

# 1 The Clustering Problem

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 data  $\mathcal{A}$  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  $x^l \in \mathbb{R}^n$ , and we are interested to determine  $k$  cluster centers  $\{x^1, x^2, \dots, x^k\}$  such that the sum of certain proximity measures from each point  $a^i$ ,  $i = 1, 2, \dots, m$ , to a nearest cluster center  $x^l$  is minimized. We define the vector of all centers by  $x = (x^1, x^2, \dots, x^k) \in \mathbb{R}^{nk}$ .

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.

## 2 Problem Reformulation and Notations

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 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

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

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

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

### 3 Clustering: The Squared Euclidean Norm Case

#### 3.1 Introduction to the PALM Theory

Presentation of PALM's requirements and of the algorithm steps ...

#### 3.2 Clustering with PALM

In this section we tackle the clustering problem, given in (2.3), with 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 PALM.
- (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 PALM.

As in the PALM algorithm, our algorithm is based on 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

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

On the other hand, with respect to  $x$  we perform exact minimization

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

It is easy to check that all subproblems, with respect to  $w^i$ ,  $i = 1, 2, \dots, m$ , and  $x$ , can be written explicitly as follows:

$$w^i(t+1) = P_{\Delta} \left( w^i(t) - \frac{d^i(x(t))}{\alpha_i(t)} \right), \quad i = 1, 2, \dots, m, \quad (3.3)$$

where  $P_{\Delta}$  is the orthogonal projection onto the set  $\Delta$ , and

$$x^l(t+1) = \frac{\sum_{i=1}^m w_l^i(t+1) a^i}{\sum_{i=1}^m w_l^i(t+1)}, \quad l = 1, 2, \dots, k. \quad (3.4)$$

Therefore we can record now the suggested 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

$$w^i(t+1) = P_\Delta \left( w^i(t) - \frac{d^i(x(t))}{\alpha_i(t)} \right). \quad (3.5)$$

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

$$x^l(t+1) = \frac{\sum_{i=1}^m w_l^i(t+1) a^i}{\sum_{i=1}^m w_l^i(t+1)}. \quad (3.6)$$

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  $\text{Conv}(\mathcal{A})$ , 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) 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

$$x^l(t) = \frac{\sum_{i=1}^m w_l^i(t) a^i}{\sum_{i=1}^m w_l^i(t)} = \sum_{i=1}^m \left( \frac{w_l^i(t)}{\sum_{j=1}^m w_l^j(t)} \right) a^i = \sum_{i=1}^m \lambda_i a^i \in \text{Conv}(\mathcal{A}). \quad (3.7)$$

Hence  $x^l(t)$  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 (3.7) 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

$$\underline{\alpha}_i \leq \alpha_i(t) \leq \overline{\alpha}_i, \quad \forall t \in \mathbb{N}. \quad (3.8)$$

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

$$2 \min_{1 \leq l \leq k} \sum_{i=1}^m w_l^i(t) := \beta(w(t)) \geq \underline{\beta}. \quad (3.9)$$

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 involved in the objective function.

**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 (3.9), whenever  $\beta(w) > 0$ .*

*Proof.* Since the function  $x \mapsto H(w(t), x) = \sum_{l=1}^k \sum_{i=1}^m w_l^i \|x^l - a^i\|^2$  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 the descent property of the KPALM algorithm.

**Proposition 3.2** (Sufficient decrease property). *Suppose that Assumption 1 holds true and let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by KPALM. 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 (3.5), 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

$$\frac{\alpha_i(t)}{2} \|w^i(t+1) - w^i(t)\|^2 \leq H^i(w(t), x(t)) - H^i(w(t+1), x(t)). \quad (3.10)$$

Denote  $\underline{\alpha} = \min_{1 \leq i \leq m} \alpha_i$ . Summing inequality (3.10) over  $i = 1, 2, \dots, m$  yields

$$\begin{aligned} \frac{\underline{\alpha}}{2} \|w(t+1) - w(t)\|^2 &= \frac{\underline{\alpha}}{2} \sum_{i=1}^m \|w^i(t+1) - w^i(t)\|^2 \\ &\leq \sum_{i=1}^m \frac{\alpha_i(t)}{2} \|w^i(t+1) - w^i(t)\|^2 \\ &\leq \sum_{i=1}^m [H^i(w(t), x(t)) - H^i(w(t+1), x(t))] \\ &= H(w(t), x(t)) - H(w(t+1), x(t)), \end{aligned} \quad (3.11)$$

where the first inequality follows from Assumption 1(i).

From Assumption 1(ii) we have that  $\beta(w(t)) \geq \underline{\beta}$ , and from Lemma 3.1.1 it follows that the function  $x \mapsto H(w(t), x)$  is strongly convex with parameter  $\beta(w(t))$ , hence it follows that

$$\begin{aligned} H(w(t+1), x(t)) - H(w(t+1), x(t+1)) &\geq \\ &\geq \langle \nabla_x H(w(t+1), x(t+1)), x(t) - x(t+1) \rangle + \frac{\beta(w(t))}{2} \|x(t) - x(t+1)\|^2 \\ &= \frac{\beta(w(t))}{2} \|x(t+1) - x(t)\|^2 \\ &\geq \frac{\underline{\beta}}{2} \|x(t+1) - x(t)\|^2, \end{aligned} \quad (3.12)$$

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 (3.11) and (3.12), 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  $G(w(t)) = 0$  for all  $t \in \mathbb{N}$  and therefore  $H(z(t)) = \Psi(z(t))$ ,  $t \in \mathbb{N}$ .  $\square$

Now, we aim to prove the subgradient lower bound for the iterates gap. The following lemma will be essential in our proof.

**Lemma 3.2.1.** *Let  $\{z(t)\}_{t \in \mathbb{N}}$  be the sequence generated by KPALM, then*

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

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

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

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

this proves the desired result.  $\square$

**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

$$\begin{aligned}
\partial\Psi(z(t+1)) &= \\
&= \left( (\nabla_{w^i} H^i(w(t+1), x(t+1)) + \partial_{w^i} \delta_\Delta(w^i(t+1)))_{i=1,2,\dots,m}, \nabla_x H(w(t+1), x(t+1)) \right) \\
&= \left( (d^i(x(t+1)) + \partial_{w^i} \delta_\Delta(w^i(t+1)))_{i=1,2,\dots,m}, \nabla_x H(w(t+1), x(t+1)) \right) \\
&= \left( (d^i(x(t+1)) + \partial_{w^i} \delta_\Delta(w^i(t+1)))_{i=1,2,\dots,m}, \mathbf{0} \right), \tag{3.13}
\end{aligned}$$

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

$$d^i(x(t)) + \alpha_i(t) (w^i(t+1) - w^i(t)) + u^i(t+1) = \mathbf{0}. \tag{3.14}$$

Setting  $\gamma(t+1) := \left( (d^i(x(t+1)) + u^i(t+1))_{i=1,2,\dots,m}, \mathbf{0} \right)$  and from (3.13) it follows that  $\gamma(t+1) \in \partial\Psi(z(t+1))$ . Using (3.14) 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

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

where the third inequality follows from Lemma 3.2.1. Define  $\rho_2 = m(4M + \bar{\alpha})$ , 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 3.1, we have used 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

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

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

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

Note that  $\Psi_{\varepsilon}(z)$  is a perturbed form of  $\Psi(z)$  for a small  $\varepsilon > 0$ , and obviously  $\Psi_0(z) = \Psi(z)$ .

Now we would like to develop an algorithm which is based on the methodology of PALM 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 and adding regularizing term

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

where

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

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

$$w^i(t+1) = P_\Delta \left( w^i(t) - \frac{d_\varepsilon^i(x(t))}{\alpha_i(t)} \right). \quad (4.3)$$

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

$$x^l(t+1) = x^l(t) - \frac{1}{L_\varepsilon^l(w(t+1), x(t))} \nabla_{x^l} H_\varepsilon(w(t+1), x(t)). \quad (4.4)$$

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

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



Before we will be able to prove the two properties 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 positive numbers  $v_1, v_2, \dots, v_m \in \mathbb{R}$  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 operator,  $L_\varepsilon : \mathbb{R}^n \rightarrow \mathbb{R}$  defined by

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

**Lemma 4.0.1** (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  $x, y \in \mathbb{R}^n$ ,

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

*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 can prove that the function  $f_\varepsilon$  has Lipschitz continuous gradient.

**Lemma 4.0.2.** *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 two functions

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

and

$$\tilde{h}_\varepsilon(x, y) = h_\varepsilon(x, y) - \langle \nabla f_\varepsilon(z), x \rangle.$$

It is clear that  $x \mapsto \tilde{h}_\varepsilon(x, y)$  is also a quadratic function with associated matrix  $L_\varepsilon(y)\mathbf{I}$ . Therefore, from Lemma 4.0.1(i) we can write

$$\begin{aligned} \tilde{h}_\varepsilon(x, y) &= \tilde{h}_\varepsilon(y, y) + \left\langle \nabla_x \tilde{h}_\varepsilon(y, y), x - y \right\rangle + L_\varepsilon(y) \|x - y\|^2 \\ &= \tilde{f}_\varepsilon(y) + \langle 2\nabla f_\varepsilon(y) - \nabla f_\varepsilon(z), x - y \rangle + L_\varepsilon(y) \|x - y\|^2. \end{aligned} \quad (4.5)$$

On the other hand, from Lemma 4.0.1(ii) we have that

$$\begin{aligned} \tilde{h}_\varepsilon(x, y) &= h_\varepsilon(x, y) - \langle \nabla f_\varepsilon(z), x \rangle \geq 2f_\varepsilon(x) - f_\varepsilon(y) - \langle \nabla f_\varepsilon(z), x \rangle \\ &= 2\tilde{f}_\varepsilon(x) - \tilde{f}_\varepsilon(y) + \langle \nabla f_\varepsilon(z), x - y \rangle, \end{aligned} \quad (4.6)$$

where the last equality follows from the definition of  $\tilde{f}_\varepsilon$ . Combining (4.5) and (4.6) yields

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

Dividing the last inequality by 2 leads to

$$\tilde{f}_\varepsilon(x) \leq \tilde{f}_\varepsilon(y) + \left\langle \nabla \tilde{f}_\varepsilon(y), x - y \right\rangle + \frac{L_\varepsilon(y)}{2} \|x - y\|^2. \quad (4.7)$$

It is clear that the optimal point of  $\tilde{f}_\varepsilon$  is  $z$  since  $\nabla \tilde{f}_\varepsilon(z) = 0$ , therefore using (4.7) 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.  $\square$

Now we get back to  $\varepsilon$ -KPALM algorithm and prove few technical results about the involved functions which are based on the auxiliary results 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  $d_{\mathcal{A}}$  is the diameter of  $\text{Conv}(\mathcal{A})$  and  $\underline{\beta}$  is given in (3.9).*

(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

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

as asserted.

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

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

as asserted.  $\square$

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  $v_i = w_l^i(t+1)$ ,  $i = 1, 2, \dots, m$ . Therefore, by applying Lemma 4.0.1(iii) with  $x = x^l(t+1)$  and  $y = x^l(t)$ , we get

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

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

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

Replacing the last term with the following compact form

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

and the result follows.  $\square$

Now we are finally ready to prove the two properties needed for guaranteeing that the sequence 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 steps with respect to  $w$  of KPALM and  $\varepsilon$ -KPALM are similar and therefore following the same arguments given at the beginning of the proof of Proposition 3.2 we have that

$$\frac{\alpha}{2} \|w(t+1) - w(t)\|^2 \leq H_\varepsilon(w(t), x(t)) - H_\varepsilon(w(t+1), x(t)), \quad (4.8)$$

where  $\underline{\alpha} = \min_{1 \leq i \leq m} \alpha_i$ . Applying Proposition 4.2 with (4.4) we get for all  $t \in \mathbb{N}$  that

$$\begin{aligned}
H_\varepsilon(w(t+1), x(t)) - H_\varepsilon(w(t+1), x(t+1)) &\geq \sum_{l=1}^k \frac{L_\varepsilon^l(w(t+1), x(t))}{2} \|x^l(t+1) - x^l(t)\|^2 \\
&\geq \frac{\underline{\beta}}{(d_{\mathcal{A}}^2 + \varepsilon^2)^{1/2}} \sum_{l=1}^k \|x^l(t+1) - x^l(t)\|^2 \\
&\geq \frac{\underline{\beta}}{(d_{\mathcal{A}}^2 + \varepsilon^2)^{1/2}} \|x(t+1) - x(t)\|^2, \tag{4.9}
\end{aligned}$$

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 (4.8) and (4.9) 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  $G(w(t)) = 0$ , for all  $t \in \mathbb{N}$ . This proves the desired result.  $\square$

The next lemma 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 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  $d_{\mathcal{A}} = \text{diam}(\text{Conv}(\mathcal{A}))$ .

*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

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

Hence,

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

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 yields that

$$\gamma(t+1) := \left( (d_\varepsilon^i(x(t+1)) + u^i(t+1))_{i=1, \dots, m}, \nabla_x H_\varepsilon(w(t+1), x(t+1)) \right) \in \partial \Psi_\varepsilon(z(t+1)), \quad (4.10)$$

where for all  $1 \leq i \leq m$ ,  $u^i(t+1) \in \partial \delta_\Delta(w^i(t+1))$  such that

$$d_\varepsilon^i(x(t)) + \alpha_i(t) (w^i(t+1) - w^i(t)) + u^i(t+1) = \mathbf{0}. \quad (4.11)$$

Plugging (4.11) into (4.10), and taking the norm yields

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

where the last inequality follows from Lemma 4.3.1 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

$$\begin{aligned}
\nabla_{x^l} H_\varepsilon(w(t+1), x(t+1)) &= \nabla_{x^l} H_\varepsilon(w(t+1), x(t+1)) - \nabla_{x^l} H_\varepsilon(w(t+1), x(t)) \\
&\quad + \nabla_{x^l} H_\varepsilon(w(t+1), x(t)) \\
&= \nabla_{x^l} H_\varepsilon(w(t+1), x(t+1)) - \nabla_{x^l} H_\varepsilon(w(t+1), x(t)) \\
&\quad + L_\varepsilon^l(w(t+1), x(t)) (x^l(t) - x^l(t+1)), \quad (4.12)
\end{aligned}$$

where the last equality follows from (4.4). Therefore,

$$\begin{aligned}
\|\nabla_x H_\varepsilon(w(t+1), x(t+1))\| &\leq \sum_{l=1}^k \|\nabla_{x^l} H_\varepsilon(w(t+1), x(t+1))\| \\
&\leq \sum_{l=1}^k L_\varepsilon^l(w(t+1), x(t)) \|x^l(t+1) - x^l(t)\| \\
&\quad + \sum_{l=1}^k \|\nabla_{x^l} H_\varepsilon(w(t+1), x(t+1)) - \nabla_{x^l} H_\varepsilon(w(t+1), x(t))\| \\
&\leq \frac{m}{\varepsilon} \sum_{l=1}^k \|x^l(t+1) - x^l(t)\| + \sum_{l=1}^k \gamma^l(t) \|x^l(t+1) - x^l(t)\|, \quad (4.13)
\end{aligned}$$

where the last inequality follows from Proposition 4.1(ii) and Lemma 4.0.2 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 (4.13), we have

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

Therefore, setting  $\rho_2 = \sqrt{m} \left( \frac{d_A}{\varepsilon} + \bar{\alpha} \right) + \frac{2m\sqrt{k}}{\varepsilon}$ , yields the result.  $\square$

The following lemma shows that the smoothed function  $H_\varepsilon(w, x)$  indeed approximates  $H(w, x)$ .

**Lemma 4.4.1** (Closeness of smooth). *For any  $(w, x) \in \Delta^m \times \mathbb{R}^{nk}$  and  $\varepsilon > 0$  the following inequalities hold true*

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

*Proof.* Applying the inequality

$$(a+b)^\lambda \leq a^\lambda + b^\lambda, \quad \forall a, b \geq 0, \lambda \in (0, 1],$$

with  $a = \|x^l - a^i\|^2$ ,  $b = \varepsilon^2$  and  $\lambda = \frac{1}{2}$ , yields

$$\left( \|x^l - a^i\|^2 + \varepsilon^2 \right)^{1/2} \leq \|x^l - a^i\| + \varepsilon, \quad \forall 1 \leq l \leq k, 1 \leq i \leq m.$$

Together with the fact that

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

yields the following inequality

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

for all  $l = 1, 2, \dots, k$  and  $i = 1, 2, \dots, m$ . 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. □



## 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) Centers 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 employ the PALM theory once more 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 breaks down 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 5.0.2.** *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 (3.12)), 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 5.0.2, 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

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

(2.2) Centers 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  $\overline{C}^l(x) = \{a \in \mathcal{A} \mid \|a - x^l\| \leq \|a - x^j\|, \quad \forall j \neq l\}$ , for all  $1 \leq l \leq k$ , and the sub-partition of  $\mathcal{A}$  by  $\underline{C}^l(x) = \{a \in \mathcal{A} \mid \|a - x^l\| < \|a - x^j\|, \quad \forall j \neq l\}$ . 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  $U = B\left(x^1(t), \frac{r(t)}{2}\right) \times B\left(x^2(t), \frac{r(t)}{2}\right) \times \dots \times B\left(x^k(t), \frac{r(t)}{2}\right)$  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  $U = B\left(x^1(t), \frac{r(t)}{2}\right) \times B\left(x^2(t), \frac{r(t)}{2}\right) \times \cdots \times B\left(x^k(t), \frac{r(t)}{2}\right) \subset \mathbb{R}^{nk}$ .*

*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.

### 6.1 Iris Dataset

We use 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.

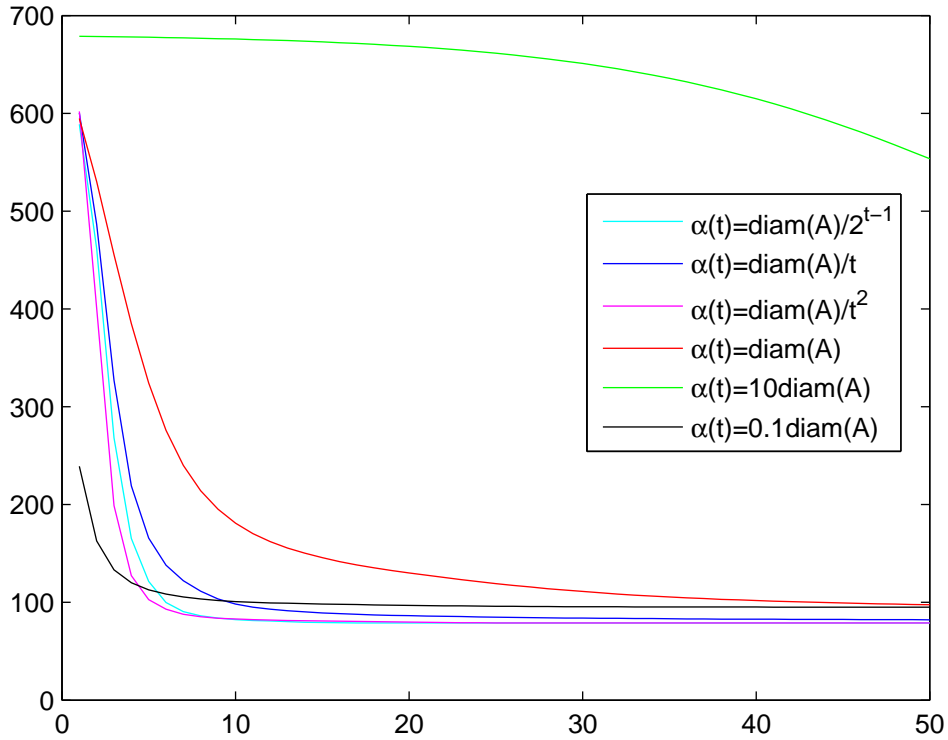


Figure 1: Comparison of the objective values for different values of alpha.

Figure 1 shows that dynamic values of the parameter  $\alpha$  which decreases fast, such as  $\alpha_i(t) = \frac{\text{diam}(\mathcal{A})}{2^{t-1}}$ , achieve smaller function values, and increase the rate of conversion.

In Figure 2 we made a comparison between KPALM with dynamic rule for choosing the parameter  $\alpha$ ,  $\alpha_i(t) = \frac{\text{diam}(\mathcal{A})}{2^{t-1}}$ , 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++.

Figure 3 shows the number of iteration needed to reach precision  $1e-3$  between consecutive

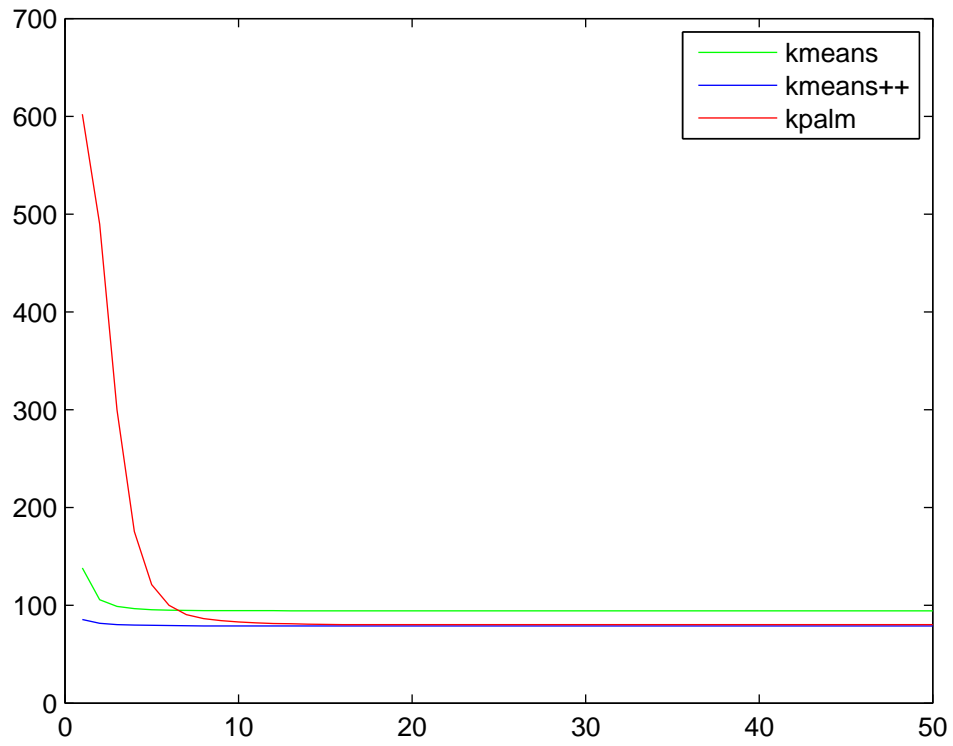
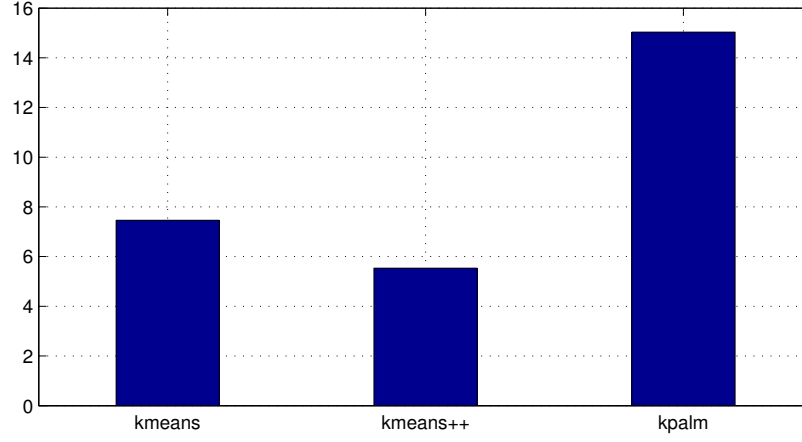
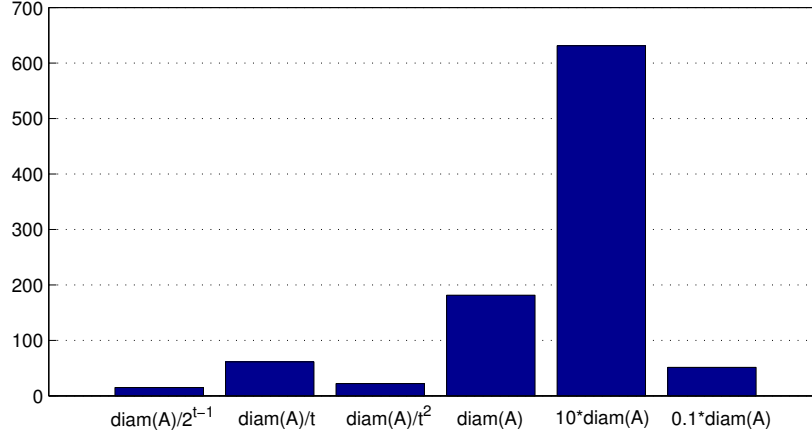


Figure 2: Comparison of objective function values for KMEANS, KMEANS++ and KPALM.

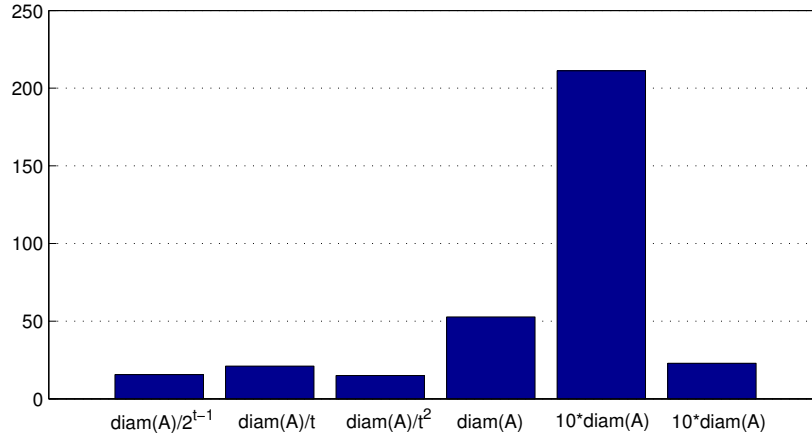
objective function values.



(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  $\Psi$ .