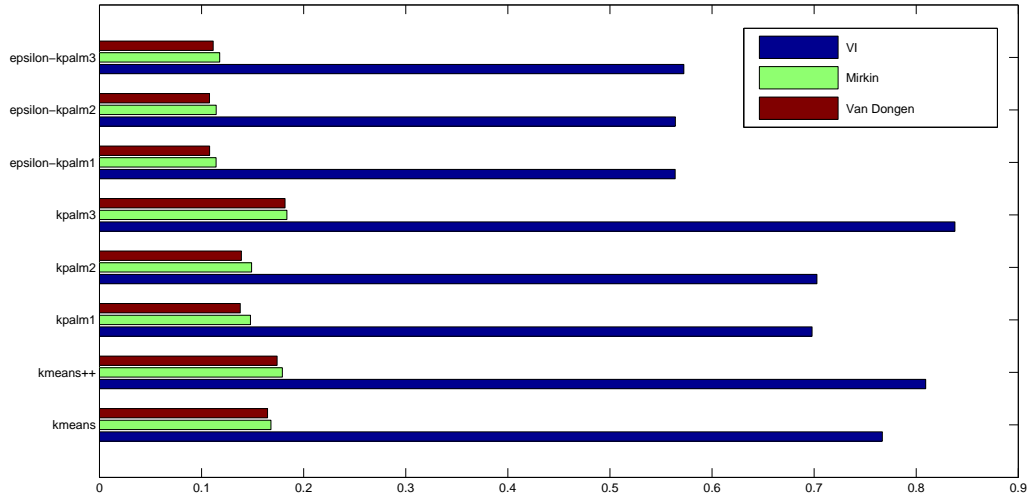


(b) Sparse Gaussians clustering.

Figure 6: Results of clustering algorithms for dense and sparse datasets.



(a) Dense Gaussians metrics comparison.



(b) Sparse Gaussians metrics comparison.

Figure 7: Comparison of metrics between clusterings for dense and sparse datasets.

## References

- [1] Attouch, H., Bolte, J.: On the convergence of the proximal algorithm for nonsmooth functions involving analytic features. *Math. Program.* **116**, 5-16 (2009)
- [2] Attouch, H., Bolte, J., Svaiter, B.F.: Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forwardbackward splitting, and regularized GaussSeidel methods. *Math. Program.* **137**, 91-129 (2013)
- [3] Beck, A., Sabach, S.: Weiszfelds Method: Old and New Results. *Journal of Optimization Theory and Applications* **164**, 1-40 (2015)
- [4] Ben-Tal, A., Teboulle, M., Yang, W.H.: A least-squares-based method for a class of nonsmooth minimization problems with applications in plasticity. *Applied Mathematics and Optimization* **24**, 273-288 (1991)
- [5] Bolte, J., Sabach, S., Teboulle, M.: Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Math. Program.* **146**, 459-494 (2014)
- [6] Meila, M.: Comparing Clusterings An Axiomatic View. In *ICML '05: Proceedings of the 22nd international conference on Machine Learning* ACM, New York, 577-584 (2005)
- [7] Selim, S.Z., Ismail, M. A.: K-Means-Type Algorithms: A Generalized Convergence Theorem and Characterization of Local Optimality. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **6**, 81-87 (1984)
- [8] Teboulle, M.: A Unified Continuous Optimization Framework for Center-Based Clustering Methods. *The Journal of Machine Learning Research* **8**, 65-102 (2007)