

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 , 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^l is minimized.

The clustering problem is given by

$$\min_{x^1, x^2, \dots, x^k \in \mathbb{R}^n} F(x^1, x^2, \dots, x^k) := \sum_{i=1}^m \min_{1 \leq l \leq k} d(x^l, a^i), \quad (1.1)$$

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

2 Problem Reformulation and Notations

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

$a = (a^1, a^2, \dots, a^m) \in \mathbb{R}^{nm}$, where $a^i \in \mathbb{R}^n, i = 1, 2, \dots, m$.

$w = (w^1, w^2, \dots, w^m) \in \mathbb{R}^{km}$, where $w^i \in \mathbb{R}^k, i = 1, 2, \dots, m$.

$x = (x^1, x^2, \dots, x^k) \in \mathbb{R}^{nk}$, where $x^l \in \mathbb{R}^n, l = 1, 2, \dots, k$.

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

$$\Delta = \left\{ u \in \mathbb{R}^k \mid \sum_{l=1}^k u_l = 1, u_l \geq 0, l = 1, 2, \dots, k \right\}.$$

Let $S \subseteq \mathbb{R}^n$. The indicator function of S is defined and denoted as follows $\delta_S(p) = \begin{cases} 0, & \text{if } p \in S, \\ \infty, & \text{if } p \notin S. \end{cases}$

Using the fact that $\min_{1 \leq l \leq k} u_l = \min \{\langle u, v \rangle \mid v \in \Delta\}$, and applying it over (1.1), gives a smooth reformulation of the clustering problem

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

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

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

Finally, introducing several more useful notations, for each $i = 1, 2, \dots, m$, we denote

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

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

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

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

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

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

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

At each step $t \in \mathbb{N}$, the KPALM algorithm alternates between cluster assignment and centers update. The explicit formulas, at step t , are given below

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

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

Assumption 1. We assume that none of the clusters get empty during this process, hence for all $1 \leq l \leq k$ and $t \in \mathbb{N}$ we have that $\sum_{i=1}^m w_l^i(t) > 0$.

Remark 1. (i) Since for all $t \in \mathbb{N}$ $w(t) \in \Delta^m$ then $\Psi(z(t)) = H(w(t), x(t)) + G(w(t)) = H(w(t), x(t))$.

(ii) For any choice of distance-like function $d(\cdot, \cdot)$, the function $x \mapsto H(w, x)$ is separable in x^l for all $l = 1, 2, \dots, k$. Thus, regardless the choice of distance-like function $d(\cdot, \cdot)$, the centers update step can be done in parallel over all centers, that is, $x^l(t+1) = \arg \min_{x^l \in \mathbb{R}^k} \left\{ \sum_{i=1}^m w_l^i d(x^l, a^i) \right\}$, $l = 1, 2, \dots, k$, and in the case of the squared Euclidean norm the result is as in (3.4).

Lemma 3.0.1 (Boundedness of KPALM sequence). Let $\{z(t)\}_{t \in \mathbb{N}} = \{(w(t), x(t))\}_{t \in \mathbb{N}}$ 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 $\text{Conv}(\mathcal{A})$, where $\text{Conv}(\mathcal{A})$ is the convex hull of \mathcal{A} .

(ii) For all $l = 1, 2, \dots, k$, the sequence $\{x^l(t)\}_{t \in \mathbb{N}}$ is bounded by $M = \max_{1 \leq i \leq m} \|a^i\|$.

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

Proof. (i) Set $\lambda_i = \frac{w_l^i(t)}{\sum_{j=1}^m w_l^j(t)}$, $i = 1, 2, \dots, m$, then $\lambda_i \geq 0$ and $\sum_{i=1}^m \lambda_i = 1$. From (3.4) we have

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

Hence $x^l(t)$ is in the convex hull of \mathcal{A} , for all $l = 1, 2, \dots, k$ and $t \in \mathbb{N}$.

(ii) Taking the norm of $x^l(t)$ yields again from (3.4) that

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

(iii) 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

Lemma 3.0.2 (Strong convexity of $H(w, x)$ in x). *The function $x \mapsto H(w, x)$ is strongly convex with parameter $\beta(w) = 2 \min_{1 \leq l \leq k} \left\{ \sum_{i=1}^m w_l^i \right\}$, if $\beta(w) > 0$.*

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 iff the smallest eigenvalue of the corresponding Hessian matrix is positive. Thus

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

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

Now we are ready to prove the decrease property of KPALM algorithm.

Proposition 3.1 (Sufficient decrease property). *Let $\{z(t)\}_{t \in \mathbb{N}} = \{w(t), x(t)\}_{t \in \mathbb{N}}$ be the sequence generated by KPALM, then there exists $\rho_1 > 0$ such that $\rho_1 \|z(t+1) - z(t)\|^2 \leq \Psi(z(t)) - \Psi(z(t+1))$ for all $t \in \mathbb{N}$.*

Proof. From (3.1) we derive 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

$$\frac{\alpha_i(t)}{2} \|w^i(t+1) - w^i(t)\|^2 \leq H_i(w(t), x(t)) - H_i(w(t+1), x(t)). \quad (3.5)$$

Denote $\alpha(t) = \min_{1 \leq i \leq m} \{\alpha_i(t)\}$. Summing inequality (3.5) over $i = 1, 2, \dots, m$ yields

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

From Assumption 1 we have that $\beta(w(t)) = 2 \min_{1 \leq l \leq k} \left\{ \sum_{i=1}^m w_l^i(t) \right\} > 0$, 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

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

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

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

where the last equality follows from Remark 1(i). \square

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 $\{z(t)\}_{t \in \mathbb{N}} = \{(w(t), x(t))\}_{t \in \mathbb{N}}$ be the sequence generated by KPALM, then*

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

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

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

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

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

this proves the desired results. \square

Proposition 3.2 (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 KPALM, then there exists $\rho_2 > 0$ and $\gamma(t+1) \in \partial\Psi(z(t+1))$ such that $\|\gamma(t+1)\| \leq \rho_2 \|z(t+1) - z(t)\|^2$ for all $t \in \mathbb{N}$.*

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

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

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

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

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

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

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

$$d^i(x(t)) + \alpha_i(t) (w^i(t+1) - w^i(t)) + u^i(t+1) = \mathbf{0}. \quad (3.7)$$

Setting $\gamma(t+1) := \left((d^i(x(t+1)) + u^i(t+1))_{i=1, \dots, m}, \mathbf{0} \right) \in \partial\Psi(z(t+1))$, and plugging (3.7) into (3.6) we have

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

where the third inequality follows from Lemma 3.1.1, and $\bar{\alpha}(t) = \max_{1 \leq i \leq m} \alpha_i(t)$. Define

$\rho_2 = m (4M + \bar{\alpha}(t))$ 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 $t = 0, 1, \dots$ generate a sequence $\{(v(t), y(t))\}_{t \in \mathbb{N}}$ as follows:
 - (2.1) Cluster Assignment: For $i = 1, 2, \dots, m$ compute

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

- (2.2) Center Update: For $l = 1, 2, \dots, k$ compute

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

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

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

meaning, both algorithms converge to the same result.

4 KMEANS with Weiszfeld step

Solving clustering problem with distance function $d(\cdot, \cdot) = \|u - v\|$

$$\min_{x^1, \dots, x^k \in \mathbb{R}^n} \left\{ \sum_{l=1}^k \sum_{i=1}^m w_l^i \|x^l - a^i\| \mid w^i \in \Delta, i = 1, 2, \dots, m \right\}$$

similarly to KMEANS it leads to the following alternation steps:

(1) Clusters Assignment: For $i = 1, 2, \dots, m$ compute

$$w^i(t+1) = \arg \min_{w^i \in \Delta} \{ \langle w^i, d^i(x(t)) \rangle \} \quad (4.1)$$

(2) Centers Update: For $l = 1, 2, \dots, k$ compute

$$x^l(t+1) = \arg \min_{x^l \in \mathbb{R}^n} \sum_{i=1}^m w_l^i(t+1) \|x^l - a^i\| \quad (4.2)$$

The center update step could be computed via Weiszfeld algorithm.