# Contents

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

#### 1 Introduction

(The K-

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

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

summation

On

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

- (i) Hierarchical versus partitioning clustering. In partitioning clustering the dataset is divided into clusters, whereas in hierarchical clustering each cluster may have subclusters, thus forming a tree which leaves are the single points of the dataset.
- (ii) Hard versus soft and fuzzy clustering In hard clustering each data point is assigned to single cluster, versus a soft clustering where each point may be assigned to more than one cluster, hence clusters may overlap. In fuzzy clustering for each point there is a distribution that describes the probability of a point to be part of a cluster.

(iii) Complete versus partial clustering. In complete clustering all points in the dataset are assigned to clusters, whereas in partial clustering some points may be intentionally skipped and are not being assigned to a cluster.

(cerbain)

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

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

This algorithm can be described as an aptimization algorithm (see precise dutails below) which minimizes the following

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

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

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

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

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

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

Soction 21 Sor mose

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

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

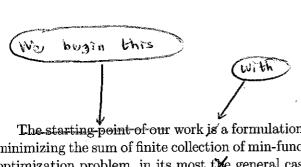
where  $p(a^i|x^l)$  is the probability of  $a^i$  given that it is generated by the Gaussian distribution with center  $x^l$  and  $p(x^l)$  is the prior probability of center  $x^l$ .

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

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

to chase

furthest point from the points chosen so far.



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

o This

xx) eR"x (inches  $x = (x^1, x^2)$ 

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

(1.1)the paragraph sho lite b norto

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

The lack of smoothness in this formulation can be overcome, yet the nonconvex nature of the clustering problem shall accompany the discussions throughout this work. Signifi-Hiw cant amount of studies have been made on convex models, even though in many cases the original optimization problem is nonconvex. To overcome the lack of convexity the common approach is usually achieved by relaxation of the original problem. Motivated by papers of Attouch et al. [1, 2] that established convergence of the sequences generated by the proximal Gauss-Seidel scheme in the general nonconvex and nonsmooth settings, and similar result

madizon. CONVOX relaxation  $\sigma \mathcal{F}$ 

for the proximal-forward-backward algorithm applied to the nonconvex and nonsmooth minimization of the sum of a nonsmooth function with a smooth one. This approach assumes that the objective function to be minimized satisfies the Kurdyka-Łojasiewicz (KL) property. The convergence results were further extended in the recent work by Bolte et al. [5], where the objective function is a function of finite blocks of variables.

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

was recently

algonthm)

Chartxoximals

mosom

in groat Eutail 5

Thon,

beabases

, wolfd

WO.

Solotokop

be discussed datasets in which the data points are spread relatively sparsely form their centers. Outline of the thesis. This work is organized as follows. In the following section we

In there is you

the format 1

write it

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

0/300 Km

to solve the approximated

model which combines iseas which

used in the

it Euclidean case

n classial 1500

In this thosis wo take ou 31680 nunt route and consider the problem in its original non convex form. Very recently complicated router become more rollowant and interesting thanks to fow papers which para the way for entities with nonconvex problems using sophisticated mathematical tools as will be explained later.

### 2 Problem Reformulation and Mathematical Tools

#### 2.1 Reformulation of the Clustering Problem

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

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

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

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

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

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

where

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

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

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

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

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

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

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

# 2.2 Introduction to the PALM Theory

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

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

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

$$\lim_{y \neq x} \inf_{y \to x} \frac{\sigma(y) - \sigma(x) - \langle u, y - x \rangle}{\|y - x\|} \ge 0.$$

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

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

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

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

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



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

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

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

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

the following inequality holds

$$\varphi'(\sigma(u) - \sigma(\overline{u}))dist(0, \partial \sigma(u)) \ge 1,$$

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

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

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

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

As montioned in [3], the PALM absorithm is nothing but alternating the classical proximal gradient (axa Proximal forward-Bockward) over

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

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

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

$$\widehat{\Psi}\left(x,y^{k}\right) = \left\langle x - x^{k}, \nabla_{x} H\left(x^{k},y^{k}\right)\right\rangle + \frac{e^{k^{k}}}{2} \left\|x - x^{k}\right\|^{2} + f(x), \ (c_{k} > 0),$$

and

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

Thus, PALM is alternating between the following-two-subproblems

com be summerized as Sollows

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

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

# 3 Clustering: The Squared Euclidean Norm Case

## 3.1 Clustering with PALM

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

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

Attouch & al. 11,21 established convergence — (continuo as in laguar account to prove convergence of contribution of [5] is the general methodogy to prove convergence of generic algorithm in the setting of nonconvex and mensmooth optimization problems

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

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

min

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

ωπροσυ

proposyb

Plansor replaces
this with
what I

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

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

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

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

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

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

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

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

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

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

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

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

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

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