

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 $l=1,2,\ldots,k$, the cluster C^l is represented by its center x^l and we want to determine k cluster centers $\{x^1, x^2, \dots, x^k\}$ such that the sum of proximity measures from each point $a^i, i = 1, 2, \dots, m$, to a nearest cluster center x^i is minimized.

The clustering problem is given by

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

We benote the voctor

cortain

<R

, which easing

in Δ and oc otherwisi

Problem Reformulation and Notations

We introduce some notations that will be used throughout this document.

 \mathbb{R}^{nm} , where a^i

Let $S \subseteq \mathbb{R}^n$. The indicator function of S is defined and denoted as follows δ_S

 $\min_{1 \le l \le k} u_l = \min \left\{ \langle u, v \rangle \mid v \in \right.$ Δ , and applying it over (1.1), gives a smooth reformulation of the clustering problem

whore $\mathcal{L}^{1}(\mathbf{x}) = \left(\mathcal{L}(\mathbf{x}^{1}, \mathbf{a}^{i}), \mathcal{L}(\mathbf{x}^{1}, \mathbf{a}^{i}), \ldots\right) \quad \min_{x \in \mathbb{R}^{nk}} \sum_{i=1}^{m} \min_{w^{i} \in \Delta} \langle w^{i}, d^{i}(x) \rangle_{\mathbf{x}} = \left(\mathcal{L}(\mathbf{x}^{1}, \mathbf{a}^{i}), \mathcal{L}(\mathbf{x}^{1}, \mathbf{a}^{i}), \ldots\right)$ $\Delta(\mathbf{x}^{K}, \mathbf{a}^{i})) \in \mathbb{R}^{m}$ i=1,2,...,m

Replacing further the constraint $w^i \in \Delta$ by adding the indicator function $\delta_{\Delta}(\cdot)$ to the objective function, results in a equivalent formulation

Nhere

 $\left\{\sum_{i=1}^{m}\langle w^{i},d^{i}(x)
angle+\delta_{\Delta}(w^{i})
ight\}$

For the simplicity of the yet to come expositions, we cofine the Finally, introducing several more useful notations is needed. For each $i=1,2,\ldots,m$, we denote following function m, we denote following functions

We begin with a reformulation the clustering problem which be the basis for our boxelopments work. The resormulation is

based on the following fact:

min ul= min {< u, v> : ve A5,

where A is the well-known, simplex defined by

(2.1)

Δ - Juer : ∑u=1, u>0}

Using this fact in Problem (1.1) and introducing

Variables in eR i=1,2, m

 $H(w,x):=\sum\limits_{i=1}^m H_i(w,x)=\sum\limits_{i=1}^m \langle w^i,d^i(x)
angle ext{ and } G(w)=\sum\limits_{i=1}^m G(w^i):=\sum\limits_{i=1}^m \delta_{\Delta}(w^i).$

ቌ፧፞የውጥ ነጥ ___(2,3),

Replacing the terms in (2.1) with the functions defined above gives a compact form of the original clustering problem

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

3 Clustering via PALM Approach

3.1 Introduction to PALM Theory

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

3.2 Clustering with PALM for Squared Euclidean Norm Distance Function

In this section we tackle the clustering problem with the classical distance function defined by $d(u,v) = ||u-v||^2$. We devise a PALM-like algorithm, based on the discussion about PALM in the previous subsection. Since the clustering problem has a specific structure, we are ought to exploit it in the following manner. First we notice that the function $w\mapsto H(w,x)$ is linear in w, so there w

is no need to linearize it. In addition, the function $x \mapsto H(w,x) = \sum_{i=1}^m \sum_{l=1}^k w_l^i ||x^l - a^i||^2$

 $\sum_{l=1}^k \sum_{i=1}^m w_l^i \| w^l - a^i \|^2 \text{ is convex-and-quadratic in } x, \text{ hence we do not need to add a proximal term as in PALM algorithm.}$

· Now we propose a PALM-like algorithm for clustering, which we call KPALM.

- (1) Initialization: Set t=0, and pick random vectors $(w(0), x(0)) \in \Delta^m \times \mathbb{R}^{nk}$.
- (2) For each $t = 0, 1, \dots$ generate a sequence $\{(w(t), x(t))\}_{t \in \mathbb{N}}$ as follows:
 - (2.1) Cluster Assignment: Take any $\alpha_i(t) > 0$ and for each i = 1, 2, ..., m compute

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

(2.2) Centers Update: For each $l=1,2,\ldots,k$ compute $x^l\in\mathbb{R}^n$ via

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

The Sunction $w \to H(w,x)$, for fixed x, is linear and therefore there is no need to linearize it as suggested in PALM.

The function $x \longrightarrow H(w,x)$, for fixed w, is quadratic and convolve. Hence, there is no need to 2 add a quadratic proximal term a suggested in PALM.

in the PALM organithm, our algorithm is based on altermiting minimization

Sollowing adaptations which are notivated

Assumption 2. For any step $t \in \mathbb{N}$, 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 $1 \le l \le k$, $t \in \mathbb{N}$, we also have that r(t) > 0 for all $t \in \mathbb{N}$.

Proposition 3.3. Let (C(t), x(t)) be the clusters and centers KMEANS returns. Denote an open neighbourhood of x(t) by $U = B\left(x^1(t), \frac{r(t)}{2}\right) \times B\left(x^2(t), \frac{r(t)}{2}\right) \times \cdots \times B\left(x^l(t), \frac{r(t)}{2}\right)$, then for any $x \in U$ we have $\underline{C}^l(x) = C^l(t)$ for all $1 \leq l \leq k$. 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 (3.12) 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

$$r(t) \le ||x^{j}(t) - a|| - ||x^{l}(t) - a|| \quad \forall j \ne l.$$
 (3.13)

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 (3.13), 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 $1 \leq l \leq k$.

Proposition 3.4 (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 3.3.

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 $1 \le l \le k$, and reaches its minimum at $(x^l)^* = \frac{1}{|C^l(t)|} \sum_{a \in C^l(t)} a = mean(C^l(t)) = x^l(t)$, and the result follows. \square

$(_{3.3})$

Similarity to KMEANS

The famous KMEANS algorithm has close proximity to KPALM algorithm. KMEANS alternates between cluster assignments and center updates as well. In detail, we can write its steps in the following manner

- (1) Initialization: Set t = 0, and pick random centers $y(0) \in \mathbb{R}^{nk}$.
- (2) For each t = 0, 1, ... generate a sequence $\{(v(t), y(t))\}_{t \in \mathbb{N}}$ as follows:
 - (2.1) Cluster Assignment: For i = 1, 2, ..., m compute

$$v^{i}(t+1) = \arg\min_{v^{i} \in \Delta} \left\{ \langle v^{i}, d^{i}(y(t)) \rangle \right\}. \tag{3.8}$$

(2.2) Center Update: For l = 1, 2, ..., k compute

$$y^{l}(t+1) = \frac{\sum_{i=1}^{m} v_{l}^{i}(t+1)a^{i}}{\sum_{i=1}^{m} v_{l}^{i}(t+1)}.$$
 (3.9)

The KMEANS algorithm obviously resemble KPALM algorithm. Denote $\overline{\alpha}(t) = \max_{1 \le i \le m} \alpha_i(t)$. Assuming same starting point x(0) = y(0) and by taking $\overline{\alpha}(t) \to 0$, we have

$$v(t) = \lim_{\overline{\alpha}(t) \to 0} w(t), \quad y(t) = \lim_{\overline{\alpha}(t) \to 0} x(t),$$

meaning, both algorithms converge to the same result.



KMEANS Convergence Proof

We start with rewriting the KMEANS algorithms, in its most familiar form

- (1) Initialization: Set t = 0, and pick random centers $x(0) \in \mathbb{R}^{nk}$.
- (2) For each t = 0, 1, ... generate a sequence $\{(C(t), x(t))\}_{t \in \mathbb{N}}$ as follows:
 - (2.1) Cluster Assignment: For i = 1, 2, ..., m compute

$$C^{l}(t+1) = \left\{ a \in \mathcal{A} \mid ||a - x^{l}(t)|| \le ||a - x^{j}(t)||, \quad \forall 1 \le l \le k \right\}.$$
 (3.10)

(2.2) Center Update: For l = 1, 2, ..., k compute

$$x^{l}(t+1) = mean(C^{l}(t)) := \frac{1}{|C^{l}(t)|} \sum_{a \in C^{l}(t)} a.$$
(3.11)

(2.3) Stopping criteria: Halt if

$$\forall 1 \le l \le k \quad C^l(t+1) = C^l(t) \tag{3.12}$$

As in KPALM, KMEANS needs Assumption 1 for step (3.11) to be well defined. In order to prove the convergence of KMEANS to local minimum, we will need to following assumption.

$$= \left[\sum_{l=1}^{k} (4M)^2 ||x^l(t+1) - x^l(t)||^2 \right]^{\frac{1}{2}} = 4M ||x(t+1) - x(t)||,$$

this proves the desired result.

Proposition 3.2 (Subgradient lower bound for iterates gap property). Let $\{z(t)\}_{t\in\mathbb{N}} = \{(z(t), z(t))\}_{t\in\mathbb{N}}$ be the sequence generated by KPALM, then there exists $\rho_2 > 0$ and $\gamma(t+1) \in \partial \Psi(z(t+1))$ such

obtain

· Thon,

$$\|\gamma(t+1)\| \leq
ho_2 \|z(t+1)-z(t)\|, \quad orall t \in \mathbb{N}.$$

Proof. By the definition of Ψ (see (2.3)) we get

on of
$$\Psi$$
 (see (2.3)) we get
$$\partial \Psi = \nabla H + \partial G = \left((\nabla_{w^i} H_i + \partial_{w^i} \delta_{\Delta})_{i=1}, \nabla_x H \right).$$

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

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

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

$$\|\underline{\partial \Psi(z(t+1))}\| \leq \sum_{i=1}^{m} \|d^{i}(x(t+1)) + \partial_{w^{i}} \delta_{\Delta}(w^{i}(t+1))\|. \tag{3.6}$$

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

 $d^{i}(x(t)) + \alpha_{i}(t) \left(w^{i}(t+1) - w^{i}(t)\right) + u^{i}(t+1) = \mathbf{0}. \tag{3.7}$ Setting $\gamma(t+1) := \left(\left(d^{i}(x(t+1)) + u^{i}(t+1)\right)_{i=1}, \gamma_{n}, \mathbf{0}\right) \in \partial \Psi(z(t+1)), \text{ and plugging (3.7) into }$ By using (3.7) we

$$\begin{split} \|\gamma(t+1)\| &\leq \sum_{i=1}^{m} \|d^{i}(x(t+1)) - d^{i}(x(t)) - \alpha_{i}(t) \left(w^{i}(t+1) - w^{i}(t)\right)\| \\ &\leq \sum_{i=1}^{m} \|d^{i}(x(t+1)) - d^{i}(x(t))\| + \sum_{i=1}^{m} \alpha_{i}(t) \|w^{i}(t+1) - w^{i}(t)\| \\ &\leq \sum_{i=1}^{m} 4M \|x(t+1) - x(t)\| + m\overline{\alpha} \|x\| \|z(t+1) - z(t)\| \\ &\leq m \left(4M + \overline{\alpha}(x)\right) \|z(t+1) - z(t)\|, \end{split}$$

where the third inequality follows from Lemma 3.1.1, and $\overline{\alpha} = \max_{1 \le i \le m} \lambda$. Define

 $\rho_2 = \max m (4M + \overline{\alpha})$, due to Remark 1(iii) it follows that ρ_2 is bounded-from above, and the result follows.

$$\delta(t+1) = \left(\left(\delta^{1}(x(t+1)) - \delta^{1}(x(t)) - \omega_{1}(t) (w'(t+1) - w'(t)) \right)_{1=1,2,...,m}, 0 \right)$$

From Assumption I we have that $\beta(w(t)) = 2 \min_{1 \le l \le k} \left\{ \sum_{i=1}^m w_l^i(t) \right\} > \emptyset$, and from Lemma 3.0.2 it follows that the function $x \mapsto H(w(t), x)$ is strongly convex with parameter $\beta(w(t))$, hence it follows that

$$egin{aligned} H(w(t+1),x(t)) - H(w(t+1),x(t+1)) &\geq \\ &\geq \langle
abla_x H(w(t+1),x(t+1)),x(t) - x(t+1)
angle + rac{eta(w(t))}{2} \|x(t) - x(t+1)\|^2 \\ &= rac{eta(w(t))}{2} \|x(t+1) - x(t)\|^2, \end{aligned}$$

where the last equality follows from (3.2), since $\nabla_x H(w(t+1), x(t+1)) = 0$. Set $\rho_1 = \frac{1}{2} \min \{\underline{\alpha}(\lambda), \underline{\beta}(\lambda)\},$ combined with the previous inequalities, we have

$$\rho_{1}\|z(t+1)-z(t)\|^{2} = \rho_{1}\left(\|w(t+1)-w(t)\|^{2} + \|x(t+1)-x(t)\|^{2}\right) \leq \\
\leq \left[H(w(t),x(t)) - H(w(t+1),x(t))\right] + \left[H(w(t+1),x(t)) - H(w(t+1),x(t+1))\right] \\
= H(z(t)) - H(z(t+1)) = \Psi(z(t)) - \Psi(z(t+1)),$$

where the last equality follows from Remark 1(i)

Now

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

Lemma 3.1.1. Let $\{z(t)\}_{t\in\mathbb{N}} = \{(x(t), x(t))\}_{t\in\mathbb{N}}$ be the sequence generated by KPALM, then $\|d^i(x(t+1)-d^i(x(t))\| \leq 4M\|x(t+1)-x(t)\|, \quad \forall i=1,2,\ldots,m,\ t\in\mathbb{N},$

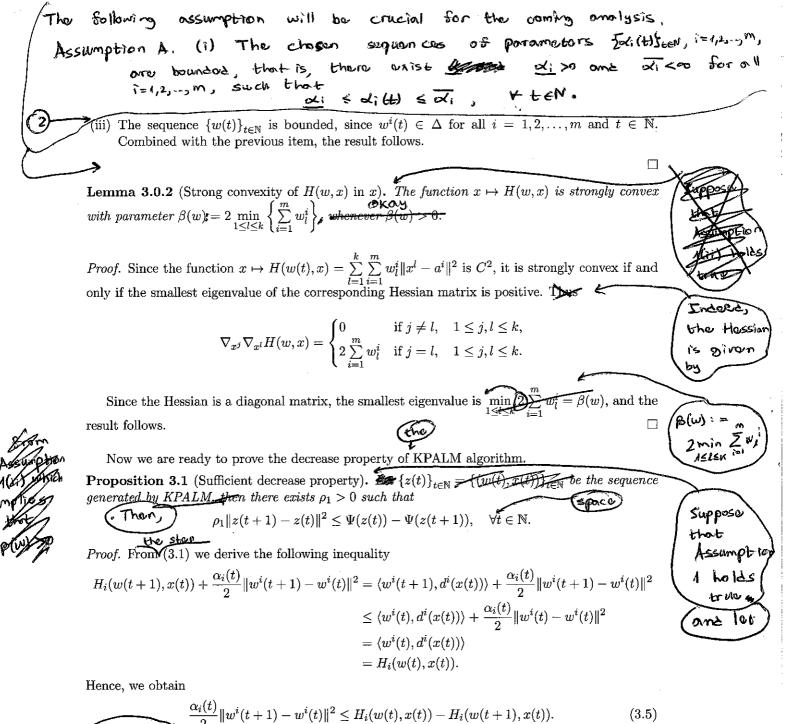
$$\|d^i(x(t+1)-d^i(x(t))\| \le 4M\|x(t+1)-x(t)\|, \quad \forall i=1,2,\ldots,m,\ t\in\mathbb{N}$$

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

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

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

the fact that G(w(t)) = 0



Denote
$$\underline{\alpha(t)} = \min_{1 \le i \le m} \{\alpha_i(t)\}$$
. Summing inequality (3.5) over $i = 1, 2, ..., m$ yields
$$\frac{\underline{\alpha(X)}}{2} \|w(t+1) - w(t)\|^2 = \frac{\underline{\alpha(X)}}{2} \sum_{i=1}^m \|w^i(t+1) - w^i(t)\|^2$$

$$\leq \sum_{i=1}^m \frac{\alpha_i(t)}{2} \|w^i(t+1) - w^i(t)\|^2$$

$$\leq \sum_{i=1}^m H_i(w(t), x(t)) - \sum_{i=1}^m H_i(w(t+1), x(t))$$

$$= H(w(t), x(t)) - H(w(t+1), x(t))$$

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

montioned above. More precisely, with respect to us we suggest b regularize the subproblem with proximal term as follows: w'(++1) = one min { < w', & (x(t)) > + xi(t) ||w'-w'(t)|| } i=1, m other hand, with respect to x we perfor exact minimizate

X(+++) = organizax H(V(+++), x). update. The explicit formulas, at step t, are given below. It is easy to c subproblems, with respect to will i=1,2,--, m, one X, com bછ simplified as $w^i(t+1) = P_{\Delta}\left(w^i(t) - \frac{d^i(x(t))}{\alpha_i(t)}\right), \quad i = 1, 2, \dots, m,$ Sollow: (3.3) $x^l(t+1) = rac{\sum_{i=1}^m w_l^i(t+1)a^i}{\sum_{i=1}^m w_l^i(t+1)}, \quad l = 1, 2, \dots, k$ (3.4)where P_{Δ} is the orthogonal projection onto the set Δ , Δ^m then G(w(t)) = 0 and therefore $\Psi(z(t)) = H(w(t), x(t)).$ (ii) For any choice of distance like function $d(\cdot,\cdot)$, the function $x\mapsto H(w,x)$ is separable in x^t for $all 1-1,2,\ldots,k$. Thus, regardless the choice of distance-like function $d(\cdot,\cdot)$ we bogin our analysis of the KPALIN algorithm cluster assignment step we can bound the choice of $\alpha_i(t)$ out of some interval with the Enllowing **Lemma 3.0.1** (Boundedness of KPALM sequence). Let $\{z(t)\}_{t\in\mathbb{N}}$ $\underbrace{(z(t))}_{t\in\mathbb{N}}$ be the seboundedness quence generated by KPALM. Then, the following statements hold true. the convex beabacta az Logallo & Ab, (i) For all $l=1,2,\ldots,k$, the sequence $\{x^l(t)\}_{t\in\mathbb{N}}$ is contained in Conv(A), where Conv(A) is ltho gunorobol and themselve soquence. the sequence $\{x^l(t)\}_{t\in\mathbb{N}}$ is bounded by $M=\max_{1\leq i\leq m}\|a^i\|$. bounded by For simplicity M : = (iii) The sequence $\{z(t)\}_{t\in\mathbb{N}}$ is bounded in $\mathbb{R}^{km}\times\mathbb{R}^{nk}$. from mon max Na'H on, we (i) Set $\lambda_i = \frac{w_i^i(t)}{\sum_{j=1}^m w_j^j(t)}$, $i = 1, 2, \dots, m$, then $\lambda_i \geq 0$ and $\sum_{j=1}^m \lambda_i = 1$. From (3.4) we have HEISM Lanota Z(t):= $x^{l}(t) = \frac{\sum_{i=1}^{m} w_{l}^{i}(t)a^{i}}{\sum_{i=1}^{m} w_{l}^{i}(t)} = \sum_{i=1}^{m} \left(\frac{w_{l}^{i}(t)}{\sum_{i=1}^{m} w_{l}^{j}(t)}\right)a^{i} = \sum_{i=1}^{m} \lambda_{i}a^{i} \in Conv(\mathcal{A}).$ $(\omega(t), X(t))$ Abb numbor (num) ten. M to this Hence $x^l(t)$ is in the convex hull of \mathcal{A} , for all $l=1,2,\ldots,k$ and $t\in\mathbb{N}$. Taking the norm of $x^l(t)$ yields again from (3.4) that OKEN THE SECOND Sormula $\|x^l(t)\| = \left\|\sum_{i=1}^m \left(\frac{w_l^i(t)}{\sum_{j=1}^m w_l^j(t)}\right) a^i\right\| \leq \sum_{i=1}^m \left(\frac{w_l^i(t)}{\sum_{i=1}^m w_l^j(t)}\right) \|a^i\| \leq \sum_{i=1}^m \lambda_i \max_{1 \leq i \leq m} \|a^i\| = M.$ $\|\mathbf{x}^{p}(\mathbf{t})\| = \|\mathbf{Z}\mathbf{\lambda}_{1}\mathbf{a}^{\dagger}\| \leq 3\mathbf{Z}\mathbf{\lambda}_{1}\|\mathbf{a}^{\dagger}\|$

4 Clustering via Alternation with Weiszfeld Step

In this section we tackle the clustering problem with distance-like function being the Euclidean norm in \mathbb{R}^n , namely

$$\min_{x^1, x^2, \dots, x^k \in \mathbb{R}^n} \left\{ \sum_{i=1}^m \min_{1 \le l \le k} \|x^l - a^i\| \right\}. \tag{4.1}$$

We are about to develop an algorithm that is based on PALM theory to treat this problem. However, first we need to discuss the Fermat-Weber problem that bears close relation with the algorithm that we will present later, and develop some useful tools.

4.1 The Smoothed Fermat-Weber Problem

Solving the smoothed Fermat-Weber plays a significant role in the algorithm that addresses the clustering problem with Euclidean norm distance-like function. The Fermat-Weber problem is formulated as follows

$$\min_{x \in \mathbb{R}^n} \left\{ f(x) := \sum_{i=1}^m w_i ||x - a^i|| \right\},\tag{4.2}$$

where $w_i > 0$, i = 1, 2, ..., m, are given positive weights and $\mathcal{A} = \{a^1, a^2, ..., a^m\} \subset \mathbb{R}^n$ are given vectors. As shown in [BS2015] this problem can be solved via the consecutive appliance of the operator $T : \mathbb{R}^n \to \mathbb{R}^n$ defined by

$$T(x) = \frac{1}{\sum_{i=1}^{m} \frac{w_i}{\|x - a^i\|}} \sum_{i=1}^{m} \frac{w_i a^i}{\|x - a^i\|}.$$

It is easily noticed that f(x) is not differentiable over \mathcal{A} . For our purposes we are interested in the smoothed Fermat-Weber problem, that can be formulated in the following manner

$$\min_{x \in \mathbb{R}^n} \left\{ f_{\epsilon}(x) := \sum_{i=1}^m w_i \left(\|x - a^i\|^2 + \epsilon^2 \right)^{1/2} \right\},\tag{4.3}$$

with $\epsilon > 0$ being some small perturbation constant. Next we introduce the operator $T_{\epsilon} : \mathbb{R}^n \to \mathbb{R}^n$ defined by

$$T_{\epsilon}(x) = \frac{1}{\sum_{i=1}^{m} \frac{w_i}{(\|x - a^i\|^2 + \epsilon^2)^{1/2}}} \sum_{i=1}^{m} \frac{w_i a^i}{(\|x - a^i\|^2 + \epsilon^2)^{1/2}}.$$

This version of the operator together with its properties that are to be discussed below are the cornerstone to prove the properties needed by PALM, and in turn to show the convergence of the sequence generated by the algorithm proposed to tackle the smooth version of the clustering problem presented later on. In order to prove some properties of T_{ϵ} , which are the same as the properties of T described in [BS2015], we also will need an auxiliary function $h_{\epsilon}: \mathbb{R}^{n} \times \mathbb{R}^{n} \to \mathbb{R}$ given by

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

Another useful function $L_{\epsilon}: \mathbb{R}^n \to \mathbb{R}$ that serves somewhat like Lipschitz function for the gradient of f_{ϵ} is defined by

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

It is easy to verify the following equality

$$T_{\epsilon}(x) = x - \frac{1}{L_{\epsilon}(x)} \nabla f_{\epsilon}(x), \quad \forall x \in \mathbb{R}^{n}.$$
 (4.4)

Lemma 4.0.1 (Properties of the auxiliary function h_{ϵ}). The following properties of h_{ϵ} hold.

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

$$h_{\epsilon}(y,y) = f_{\epsilon}(y).$$

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

$$h_{\epsilon}(x,y) \ge 2f_{\epsilon}(x) - f_{\epsilon}(y).$$

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

$$T_{\epsilon}(y) = \arg\min_{x \in \mathbb{R}^n} h_{\epsilon}(x, y).$$

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

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

Proof. (i) Follows by substituting x = y in h(x, y).

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

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

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

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

Multiplying the last inequality by w_i and summing over i = 1, 2, ..., m, the results follows.

(iii) The function $x \mapsto h_{\epsilon}(x,y)$ is strongly convex and its unique minimizer is determined by the optimality equation

$$\nabla_x h_{\epsilon}(x,y) = \sum_{i=1}^m \frac{2w_i (x - a^i)}{(\|y - a^i\|^2 + \epsilon^2)^{1/2}} = 0.$$

Simple algebraic manipulation leads to the relation

$$x = T_{\epsilon}(y),$$

and the desired results follows.

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

The following proofs are based on the properties of the auxiliary function h_{ϵ} , and they are similar to the proofs in [BS2015], hence we will just state them here. Lemma 4.0.5 does not appear in that paper, and its proof is given here.

Lemma 4.0.2 (Monotonicity property of T_{ϵ} , similar to (BS2015, Lemma 3.2, page 7)). For every $y \in \mathbb{R}^n$ we have

$$f_{\epsilon}(T_{\epsilon}(y)) \leq f_{\epsilon}(y).$$

Lemma 4.0.3 (Decent lemma for function f_{ϵ} , similar to (BS2015, Lemma 5.1, page 10)). For every $y \in \mathbb{R}^n$ we have

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

Lemma 4.0.4 (Similar to (BS2015, Lemma 5.2, page 12)). For every $x, y \in \mathbb{R}^n$ we have

$$f_{\epsilon}(T_{\epsilon}(y)) - f_{\epsilon}(x) \leq \frac{L_{\epsilon}(y)}{2} \left(\|y - x\|^2 - \|T_{\epsilon}(y) - x\|^2 \right).$$

Lemma 4.0.5. For all $y^0, y \in \mathbb{R}^n$ the following statement holds true

$$\|\nabla f_{\epsilon}(y) - \nabla f_{\epsilon}(y^{0})\| \leq \frac{2L_{\epsilon}(y^{0})L_{\epsilon}(y)}{L_{\epsilon}(y^{0}) + L_{\epsilon}(y)}\|y^{0} - y\|.$$

Proof. Let $\mathcal{P} \in \mathbb{R}^n$ be a fixed vector. Define the following two functions

and $\widetilde{h_{\epsilon}}(x,y) = h_{\epsilon}(x,y) - \left\langle \nabla f_{\epsilon}(y^{0}), x \right\rangle.$ It is clear that $x \mapsto \widetilde{h_{\epsilon}}(x,y)$ is sail quadratic function with associated matrix $L_{\epsilon}(y)\mathbf{I}$. Therefore, from 4.0.1(i) we can write

$$\widetilde{h_{\epsilon}}(x,y) = \widetilde{h_{\epsilon}}(y,y) + \left\langle \nabla_{x}\widetilde{h_{\epsilon}}(y,y), x - y \right\rangle + L_{\epsilon}(y)\|x - y\|^{2}$$

$$= \widetilde{f_{\epsilon}}(y) + \left\langle 2\nabla f_{\epsilon}(y) - \nabla f_{\epsilon}(y^{0}), x - y \right\rangle + L_{\epsilon}(y)\|x - y\|^{2}.$$
(4.5)

(114-bill 2+ e2)

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

$$\widetilde{h_{\epsilon}}(x,y) = h_{\epsilon}(x,y) - \langle \nabla f_{\epsilon}(y^{0}), x \rangle \ge 2f_{\epsilon}(x) - f_{\epsilon}(y) - \langle \nabla f_{\epsilon}(y^{0}), x \rangle$$

$$= 2\widetilde{f_{\epsilon}}(x) - \widetilde{f_{\epsilon}}(y) + \langle \nabla f_{\epsilon}(y^{0}), x - y \rangle,$$
(4.6)

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

$$2\widetilde{f}_{\epsilon}(x) \leq 2\widetilde{f}_{\epsilon}(y) + 2\left\langle \nabla f_{\epsilon}(y) - \nabla f_{\epsilon}(y^{0}), x - y \right\rangle + L_{\epsilon}(y)\|x - y\|^{2}$$
$$= 2\widetilde{f}_{\epsilon}(y) + 2\left\langle \nabla \widetilde{f}_{\epsilon}(y), x - y \right\rangle + L_{\epsilon}(y)\|x - y\|^{2}.$$

Dividing the last inequality by 2 leads to

$$\widetilde{f}_{\epsilon}(x) \le \widetilde{f}_{\epsilon}(y) + \left\langle \nabla \widetilde{f}_{\epsilon}(y), x - y \right\rangle + \frac{L_{\epsilon}(y)}{2} \|x - y\|^2.$$
 (4.7)

It is clear that the optimal point of \widetilde{f}_{ϵ} is y^0 since $\nabla \widetilde{f}_{\epsilon}(y^0) = 0$, therefore using (4.7) with $x = y - \frac{1}{L_{\epsilon}(y)} \nabla \widetilde{f}_{\epsilon}(y)$ yields

$$\begin{split} \widetilde{f}_{\epsilon}(y^{0}) &\leq \widetilde{f}_{\epsilon}\left(y - \frac{1}{L_{\epsilon}(y)}\nabla\widetilde{f}_{\epsilon}(y)\right) \leq \widetilde{f}_{\epsilon}(y) + \left\langle\nabla\widetilde{f}_{\epsilon}(y), -\frac{1}{L_{\epsilon}(y)}\nabla\widetilde{f}_{\epsilon}(y)\right\rangle + \frac{L_{\epsilon}(y)}{2}\left\|\frac{1}{L_{\epsilon}(y)}\nabla\widetilde{f}_{\epsilon}(y)\right\|^{2} \\ &= \widetilde{f}_{\epsilon}(y) - \frac{1}{2L_{\epsilon}(y)}\left\|\nabla\widetilde{f}_{\epsilon}(y)\right\|^{2}. \end{split}$$

Thus, using the definition of \widetilde{f}_{ϵ} and the fact that $\nabla \widetilde{f}_{\epsilon}(y) = \nabla f_{\epsilon}(y) - \nabla f_{\epsilon}(y^{0})$, yields that

$$f_{\epsilon}(y^0) \le f_{\epsilon}(y) + \langle \nabla f_{\epsilon}(y^0), y^0 - y \rangle - \frac{1}{2L_{\epsilon}(y)} \|\nabla f_{\epsilon}(y) - \nabla f_{\epsilon}(y^0)\|^2.$$

Now, following the same arguments we can show that

$$f_{\epsilon}(y) \leq f_{\epsilon}(y^0) + \left\langle \nabla f_{\epsilon}(y), y - y^0 \right\rangle - \frac{1}{2L_{\epsilon}(y^0)} \|\nabla f_{\epsilon}(y^0) - \nabla f_{\epsilon}(y)\|^2,$$

and combining last two inequalities yields that

$$\left(\frac{1}{2L_{\epsilon}(y^0)} + \frac{1}{2L_{\epsilon}(y)}\right) \|\nabla f_{\epsilon}(y) - \nabla f_{\epsilon}(y^0)\|^2 \leq \left\langle \nabla f_{\epsilon}(y^0) - \nabla f_{\epsilon}(y), y^0 - y \right\rangle,$$

that is,

$$\|\nabla f_{\epsilon}(y) - \nabla f_{\epsilon}(y^0)\| \leq \frac{2L_{\epsilon}(y^0)L_{\epsilon}(y)}{L_{\epsilon}(y^0) + L_{\epsilon}(y)}\|y^0 - y\|,$$

for all $y^0, y \in \mathbb{R}^n$.

Stort

4.2 Algorithm to the Smoothed Clustering Problem

In the previous section we showed that (4.1) has 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

$$H(w,x) = \sum_{i=1}^m \left\langle w^i, d^i(x) \right
angle = \sum_{i=1}^m \sum_{l=1}^k w^i_l \|x^l - a^i\|,$$

and

$$G(w) = \sum_{i=1}^m \delta_{\Delta}(w^i).$$

which is not the case now However, in order to be able to use the theory of PALM, we need the coupled function H(w, x) to be smooth, and in our case it is not. Therefore, it leads us to the following smoothed form of the clustering problem

E>O,

montioned

$$\min\left\{\Psi_{\epsilon}(z) := H_{\epsilon}(w, x) + G(w) \mid z := (w, x) \in \mathbb{R}^{km} \times \mathbb{R}^{nk}\right\},\tag{4.8}$$

where

$$H_{\epsilon}(w,x) = \sum_{i=1}^m \left\langle w^i, d^i_{\epsilon}(x)
ight
angle = \sum_{i=1}^m \sum_{l=1}^k w^i_l \left(\|x^l - a^i\|^2 + \epsilon^2
ight)^{1/2},$$

12

Please use vareption one not lapsilan!

 $\lambda_{\epsilon}'(x) = ($

ER" =1,2, -- , m.

with $d_{\epsilon}^{i}(x) = \left(\left(\|x^{1} - a^{i}\|^{2} + \epsilon^{2} \right)^{1/2}, \left(\|x^{2} - a^{i}\|^{2} + \epsilon^{2} \right)^{1/2}, \dots, \left(\|x^{k} - a^{i}\|^{2} + \epsilon^{2} \right)^{1/2} \right) \in \mathbb{R}^{k}$, for $i = 1, 2, \dots, m$. Note that $\Psi_{\epsilon}(z)$ is a perturbed form of $\Psi(z)$ for some small $\epsilon > 0$,

40(z)=4(z)

Next we extend the notations of the previous subsection, so that the functions and operators defined there are to be dependent on the weights w. For each $1 \le l \le k$, denote $w_l = (w_l^1, w_l^2, \dots, w_l^m)$ $\in \mathbb{R}^m_+$ and define

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

and

$$-\frac{1}{T_{\epsilon}^{w_l}(x^l)} - \frac{1}{L_{\epsilon}^{w_l}(x^l)} \sum_{i=1}^{m} \frac{w_l^i a^i}{(\|x^l - a^i\|^2 + \epsilon^2)^{1/2}}.$$

we define $H^{w_l}:\mathbb{R}^n\to\mathbb{R}$ as follows

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

thus we have

$$H_c(w,x) = \sum_{l=1}^k H^{w_l}(x^l).$$

Now we present our algorithm for solving broblem (4.8), we call it ε -KPALM. The algorithm alternates between cluster assignment step, similar to that as in KPALM, and centers update step that is based on a Toperator.
Cortain one gradient stop

- (1) Initialization: Set t=0, and pick random vectors $(w(0), x(0)) \in \Delta^m \times \mathbb{R}^{nk}$
- (2.1) Cluster Assignment: Pake by $\alpha_i(t) > 0$ and for each i = 1, 2, ..., m compute

$$w^{i}(t+1) = \arg\min_{\boldsymbol{w}^{i} \in \Delta} \left\{ \langle \boldsymbol{w}^{i}, \boldsymbol{d}_{\epsilon}^{i}(\boldsymbol{x}(t)) \rangle + \frac{\alpha_{i}(t)}{2} \| \boldsymbol{w}^{i} - \boldsymbol{w}^{i}(t) \|^{2} \right\}$$

$$= P_{\Delta} \left(\boldsymbol{w}^{i}(t) - \frac{\boldsymbol{d}_{\epsilon}^{i}(\boldsymbol{x}(t))}{\alpha_{i}(t)} \right). \tag{4.9}$$

(2.2) Center \mathcal{U} pdate: For each l = 1, 2, ..., k compute

$$x^{l}(t+1) - \frac{1}{T_{\epsilon}^{w_{l}(t+1)}(x^{l}(t))} \times \frac{1}{L_{\epsilon}(w(t+1), x(t)) \cdot 10)}$$

(i) Assumption I is still valid, hence the center update step in (4.10) is well defined.

(ii) It is easy to verify that for all $1 \le l \le k$ the following equations hold true:

$$\nabla H_{\epsilon}^{w_{l}}(x^{l}) = \sum_{i=1}^{m} w_{l}^{i} \frac{x^{l} - a^{i}}{(\|x^{l} - a^{i}\|^{2} + \epsilon^{2})^{1/2}}, \quad \forall x^{l} \in \mathbb{R}^{n}, \tag{4.11}$$

$$T_{\epsilon}^{w_l}(\underline{x^l}) = \underline{x^l} - \frac{1}{L_{\epsilon}^{w_l}(\underline{x^l})} \nabla H_{\epsilon}^{w_l}(\underline{x^l}), \quad \forall \, \underline{x^l} \in \mathbb{R}^n.$$

$$\tag{4.12}$$

As in KPALM case, the sequence that is generated by ε -KPALM is contained within the convex-hull of A. Indeed,

$$x^l(t+1) \equiv T^{w_l(t+1)}_{\epsilon}(x^l(t)) = \underbrace{\sum_{i=1}^m \frac{w_l^i(t+1)a^i}{\left(\|x^l(t)-a^i\|^2+\epsilon^2\right)^{1/2}}}_{i=1} \underbrace{\sum_{i=1}^m \frac{w_l^i(t+1)}{\left(\|x^l(t)-a^i\|^2+\epsilon^2\right)^{1/2}}}_{i=1} \underbrace{\sum_{i=1}^m \frac{w_l^i(t+1)}{\left(\|x^l(t)-a^i\|^2+\epsilon^2\right)^{1/2}}}_{j=1} \underbrace{\sum_{i=1}^m \frac{w_l^i(t+1)}{\left(\|x^l(t)-a^i\|^2+\epsilon^2\right)^{1/2}}}_{j=1} \underbrace{a^i \in Conv(\mathcal{A})}_{i=1},$$

hence the sequence generated by e-KPALM is bounded as well.

Now we are finally ready to prove the properties needed by PALM, and deduce that the sequence that is generated by ε -KPALM converge to critical point of Ψ_{ε} .

Proposition 4.1 (Sufficient decrease property). Let $\{z(t)\}_{t\in\mathbb{N}} = \{(w(t), x(t))\}_{t\in\mathbb{N}}$ be the sequence generated by ε -KPALM, then there exists $\rho_1 > 0$ such that

$$ho_1 \|z(t+1) - z(t)\|^2 \le \Psi_{\epsilon}(z(t)) - \Psi_{\epsilon}(z(t+1)) \quad \forall t \in \mathbb{N}.$$

Proof. Similar steps to the ones in the proof of sufficient decrease property of KPALM lead to

$$\frac{\underline{\alpha}(t)}{2} \| w(t+1) - w(t) \|^2 \le H_{\epsilon}(w(t), x(t)) - H_{\epsilon}(w(t+1), x(t)), \tag{4.13}$$

where $\underline{\alpha}(t) = \min_{1 \le i \le m} \{\alpha_i(t)\}.$

Applying Lemma 4.0.4 with respect to $H_{\epsilon}^{w_l(t+1)}(\cdot)$ yields

$$H^{w_l(t+1)}_{\epsilon}(x^l(t+1)) - H^{w_l(t+1)}_{\epsilon}(x^l) \leq \frac{L^{w_l(t+1)}_{\epsilon}(x^l(t))}{2} \left(\|x^l(t) - x^l\|^2 - \|x^l(t+1) - x^l\|^2 \right), \quad \dot{\forall} x^l \in \mathbb{R}^n,$$

for all $l=1,2,\ldots,k$. Setting $x^l=x^l(t)$ and rearranging yields

$$\frac{L_{\epsilon}^{w_l(t+1)}(x^l(t))}{2} \|x^l(t+1) - x^l(t)\|^2 \le H_{\epsilon}^{w_l(t+1)}(x^l(t)) - H_{\epsilon}^{w_l(t+1)}(x^l(t+1)), \quad \forall \, 1 \le l \le k. \quad (4.14)$$

Denote $\underline{L}(t) = \min_{1 \leq l \leq k} \left\{ L_{\epsilon}^{w_l(t+1)}(x^l(t)) \right\}$. Summing (4.14) over $l = 1, 2, \dots, k$ leads to

$$\frac{\underline{L}(t)}{2} \|x(t+1) - x(t)\|^{2} = \frac{\underline{L}(t)}{2} \sum_{l=1}^{k} \|x^{l}(t+1) - x^{l}(t)\|^{2}$$

$$\leq \sum_{l=1}^{k} \frac{L_{\epsilon}^{w_{l}(t+1)}(x^{l}(t))}{2} \|x^{l}(t+1) - x^{l}(t)\|^{2}$$

$$\leq \sum_{l=1}^{k} \left(H_{\epsilon}^{w_{l}(t+1)}(x^{l}(t)) - H_{\epsilon}^{w_{l}(t+1)}(x^{l}(t+1))\right)$$

$$= H_{\epsilon}(w(t+1), x(t)) - H_{\epsilon}(w(t+1), x(t+1)).$$
(4.15)

Set $\rho_1 = \frac{1}{2} \min_{t \in \mathbb{N}} \{\underline{\alpha}(t), \underline{L}(t)\}$, and note that since $x^l(t) \in Conv(\mathcal{A})$ for all $1 \leq l \leq k$, then

$$L_{\epsilon}^{w_l(t+1)}(x^l(t)) = \sum_{i=1}^m \frac{w_l^i(t+1)}{\left(\|x^l(t) - a^i\|^2 + \epsilon^2\right)^{1/2}} \ge \frac{\sum_{i=1}^m w_l^i(t+1)}{\left(d_{\mathcal{A}}^2 + \epsilon^2\right)^{1/2}}$$

where $d_{\mathcal{A}} = diam(Conv(\mathcal{A}))$, hence together with Remark 1(iii) and Assumption 1 assures that $\rho_1 > 0$. Combining (4.13) and (4.15) yields

$$\rho_1 \|z(t+1) - z(t)\|^2 = \rho_1 \left(\|w(t+1) - w(t)\|^2 + \|x(t+1) - x(t)\|^2 \right) \le \\
\le \left[H_{\epsilon}(w(t), x(t)) - H_{\epsilon}(w(t+1), x(t)) \right] + \left[H_{\epsilon}(w(t+1), x(t)) - H_{\epsilon}(w(t+1), x(t+1)) \right] \\
= H_{\epsilon}(z(t)) - H_{\epsilon}(z(t+1)) = \Psi_{\epsilon}(z(t)) - \Psi_{\epsilon}(z(t+1)),$$

which proves the desired result.

The next lemma will be useful in proving the subgradient lower bounds for iterates gap property of the sequence generated by ε -KPALM.

Lemma 4.1.1. For any $x, y \in \mathbb{R}^{nk}$ such that $x^l, y^l \in Conv(\mathcal{A})$ for all $1 \leq l \leq k$ the following inequality holds

$$\|d^i_\epsilon(x)-d^i_\epsilon(y)\| \leq rac{d_{\mathcal{A}}}{\epsilon}\|x-y\|, \quad orall\, i=1,2,\ldots,m,$$

with $d_{\mathcal{A}} = diam(Conv(\mathcal{A}))$.

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

$$\frac{\psi(a) - \psi(b)}{a - b} = \psi'(c) = \frac{1}{2\sqrt{c + \epsilon^2}} \le \frac{1}{2\epsilon}$$

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

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

Hence.

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

as asserted.

Lemma 4.1.2 (Upper bound of the sequence $\{\overline{L}(x(t))\}_{t\in\mathbb{N}}$). Let $\{z(t)\}_{t\in\mathbb{N}}=\{(w(t),x(t))\}_{t\in\mathbb{N}}$ be the sequence generated by ε -KPALM, then for any $t\in\mathbb{N}$ we have

$$\overline{L}(x(t)) = \max_{1 \le l \le k} \left\{ L_{\epsilon}^{w_l(t+1)}(x^l(t)) + \frac{2L_{\epsilon}^{w_l(t+1)}(x^l(t))L_{\epsilon}^{w_l(t+1)}(x^l(t+1))}{L_{\epsilon}^{w_l(t+1)}(x^l(t)) + L_{\epsilon}^{w_l(t+1)}(x^l(t+1))} \right\} \le \frac{2m}{\epsilon}.$$

Proof. For any $w_l \in [0,1]^m$ and $x^l \in \mathbb{R}^n$ we have

$$L_{\epsilon}^{w_l}(x^l) = \sum_{i=1}^m \frac{w_l^i}{\left(\|x^l - a^i\|^2 + \epsilon^2\right)^{1/2}} \le \sum_{i=1}^m \frac{1}{\epsilon} = \frac{m}{\epsilon}.$$

Therefore,

$$\overline{L}(x(t)) = \max_{1 \le l \le k} \left\{ L_{\epsilon}^{w_l(t+1)}(x^l(t)) + \frac{2}{L_{\epsilon}^{w_l(t+1)}(x^l(t))} + \frac{1}{L_{\epsilon}^{w_l(t+1)}(x^l(t+1))} \right\} \le \frac{m}{\epsilon} + \frac{2}{\frac{2\epsilon}{m}} = \frac{2m}{\epsilon},$$

this proves the desired result.

Proposition 4.2 (Subgradient lower bound for iterates gap property). Let $\{z(t)\}_{t\in\mathbb{N}} = \{(w(t), x(t))\}_{t\in\mathbb{N}}$ be the sequence generated by ε -KPALM, then there exists $\rho_2 > 0$ and $\gamma(t+1) \in \partial \Psi_{\epsilon}(z(t+1))$ such that

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

Proof. Repeating the steps of the proof in the case of KPALM yields that

$$\gamma(t+1) := \left(\left(d_{\epsilon}^{i}(x(t+1)) + u^{i}(t+1) \right)_{i=1,\dots,m}, \nabla_{x} H_{\epsilon}(w(t+1), x(t+1)) \right) \in \partial \Psi_{\epsilon}(z(t+1)), \quad (4.16)$$

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

$$d_{\epsilon}^{i}(x(t)) + \alpha_{i}(t) \left(w^{i}(t+1) - w^{i}(t) \right) + u^{i}(t+1) = \mathbf{0}. \tag{4.17}$$

Plugging (4.17) into (4.16), and taking norm yields

$$\begin{aligned} \|\gamma(t+1)\| &\leq \sum_{i=1}^{m} \|d_{\epsilon}^{i}(x(t+1)) - d_{\epsilon}^{i}(x(t)) - \alpha_{i}(t) \left(w^{i}(t+1) - w^{i}(t)\right)\| + \|\nabla_{x}H_{\epsilon}(w(t+1), x(t+1))\| \\ &\leq \sum_{i=1}^{m} \|d_{\epsilon}^{i}(x(t+1)) - d_{\epsilon}^{i}(x(t))\| + \sum_{i=1}^{m} \alpha_{i}(t)\|w^{i}(t+1) - w^{i}(t)\| + \|\nabla_{x}H_{\epsilon}(w(t+1), x(t+1))\| \\ &\leq \frac{md_{\mathcal{A}}}{\epsilon} \|x(t+1) - x(t)\| + m\overline{\alpha}(t)\|w(t+1) - w(t)\| + \|\nabla_{x}H_{\epsilon}(w(t+1), x(t+1))\|, \end{aligned}$$

where the last inequality follows from Lemma 4.1.1 and the fact that $\overline{\alpha}(t) = \max_{1 \le i \le m} \alpha_i(t)$.

Next we bound $\|\nabla_x H_{\epsilon}(w(t+1), x(t+1))\| \le c\|x(t+1) - x(t)\|$, for some constant c > 0. Indeed, we have

$$\|\nabla_{x} H_{\epsilon}(w(t+1), x(t+1))\| \leq \sum_{l=1}^{k} \|\nabla H_{\epsilon}^{w_{l}(t+1)}(x^{l}(t+1))\|$$

$$\leq \sum_{l=1}^{k} \|\nabla H_{\epsilon}^{w_{l}(t+1)}(x^{l}(t))\| + \sum_{l=1}^{k} \|\nabla H_{\epsilon}^{w_{l}(t+1)}(x^{l}(t+1)) - \nabla H_{\epsilon}^{w_{l}(t+1)}(x^{l}(t))\|.$$

$$(4.18)$$

From (4.10) and (4.12) we have

$$\nabla H_{\epsilon}^{w_l(t+1)}(x^l(t)) = L_{\epsilon}^{w_l(t+1)}(x^l(t)) \left(x^l(t+1) - x^l(t) \right), \quad \forall \ 1 \le l \le k,$$

applying Lemma 4.0.5 with respect to $H_{\epsilon}^{w_l(t+1)}(\cdot)$ and plugging into (4.18) yields

$$\|\nabla_x H(w(t+1), x(t+1))\| \le$$

$$\leq \sum_{l=1}^k \left(L_{\epsilon}^{w_l(t+1)}(x^l(t)) + \frac{2L_{\epsilon}^{w_l(t+1)}(x^l(t))L_{\epsilon}^{w_l(t+1)}(x^l(t+1))}{L_{\epsilon}^{w_l(t+1)}(x^l(t)) + L_{\epsilon}^{w_l(t+1)}(x^l(t+1))} \right) \|x^l(t+1) - x^l(t)\|.$$

Therefore, denote
$$\overline{L}(x(t)) = \max_{1 \le l \le k} \left\{ L_{\epsilon}^{w_l(t+1)}(x^l(t)) + \frac{2L_{\epsilon}^{w_l(t+1)}(x^l(t))L_{\epsilon}^{w_l(t+1)}(x^l(t))}{L_{\epsilon}^{w_l(t+1)}(x^l(t)) + L_{\epsilon}^{w_l(t+1)}(x^l(t))} \right\}$$
, and set $\rho_2 = m \left(\frac{dA}{L} + \overline{\alpha}(t) \right) + k\overline{L}(x(t))$, note that Lemma 4.1.2 together with Assumption 1 imply that

 $\rho_2 = m \left(\frac{d_A}{\epsilon} + \overline{\alpha}(t) \right) + k \overline{L}(x(t)),$ note that Lemma 4.1.2 together with Assumption 1 imply that ρ_2 is bounded from above, and the result follows.

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

Lemma 4.2.1 (Closeness of smooth). For any $(w,x) \in \Delta^m \times \mathbb{R}^{nk}$ and $\epsilon > 0$ the following inequalities hold true

$$H(w,x) \le H_{\epsilon}(w,x) \le H(w,x) + m\epsilon.$$

Proof. Applying the inequality

$$(a+b)^{\lambda} \le a^{\lambda} + b^{\lambda}, \quad \forall a, b \ge 0, \ \lambda \in (0,1],$$

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

$$(\|x^l - a^i\|^2 + \epsilon^2)^{1/2} \le \|x^l - a^i\| + \epsilon, \quad \forall \ 1 \le l \le k, \ 1 \le i \le m.$$

Together with the fact that

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

yields the following inequality

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

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

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

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

5 Clustering via ADMM Approach

Introducing some new variable into the problem leads to the following clustering problem notation

$$\begin{aligned} & \min_{x \in \mathbb{R}^{nk}} \min_{w \in \mathbb{R}^{km}} \left\{ \sum_{i=1}^{m} \sum_{l=1}^{k} w_{l}^{i} d(x^{l}, a^{i}) \mid w^{i} \in \Delta, i = 1, 2, \dots, m \right\} \\ & = \min_{x \in \mathbb{R}^{nk}, w \in \mathbb{R}^{km}, z \in \mathbb{R}^{km}} \left\{ \sum_{i=1}^{m} \sum_{l=1}^{k} w_{l}^{i} \mid v^{i} \in \Delta, & i = 1, 2, \dots, m, \\ z_{l}^{i} = d(x^{l}, a^{i}), & i = 1, 2, \dots, m, \\ l = 1, 2, \dots, k \end{array} \right\}.$$

The augmented Lagrangian that is associated with this problem is

$$L_{\rho}(w, x, z, y) = \sum_{i=1}^{m} \sum_{l=1}^{k} w_{l}^{i} z_{l}^{i} + \sum_{i=1}^{m} \sum_{l=1}^{k} y_{l}^{i} (z_{l}^{i} - d(x^{l}, a^{i})) + \frac{\rho}{2} \sum_{i=1}^{m} \sum_{l=1}^{k} \left(z_{l}^{i} - d(x^{l}, a^{i}) \right)^{2}.$$
 (5.1)

Thus the ADMM formulas for (4.2) are as follows

$$w(t+1) = \arg\min_{w \in \Delta^m} L_{\rho}(w, x(t), z(t), y(t)),$$

$$\Rightarrow w^{i}(t+1) = \arg\min_{w^{i} \in \Delta} \sum_{l=1}^{k} w_{l}^{i} z_{l}^{i}(t) = \arg\min_{w^{i} \in \Delta} \left\langle w^{i}, z^{i}(t) \right\rangle, \quad 1 \leq i \leq m,$$

$$x(t+1) = \arg\min_{w^{i} \in \Delta} L_{i}(w(t+1), x, z(t), y(t))$$

$$\Rightarrow x^{l}(t+1) = \arg\min_{x^{l} \in \mathbb{R}^{nk}} \sum_{i \neq 1}^{m} y_{l}^{i}(t)d(x^{l}, a^{i}) + \frac{\rho}{2} \sum_{i=1}^{m} \left(z_{l}^{i} - d(x^{l}, a^{i})\right)^{2}, \quad 1 \leq l \leq k,$$

$$z(t+1) = \arg\min_{z \in \mathbb{R}^{km}} L\left(w(t+1), x(t+1), z, y(t)\right),$$

$$\Rightarrow z^{i}(t+1) = \arg\min_{z^{i} \in \mathbb{R}^{km}} \left\langle w^{i}(t+1), z^{i} \right\rangle + \left\langle y^{i}(t), z^{i} \right\rangle + \frac{\rho}{2} \left\| z^{i} - \left(d(x^{l}(t+1), a^{i})_{l=1,\dots,k} \right)^{2} \right\|$$

$$= \left(d(x^{l}(t+1), a^{i})_{l=1,\dots,k} - \frac{1}{\rho} \left(w^{i}(t+1) + y^{i}(t) \right), \quad 1 \leq i \leq m,$$

$$y^{i}_{l}(t+1) = y^{i}_{l}(t) + \rho(z^{i}_{l}(t+1) - d(x^{l}(t+1), a^{i}), \quad 1 \leq i \leq m, \quad 1 \leq l \leq k.$$

KPALM

- 1 Initialization: $(w(0), x(0)) \in \Delta^m \times R^{nk}$
- (2) General stop (+=0,1,2,...):
 - (2.1) Cluster assignment: choose contain di(+)>0, i=1,2,..., m, and compute

$$\omega'(b+1) = P_{\Delta}\left(\omega'(t) - \frac{d'(x(t))}{\alpha_{i}(t)}\right).$$

$$\times^{2}(t+1) = \frac{\sum_{i=1}^{m} W_{2}^{i}(t+1)a^{i}}{\sum_{i=1}^{m} W_{2}^{i}(t+1)}.$$

bottomers and therefore there to subsugate new which converges to a surplied of the surplied o

(ii) For all tEN. Vara Torre that

It should be noted that Assumption I(i) is very milk since the potentials of d(t), $I \in i \in M$ and $b \in N$, can be chosen aribtrarily by the user and therefore it can be controlled such that the boundedness property holds true. Assumption I(ii) is essential since if it is not true then $W_2(t) = 0$ for all $I \le i \le M$, which means that the center X^{i} does not appropriate in the objective Sunction.

Now we would like to develope on algorithm which based on the methology of PALM to solve Problem (4.8). It is easy to see that with respect to w, the objective function the keeps on the some structure as 4 and therefore we apply the some stept. More precisely, for all i=1,2,..., m, we have

$$w'(t+1) = \underset{w \in \Delta}{\operatorname{argmin}} \left\{ < w', d_{e}^{i}(x(t)) > + \frac{d_{e}^{i}(t)}{2} ||w' - w'(t)||^{2} \right\}$$

$$= P_{\Delta} \left(w'(t) - \frac{d_{e}^{i}(x(t))}{d_{e}^{i}(t)} \right), \quad \forall t \in \mathbb{N},$$

where $x_i(t)$, i=1,2,..., m, is arbitrarly chosen. On the other hand, with respect to x we tackle the subproblem differently than kPALA. Here, we taken the follow exactly the idea of PALA

Whore was

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

(y) Similarly to the KPALM algorithm, the sequence generated by E-KPALM is also bounded, Since here we also have that

$$\times^{2}(t+1) = \times^{2}(t) - \frac{1}{L_{\epsilon}(\omega(t+1), \times^{2}(t))} \nabla_{x^{2}} H(\omega(t), \times(t))$$

use Valign

$$= \times^{2}(t) - \frac{1}{L_{\epsilon}(\omega(t+1), x^{2}(t))} \geq \omega_{2}(t^{\frac{1}{2}}+1) \circ \frac{x^{2}(t) - \alpha^{\frac{1}{2}}}{(\|x^{2}(t) - \alpha^{\frac{1}{2}}\|_{2}^{2} + \epsilon^{2})^{\frac{1}{2}}}$$

$$= \frac{1}{L_{\epsilon}(w(\pm n), x^{\ell}(\pm 1))} = \frac{w_{\epsilon}(\pm 1) o^{i}}{(\|x^{\ell}(\pm 1) - a^{i}\|^{2} + \epsilon^{2})^{1/2}} \in conv(\mathcal{A}),$$

Botore we will be able to prove the two proportions neaded for showing convergence of the sequence [z(t)] ten generated by E-KPALM, we will need several auxiliary results. For the simplicity of the expositions we define the Sollowing Sunction

for fixed W1, W23--, WM ER, me one bieR, i=1,2,.., m.

Lemma. The gradient of fe(.) is Lipschitz continuous, that is

$$\|\nabla f_{\varepsilon}(x) - \nabla f_{\varepsilon}(x)\| \leq \frac{2L_{\varepsilon}(x)L_{\varepsilon}(x)}{L_{\varepsilon}(x)+L_{\varepsilon}(x)} \|x-x\|, \forall x, x \in \mathbb{R}^n,$$

whore

$$L_{\varepsilon}(x) = \sum_{n=1}^{\infty} \frac{w_{i}}{(1 \times -\mathbf{k}^{2})^{2} + \varepsilon^{2}} \frac{1}{1/2}$$

Proof. ...

We also need the following function $h_{\varepsilon}(x,y) = - - .$

Lemma 4.0.1 ...