

Robust k-means: a Theoretical Revisit

Authored by:

ALEXANDROS GEORGOGIANNIS

Abstract

Over the last years, many variations of the quadratic k-means clustering procedure have been proposed, all aiming to robustify the performance of the algorithm in the presence of outliers. In general terms, two main approaches have been developed: one based on penalized regularization methods, and one based on trimming functions. In this work, we present a theoretical analysis of the robustness and consistency properties of a variant of the classical quadratic k-means algorithm, the robust k-means, which borrows ideas from outlier detection in regression. We show that two outliers in a dataset are enough to breakdown this clustering procedure. However, if we focus on “well-structured” datasets, then robust k-means can recover the underlying cluster structure in spite of the outliers. Finally, we show that, with slight modifications, the most general non-asymptotic results for consistency of quadratic k-means remain valid for this robust variant.

1 Paper Body

Let $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}_+$ be a lower semi-continuous (lsc) and symmetric function with minimum value $\varphi(0)$. Given a set of points $X_n = \{x_1, \dots, x_n\} \subset \mathbb{R}^p$, consider the generalized k-means problem (GKM) [7]
$$R_n(c_1, \dots, c_k) = \min_{c_1, \dots, c_k} \sum_{i=1}^n \varphi(x_i - c_{l_i}) \quad (\text{GKM})$$
 subject to $c_l \in \mathbb{R}^p$, $l \in \{1, \dots, k\}$. Our aim is to find a set of k centers $\{c_1, \dots, c_k\}$ that minimize the clustering risk R_n . These centers define a partition of X_n into k clusters $A = \{A_1, \dots, A_k\}$, defined as “ $\# A_l = x \in X_n : l = \arg\min_{1 \leq j \leq k} \varphi(x - c_j)$ ”, (1)

where ties are broken randomly. Varying φ beyond the usual quadratic function ($\varphi(t) = t^2$) we expect to gain some robustness against the outliers [9]. When φ is upper bounded by φ^* , the clusters are defined as follows. For $l \in \{1, \dots, k\}$, let “ $\# A_l = x \in X_n : l = \arg\min_{1 \leq j \leq k} \varphi(x - c_j)$ ” and “ $\# A_{k+1} = x \in X_n : \min_{1 \leq j \leq k} \varphi(x - c_j) \geq \varphi^*$ ”, (2) and define the extra cluster

$$\# A_{k+1} = x \in X_n : \min_{1 \leq j \leq k} \varphi(x - c_j) \geq \varphi^* \quad (3)$$

This extra cluster contains points whose distance from their closest center, when measured according to $\varphi(x - c_l)$, is larger than φ^* and, as will

become clear later, it represents the set of outliers. From now on, given a set of centers $\{c_1, \dots, c_k\}$, we write just $A = \{A_1, \dots, A_k\}$ and implicitly mean $A \neq A_{k+1}$ when k is bounded.

For a similar definition for the set of clusters induced by a bounded γ see also Section 4 in [2].

30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

Now, consider the following instance of (GKM), for the same set of points $X \subset \mathbb{R}^n$, $n \geq 1$, $\min_{1 \leq l \leq k} |R_l| \geq 1$, $(c_1, \dots, c_k) = \min_{1 \leq i \leq k} \min_{1 \leq j \leq k} \|x_i - c_j\|_2^2 + f(\|x_i - c_j\|_2^2)$ $c_1, \dots, c_k \in \mathbb{R}^n$, $i = 1, \dots, k$, $j = 1, \dots, k$

subject to $c_l \in \mathbb{R}^n$, $l = 1, \dots, k$, $x_i \in \mathbb{R}^n$, $i = 1, \dots, n$,
(RKM)

where $f: \mathbb{R} \rightarrow \mathbb{R}_+$ is a symmetric, lsc, proper and bounded from below function, with minimum value $f(0)$, and γ a non-negative parameter. This problem is called robust k-means (RKM) and, as we show later, it takes the form of (GKM) when γ equals the Moreau envelope of f . The problem (RKM) [5, 24] describes the following simple model: we allow each observation x_i to take on an "error" term o_i and we penalize the errors, using a group penalty, in order to encourage most of the observations' errors to be equal to zero. We consider functions f where the parameter $\gamma \geq 0$ has the following effect: for $\gamma = 0$, all o_i 's may become arbitrary large (all observations are outliers), while, for $\gamma > 0$, all o_i 's become zero (no outliers); non-trivial cases occur for intermediate values $0 < \gamma < \infty$. Our interest is in understanding the robustness and consistency properties of (RKM). Robustness: Although robustness is an important notion, it has not been given a standard technical definition in the literature. Here, we focus on the finite sample breakdown point [18], which counts how many outliers a dataset may contain without causing significant damage in the estimates of the centers. Such damage is reflected to an arbitrarily large magnitude of at least one center. In Section 3, we show that two outliers in a dataset are enough to breakdown some centers. On the other hand, if we restrict our focus on some "well structured" datasets, then (RKM) has some remarkable robustness properties even if there is a considerable amount of contamination.

Consistency: Much is known about the consistency of (GKM) when the function f is lsc and increasing [11, 15]. It turns out that this case also includes the case of (RKM) when f is convex (see Section 3.1 for details). In Section 4, we show that the known non-asymptotic results about consistency of quadratic k-means may remain valid even when f is non-convex.

2

Preliminaries and some technical remarks

We start our analysis with a few technical tools from variational analysis [19]. Here, we introduce the necessary notation and a lemma (the proofs are in the appendix). The Moreau envelope $e_\gamma f(x)$ with parameter $\gamma \geq 0$ (Definition 1.22 in [19]) of an lsc, proper, and bounded from below function $f: \mathbb{R}^p \rightarrow \mathbb{R}$ and the associated (possibly multivalued) proximal map $P_\gamma f: \mathbb{R}^p \rightarrow \mathbb{R}^p$ are $e_\gamma f(x) = \min_{z \in \mathbb{R}^p} \gamma \|x - z\|^2 + f(z)$

$$\begin{aligned} & \frac{1}{2} \|x - z\|^2 + f(z) \text{ and } \text{Pf}^{\gamma}(x) = \underset{z \in \mathbb{R}^n}{\text{argmin}} \frac{1}{2} \|x - z\|^2 + \\ & f(z), \end{aligned} \quad (4)$$

respectively. In order to simplify the notation, in the following, we fix γ to 1 and suppress the superscript. The Moreau envelope is a continuous approximation from below of f having the same set of minimizers while the proximal map gives the (possibly non-unique) minimizing arguments in (4). For (GKM), we define $\text{eF} : \mathbb{R}^n \rightarrow \mathbb{R}$ as $\text{eF}(x) := \inf_{o \in \mathbb{R}^n} \left(\frac{1}{2} \|x - o\|^2 + f(o) \right)$. Accordingly, for (RKM), we define $\text{PF} : \mathbb{R}^n \rightarrow \mathbb{R}$ as $\text{PF}(x) := \text{f}^{\gamma} \left(\frac{x - \text{Pf}(x)}{2} \right)$. Thus, we obtain the following pairs: i) $(x, o) \mapsto \frac{1}{2} \|x - o\|^2 + f(o)$, $\text{Pf}(x) := \underset{o \in \mathbb{R}^n}{\text{argmin}} \frac{1}{2} \|x - o\|^2 + f(o)$, $x \mapsto \mathbb{R}$ (5a) ii) $\text{eF}(x) := \min_{o \in \mathbb{R}^n} \frac{1}{2} \|x - o\|^2 + F(o)$, $\text{PF}(x) := \underset{o \in \mathbb{R}^n}{\text{argmin}} \frac{1}{2} \|x - o\|^2 + F(o)$. (5b) Obviously, (RKM) is equivalent to (GKM) when $\text{PF}(x) = \text{eF}(x)$. Every map $P : \mathbb{R}^n \rightarrow \mathbb{R}^n$ throughout the text is assumed to be i) odd, i.e., $P(-x) = -P(x)$, ii) compact-valued, iii) non-decreasing, and iv) have a closed graph. We know that for any such map there exists at least one function f such that $P = \text{Pf}$ (Proposition 3 in [26]).³ Finally, for our purposes (outlier detection), it is natural $\text{ef}(x) := \min_{o \in \mathbb{R}^n} \frac{1}{2} \|x - o\|^2 + f(o)$.

² We call f proper if $f(x) < \infty$ for at least one $x \in \mathbb{R}^n$, and $f(x) \leq \infty$ for all $x \in \mathbb{R}^n$; in words, if the domain of f is a nonempty set on which f is finite (see page 5 in [19]).³ Accordingly, for a general function $f : \mathbb{R}^n \rightarrow [0, \infty]$ to be a Moreau envelope, i.e., $\text{PF}(x) = \text{ef}(x)$ as defined in (5a) for some function f , we require that $\text{PF}(x) - \frac{1}{2} \|x\|^2$ is a concave function (Proposition 1 in [26]).

² to require that v) P is a shrinkage rule, i.e., $P(x) \preceq x$, $\forall x \succeq 0$. The following corollary is quite straightforward and useful in the sequel. Corollary 1. Using the notation in definitions (5a) and (5b), we have $x \in \text{PF}^{\gamma}(x) = \text{Pf} \left(\frac{x - \text{PF}(x)}{2} \right)$ and $\text{eF}(x) = \text{ef} \left(\frac{x - \text{eF}(x)}{2} \right)$. (6) $\frac{1}{2} \|x - \text{PF}(x)\|^2$ Passing from a model of minimization in terms of a single problem, like (GKM), to a model in which a problem is expressed in a particular parametric form, like (RKM) with the Moreau envelope, the description of optimality conditions is opened to the incorporation of the multivalued map PF^{γ} . The next lemma describes the necessary conditions for a center cl to be (local) optimal for (RKM). Since we deal with the general case, well known results, such as smoothness of the Moreau envelope or convexity of its subgradients, can no longer be taken for granted. Remark 1. Let $\text{PF}(x) = \text{eF}(x)$. The usual subgradient, denoted as $\partial f(x)$, is not sufficient to characterize the differentiability properties of \mathbb{R}^n in (RKM). Instead, we use the (generalized) subdifferential $\partial^{\gamma} f(x)$ (Definition 8.3 in [19]). For all x , we have $\partial f(x) \subseteq \partial^{\gamma} f(x)$. Usually, the previous two sets coincide at a point x . In this case, f is called regular at x . However, it is common in practice that the sets $\partial f(x)$ and $\partial^{\gamma} f(x)$ are different (for a detailed exposition on subgradients see Chapter 8 in [19]; see also Example 1 in Appendix A.9). Lemma 1. Let $\text{PF}^{\gamma} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a proximal map and set $\text{PF}(x) = \text{eF}(x)$. The necessary (generalized) first order conditions for the centers $\{cl_1, \dots, cl_k\} \subset \mathbb{R}^n$ to be optimal for (RKM) are $\nabla f^{\gamma}(xi - cl_i) \cap \nabla f^{\gamma}(xi - cl_i) \cap (cl_i - xi + \text{PF}^{\gamma}(xi - cl_i)) \neq \emptyset$, $i \in \{1, \dots, k\}$. i) All

(7)

The interpretation of the set inclusion above is the following: for any center $c \in \mathbb{R}^p$, every subgradient vector in $\partial f^*(x) \cap \partial f^*(c)$ must be a vector associated with a vector in $\text{PF}^*(x) \cap \text{PF}^*(c)$ (Theorem 10.13 in [19]). However, in general, the converse does not hold true. We note that when the proximal map is single-valued and continuous, which happens for example not only when f is convex, but also for many popular non-convex penalties, both set inclusions become equalities and the converse holds, i.e., every vector in $\text{PF}^*(x) \cap \text{PF}^*(c)$ is a vector associated with a subgradient in $\partial f^*(x) \cap \partial f^*(c)$ (Theorem 10.13 in [19] and Proposition 7 in [26]). We close this section with some useful remarks on the properties of the Moreau envelope as a map between two spaces of functions. There exist cases where two different functions, $f \neq g$, have equal Moreau envelopes, $ef = eg$ (Proposition 1 in [26]), implying that two different forms of (RKM) correspond to the same \hat{f} in (GKM). For example, the proximal hull of f , defined as $h_f(x) := e^*(e^*(x))$, is a function different from f but has the same Moreau envelope as f .

(see also Example 1.44 in [19], Proposition 2 and Example 3 in [26]). This is the main reason we preferred the proximal map instead of the penalty function point of view for the analysis of (RKM).

3

On the breakdown properties of robust k-means

In this section, we study the finite sample breakdown point of (RKM) and, more specifically, its universal breakdown point. Loosely speaking, the breakdown point measures the minimum fraction of outliers that can cause excessive damage in the estimates of the centers. Here, it will become clear how the interplay between the two forms, (GKM) and (RKM), helps the analysis. Given a dataset $X \subset \mathbb{R}^n = \{x_1, \dots, x_n\}$ and a nonnegative integer $m \leq n$, we say that X_m is an m -modification if it arises from X after replacing m of its elements by arbitrary elements $x_i \in \mathbb{R}^p$ [6]. Denote as $r(\hat{f})$ the non-outlier samples, as counted after solving (RKM), for a dataset X and some $\hat{f} \in \mathcal{F}$, i.e., $r(\hat{f}) := \{x_i \in X : \inf_{c \in \mathcal{C}} \|x_i - c\| = 0, i = 1, \dots, n\}$. (8)

Then, the number of estimated outliers is $q(\hat{f}) = n - r(\hat{f})$. In order to simplify notation, we drop the dependence of r and q on \hat{f} . With this notation, we proceed to the following definition. 4

More than one \hat{f} can yield the same r , but this does not affect our analysis. 3

Definition 1 (universal breakdown point for the centers [6]). Let n, r, k be such that $n \geq r \geq k + 1$. Given a dataset X_m in \mathbb{R}^p , let $\{c_1, \dots, c_k\}$ denote the (global) optimal set of centers for (RKM). The universal breakdown value of (RKM) is $\beta_m(n, r, k) := \min_{X_m} \max_{\{c_1, \dots, c_k\}} \frac{\sum_{i=1}^k \|c_i - x_i\|}{n} = ?$. (9)

Here, $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^p$ while $X_m \subset \mathbb{R}^p$ runs over all m -modifications of X .

According to the concept of universal breakdown point, (RKM) breaks down at the first integer m for which there exists a set X such that the estimates of the cluster centers become arbitrarily bad for a suitable modification X_m . Our

analysis is based on Pf^* and considers two cases: those of biased and unbiased proximal maps. The former corresponds to the class of convex functions f^* , while the latter corresponds to a class of non-convex f^* . 3.1

Biased proximal maps: the case of convex f^*

If f^* is convex, then $f^* = eF^*$ is also convex while PF^* is continuous, single-valued, and satisfies [19] $\|x - PF^*(x)\|_2 \leq \|x\|_2$ as $\|x\|_2 \leq \frac{1}{2}$. (10) Proximal maps with this property are called biased since, as the l_2 -norm of x increases, so does the norm of the difference in (10). In this case, for each $x_i \in A_i$, from Lemma 1 and expression (10), we have $\|x_i - c_i\|_2 = \|eF^*(x_i - c_i)\|_2 = \|c_i - x_i + PF^*(x_i - c_i)\|_2 \leq \|x_i - c_i\|_2$. (11) The supremum value of $\|x - c_i\|_2$ is closely related to the gross error sensitivity of an estimator [9]. It is interpreted as the worst possible influence which a sample x can have on c_i [7]. In view of (11) and the definition of the clusters in (1), (RKM) is extremely sensitive. Although it can detect an outlier, i.e., a sample x_i with a nonzero estimate for c_i , it does not reject it since the influence of x_i on its closest center never vanishes. The l_1 -norm, $f^*(x) = \|x\|_1$, which has Moreau envelope equal to the Huber loss-function [24], is the limiting case for the class of convex penalty functions that, although it keeps the difference $\|x - PF^*(x)\|_2$ in (10) constant and equal to $\frac{1}{2}$, introduces a bias term proportional to $\frac{1}{2}$ in the estimate c_i . The following proposition shows that (RKM) with a biased PF^* has breakdown point equal to $\frac{1}{n}$, i.e., one outlier suffices to breakdown a center. Proposition 1. Assume $k \geq 2$, $k + 1 \leq r \leq n$. Given a biased proximal map, there exist a dataset $X \subset \mathbb{R}^n$ and a modification $X_1 \subset \mathbb{R}^n$ such that (RKM) breaks down. 3.2

Unbiased proximal maps: the case of non-convex f^*

Consider now the l_0 -(pseudo)norm on \mathbb{R} , $f^*(z) := \|z\|_0 = \begin{cases} 0 & \text{if } z=0 \\ 1 & \text{if } z \neq 0 \end{cases}$, thresholding proximal operator $Pf^* : \mathbb{R} \rightarrow \mathbb{R}$, $Pf^*(x) = \arg \min_z \|x - z\|_0 + f^*(z) = \{0, x\}$, $\forall x$,

According to Lemma 1, for $p = 1$ (scalar case), we have and the associated hard-thresholding operator $H_\tau : \mathbb{R} \rightarrow \mathbb{R}$, $H_\tau(x) = \begin{cases} x & \text{if } |x| \geq \tau \\ 0 & \text{if } |x| < \tau \end{cases}$.

(12)

$\|x_i - c_i\|_0 = \|x_i - c_i\|_0 + P_{\tau}^*(x_i - c_i) = \{0\}$ for $|x_i - c_i| < \tau$, $x_i \in A_i$,

(13)

implying that $\|x_i - c_i\|_0$, as a function of c_i , remains constant for $|x_i - c_i| < \tau$. As a consequence of (13), if c_i is local optimal, then $0 \leq \|x_i - c_i\|_0 \leq \|x_i - c_i\|_0$ and $\|x_i - c_i\|_0 = \|x_i - c_i\|_0 + P_{\tau}^*(x_i - c_i) = \|x_i - c_i\|_0$. (14) $\forall x_i \in A_i$, $\forall c_i \in A_i$, $\|x_i - c_i\|_0 = \|x_i - c_i\|_0$.

$\|x_i - c_i\|_0 = \|x_i - c_i\|_0$

Depending on the value of τ , (RKM) with the l_0 -norm is able to ignore samples with distance from their closest center larger than τ . This is done since $P_{\tau}^*(x_i - c_i) = x_i - c_i$ whenever $|x_i - c_i| < \tau$. 5

See the analysis in [7] about the influence function of (GKM) when f^* is convex.

4

and the influence of x_i vanishes. In fact, there is a whole family of non-convex f_γ 's whose proximal map Pf_γ satisfies $\text{Pf}_\gamma(x) = x$, for all $\|x\| \leq \gamma$,
(15)

for some $\gamma \geq 0$. These are called unbiased proximal maps [13, 20] and have the useful property that, as one observation is arbitrarily modified, all estimated cluster centers remain bounded by a constant that depends only on the remaining unmodified samples. Under certain circumstances, the proof of the following proposition reveals that, if there exists one outlier in the dataset, then robust k -means will reject it. Proposition 2. Assume $k \geq 2$, $k + 1 \leq r \leq n$, and consider the dataset $X_n = \{x_1, \dots, x_n\}$ along with its modification by one replacement y , $X_{1n} = \{x_1, \dots, x_{n-1}, y\}$. If we solve (RKM) with X_{1n} and an unbiased proximal map satisfying (15), then all estimates for the cluster centers remain bounded by a constant that depends only on the unmodified samples of X_n . Next, we show that, even for this class of maps, there always exists a dataset that causes one of the estimated centers to breakdown as two particular observations are suitably replaced. Theorem 1 (Universal breakdown point for (RKM)). Assume $k \geq 2$ and $n \geq r \geq k + 2$. Given an unbiased proximal map satisfying (15), there exist a dataset X_n and a modification X_{2n} , such that (RKM) breaks down. Hence, the universal breakdown point of (RKM) with an unbiased proximal map is $n/2$. In Figure 1, we give a visual interpretation of Theorem 1. The top subfigure depicts the unmodified initial dataset $X_9 = \{x_1, \dots, x_9\}$ (black circles) with a clear two-cluster structure; the bottom subfigure shows the modification X_{29} (dashed line arrows). Theorem 1 states that (RKM) on X_{29} fails to be robust since, every subset of X_{29} with $r = 8$ points has a cluster containing an outlier. 3.3

Figure 1: The top subfigure is the unmodified dataset X_9 . Theorem 1 states that every subset of the modification X_{29} (bottom subfigure) with size 8 contains an outlier.

Restricted robustness of robust k -means for well-clustered data

The result of Theorem 1 is disappointing but it is not (RKM) to be blamed for the poor performance but the tight notion of the definition about the breakdown point [6, 7]; allowing any kind of contamination in a dataset is a very general assumption. In this section, we place two restrictions: i) we consider datasets where inlier samples can be covered by unions of balls with centers that are "far apart" each other, and ii) we ask a question different from the finite sample breakdown point. We want to exploit as much as possible the results of [2] concerning a new quantitative measure of noise robustness which compares the output of (RKM) on a contaminated dataset to its output on the uncontaminated version of the dataset. Our aim is to show that (RKM), with a certain class of proximal maps and datasets that are well-structured ignores the influence of outliers when grouping the inliers. First, we introduce Corollary 2 which states the form that Pf_γ should have in order the results of [2] to apply to (RKM) and, second, we give details about the datasets which we consider as wellstructured. Using this corollary we are able to design proximal maps for which Theorems 3, 4, and 5 in [2] apply; otherwise, it is not clear how the

analysis of [2] is valid for (RKM). Let $h : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous function with the following properties: 1. h is odd and non-decreasing ($h_+ (?)$ is used to denote its restriction on $[0, ?)$); 2. h is a shrinkage rule: $0 \leq h_+(x) \leq x$, $\forall x \in [0, ?)$; 3. the difference $x - h_+(x)$ is non-decreasing, i.e., for $0 \leq x_1 \leq x_2$ we have $x_1 - h_+(x_1) \leq x_2 - h_+(x_2)$. 5

Define the map

$$\begin{aligned} \varphi(h(x), \text{Pf}(x)) &:= \{h(x), x\}, \quad \forall x, \\ -x- &\leq \varphi, \quad -x- = \varphi, \quad -x- \leq \varphi. \end{aligned}$$

(16)

Multivaluedness of $\text{Pf}(x)$ at $-x- = \varphi$ signals that $\text{ef}(x)$ is non-smooth at these points. An immediate consequence for the Moreau envelope associated with the previous map is the following. Corollary 2. Let the function $g : [0, ?) \rightarrow [0, ?)$ be defined as $\int_0^x g(u) du := (u - h_+(u))du$, $x \in [0, ?)$. (17) 0

Then, the Moreau envelope associated with $\text{Pf}(x)$ in (16) is

$$\text{ef}(x) = \min\{g(-x-), g(\varphi)\} = g(\min\{-x-, \varphi\}).$$

(18)

Next, we define what it means for a dataset to be (φ_1, φ_2) -balanced; this is the class of datasets which we consider to be well-structured. Definition 2 ((φ_1, φ_2) balanced dataset [2]). Assume that a set $X \subset \mathbb{R}^p$ has a subset I (inliers), with at least n_2 samples, and the following properties: 1. $I =$

$\bigcup_{l=1}^k$

I_l

B_l , where $B_l = B(\text{bl}_l, r)$ is a ball in \mathbb{R}^p with bounded radius r and center bl_l ;

2. $\varphi_1 - I_l \leq \varphi_2 - I_l$ for every l , where $-B_l -$ is the number of samples in B_l and $\varphi_1, \varphi_2 \geq 0$; 3. $\|\text{bl}_l - \text{bl}_{l'}\| \geq v$ for every $l \neq l'$, i.e., the centers of the balls are at least $v \geq 0$ apart. Then, X is a (φ_1, φ_2) -balanced dataset. We now state the form that Theorem 3 in [2] takes for (RKM). Theorem 2 (Restricted robustness of (RKM)). If i) $\text{ef}(x)$ is as in Corollary 2, i.e., $\text{ef}(x) = g(\min\{-x-, \varphi\})$, ii) X has a (φ_1, φ_2) -balanced subset of samples I with k balls, and iii) the centers of the balls are at least $v \geq 4r + 2g(\varphi_1 + \varphi_2 g(r))$ apart, then for $\varphi \in [1, 2]$ $\varphi_1 - I - \varphi \varphi_2 - I$ the set of outliers $X \setminus I$ has no effect on the 2, $g(X \setminus I) = (\varphi_1 g(\varphi_2 - 2r) + (\varphi_1 + \varphi_2)g(r))$

grouping of inliers I . In other words, if $\{x, y\} \in B_l$ and $\{c_1, \dots, c_k\}$ are the optimal centers when solving (RKM) for a φ as described before, then $l = \arg\min_{1 \leq j \leq k} \text{ef}(\|x - c_j\| - 2) = \arg\min_{1 \leq j \leq k} \text{ef}(\|y - c_j\| - 2)$. For the sake of completeness, we give a proof of this theorem in the appendix. In a similar way, we can recast the results of Theorems 4 and 5 in [2] to be valid for (RKM).

4

On the consistency of robust k-means

Let X be a set with n independent and identically distributed random samples x_i from a fixed but unknown probability distribution φ . Let C be the empirical optimal set of centers, i.e., $C := \arg\min_{c_1, \dots, c_k \in \mathbb{R}^p} \sum_{i=1}^n (c_1, \dots, c_k)$.

(19)

$$C := \operatorname{argmin}_{c_1, \dots, c_k} R(c_1, \dots, c_k),$$

(20)

The population optimal set of centers is the set

$$C := \min_{c_1, \dots, c_k} \int \sum_{i=1}^k \|x - c_i\|^2 + f(\|x - c_i\|^2) dx.$$

(21)

$$\int \|x - c_i\|^2 = \int \|x - c_i\|^2$$

Loss consistency and (simply) consistency for (RKM) require, respectively, that $R_n(C) \rightarrow R(C)$ and $C_n \rightarrow C$.

6

(22)

C_n converges In words, as the size n of the dataset X_n increases, the empirical clustering risk $R_n(C)$ almost surely to the minimum population risk $R(C)$ and (for n large enough) C_n can effectively replace the optimal set C in quantizing the unknown probability measure μ . For the case of convex f , non-asymptotic results describing the rate of convergence of R_n to R in (22) are already known ([11], Theorem 3). Noting that the Moreau envelope of a non-convex f belongs to a class of functions with polynomial discrimination [16] (the shatter coefficient of this class is bounded by a polynomial) we give a sketch proof of the following result. Theorem 3 (Consistency of (RKM)). Let the samples $x_i \in X_n$, $i = 1, \dots, n$, come from a fixed but unknown probability measure μ . For any $k \geq 1$ and any unbiased proximal map, we have $R_n(C_n) \rightarrow R(C)$ and $\lim E R_n(C) = R(C)$.

$n \rightarrow \infty$

$$\lim_{n \rightarrow \infty} C_n = C \quad (\text{convergence in probability}).$$

$n \rightarrow \infty$

(23)

Theorem 3 reads like an asymptotic convergence result. However, its proof (given in the appendix) uses combinatorial tools from Vapnik-Chervonenkis theory, revealing that the non-asymptotic rate $R_n(C_n) - R(C)$ is of order $O(\log n/n)$ (see Corollary 12.1 in [4]).

5

Relating (RKM) to trimmed k-means

As the effectiveness of robust k-means on real world and synthetic data has already been evaluated [5, 24], the purpose of this section is to relate (RKM) to trimmed k-means (TKM) [7]. Trimmed kmeans is based on the methodology of "impartial trimming", which is a combinatorial problem fundamentally different from (RKM). Despite their differences, the experiments show that, both (RKM) and (TKM) perform remarkably similar in practice. The solution of (TKM) (which is also a set of k centers) is the solution of quadratic k-means on the subsample containing $n(1 - \alpha)$ points with the smallest mean deviation ($0 \leq \alpha \leq 1$). The only common characteristic of (RKM) and (TKM) is that they both have the same universal breakdown point, i.e., $n/2$, for arbitrary datasets. Trimmed k-means takes as input a dataset X_n , the number of clusters k , and

a proportion of outliers $\alpha \in (0, 1)$ to remove.⁶ A popular heuristic algorithm for (TKM) is the following. After the initialization, each iteration of (TKM) consists of the following steps: i) the distance of each observation from its closest center is computed, ii) the top αn observations with larger distance from its closest center are removed, iii) the remaining points are used to update the centers. The previous three steps are repeated until the centers converge.⁷ As for robust k-means, we solve the (RKM) problem with a coordinate optimization procedure (see Appendix A.9 for details). The synthetic data for the experiments come from a mixture of Gaussians with 10 components and without any overlap between them.⁸ The number of inlier samples is 500 and each inlier $x_i \in [1, 1]^{10}$ for $i \in \{1, \dots, 500\}$. On top of the inliers lie 150 outliers in \mathbb{R}^{10} distributed uniformly in general positions over the entire space. We consider two scenarios: in the first, the outliers lie in $[3, 3]^{10}$ (call it mild-contamination), while, in the second, the outliers lie in $[6, 6]^{10}$ (call it heavy-contamination). The parameter α in trimmed k-means (the percentage of outliers) is set to $\alpha = 0.3$, while the value of the parameter γ for which (RKM) yields 150 outliers is found through a search over a grid on the set $\gamma \in (0, \gamma_{\max})$ (we set γ_{\max} as the maximum distance between two points in a dataset). Both algorithms, as they are designed, require as input an initial set of k points; these points form the initial set of centers. In all experiments, both (RKM) and (TKM) take the same k vectors as initial centers, i.e., k points sampled randomly from the dataset. The statistics we use for the comparison are: i) the rand-index for clustering accuracy [17] ii) the cluster estimation error, i.e., the root mean square error between the estimated cluster centers and the sample mean of each cluster, iii) the true positive outlier detection rate, and finally, iv) the false positive outlier detection rate. In Figures 2-3, we plot the results for a proximal map P_f like the one in (16) with $h(x) = \|x\|$ and $\eta = 0.005$; with this choice for h , we mimic the hard-thresholding operator. The results for each scenario (accuracy, cluster estimation error, etc) are averages over 150 runs of the experiment. As seen, both algorithms share almost the same statistics in all cases. ⁶

We use the implementation of trimmed k-means in the R package `trimcluster` [10]. The previous three steps are performed also by another robust variant of k-means, the k-means⁺ (see [3]). ⁸ We use the R toolbox `MixSim` [14] that guarantees no overlap among the 10 mixtures. ⁷

7
0.5
0.005
0.3
7.5
5.0
0.950
0.925
robust k-means trimmed k-means

robust k?means trimmed k?means
 0.015
 0.010
 ?
 0.005
 ?
 ?
 0.000
 robust k?meanstrimmed k?means
 Cluster Radius Estimation Error
 ? ? ? ? ? ?
 10.0
 0.975
 False Positive Error Rate
 ? ?
 ?
 True Positive Error Rate
 Accuracy
 0.7
 12.5
 ?
 Center Estimation Error
 ? ?
 0.9
 ? ? ?
 9
 ? ? ?
 6
 3
 ? ? ?
 ?
 0
 robust k?meanstrimmed k?means
 robust k?means trimmed k?means
 0.4
 ? ? ?
 ?
 ? ? ? ? ? ?
 0.2
 ? ?
 15.0
 12.5
 10.0
 ?
 robust k?means trimmed k?means
 ?

0.75 ? ?
 0.50 ? ?
 0.25 ? ? ? ?
 0.00
 ?
 0.3
 ? ?
 ? ?
 False Positive Error Rate
 0.6
 1.00
 ?
 True Positive Error Rate
 Accuracy
 0.8
 Center Estimation Error
 17.5
 robust k?means trimmed k?means
 ? ?
 0.2
 Cluster Radius Estimation Error
 Figure 2: Performance of robust and trimmed k-means on a mixture of 10
 Gaussians without overlap. On top of the 500 samples from the mixture there
 are 150 outliers uniformly distributed in $[-1, 1]^{10}$. ? ? ? ? ? ?
 ? ?
 ? ?
 0.1 ?
 0.0
 robust k?means trimmed k?means
 ? ? ?
 4
 3
 2
 1
 robust k?means trimmed k?means
 robust k?means trimmed k?means
 ? ? ? ? ? ? ? ?
 ? ? ? ? ? ? ?
 ? ?
 20
 1.000
 ?
 ? ?
 10
 ?
 robust k?means trimmed k?means

?
 ?
 ? ?
 ?
 0.975
 ?
 ?
 0.950
 0.925
 0.900 robust k?means trimmed k?means
 ?
 ? ?
 ?
 ?
 ?
 ? ?
 ?
 ?
 robust k?meanstrimmed k?means
 False Positive Error Rate
 0.6
 Center Estimation Error
 Accuracy
 0.8
 ?
 True Positive Error Rate
 30
 1.0
 0.04 ? ?
 ?
 0.03 ? ?
 0.02 0.01
 ? ?
 ?
 ?
 ?
 ?
 ? ? ?
 0.00 robust k?means trimmed k?means

Figure 4: Results on two spherical clusters with equal radius r , each one with 150 samples, and centers are at least $4r$ apart. On top of the samples lie 150 outliers uniformly distributed in $[?6, 6]10$. In Figure 4, we plot the results for the case of two spherical clusters in $R10$ with equal radius r , each one with 150 samples, and centers that are at least $4r$ apart from each other. The inlier samples are in $[?3, 3]10$. The outliers are 150 (half of the dataset is contaminated) and are uniformly distributed in $[?6, 6]10$. The results (accuracy,

cluster estimation error, etc) are averages over 150 runs of the experiment. This configuration is a heavy contamination scenario but, due to the structure of the dataset, as expected from Theorem 2, (RKM) performs remarkably well; the same holds for (TKM).

6

Conclusions

We provided a theoretical analysis for the robustness and consistency properties of a variation of the classical quadratic k-means called robust k-means (RKM). As a by-product of the analysis, we derived a detailed description of the optimality conditions for the associated minimization problem. In most cases, (RKM) shares the computational simplicity of quadratic k-means, making it a ?computationally cheap? candidate for robust nearest neighbor clustering. We show that (RKM) cannot be robust against any type of contamination and any type of datasets, no matter the form of the proximal map we use. If we restrict our attention to ?well-structured? datasets, then the algorithm exhibits some desirable noise robustness. As for the consistency properties, we showed that most general results for consistency of quadratic k-means still remain valid for this robust variant. Acknowledgments The author would like to thank Athanasios P. Liavas for useful comments and suggestions that improved the quality of the article.

8

Cluster Radius Estimation Error

Figure 3: The same setup as in Figure 2 except that the coordinates of each outlier lie in $[?3, 3]^{10}$. 7.5

5.0

2.5

0.0 robust k-means trimmed k-means

2 References

- [1] Anestis Antoniadis and Jianqing Fan. Regularization of wavelet approximations. *Journal of the American Statistical Association*, 2011.
- [2] Shai Ben-David and Nika Haghtalab. Clustering in the presence of background noise. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 280?288, 2014.
- [3] Sanjay Chawla and Aristides Gionis. k-means-: A unified approach to clustering and outlier detection. *SIAM*.
- [4] L. Devroye, L. Gy?orfi, and G. Lugosi. *A Probabilistic Theory of Pattern Recognition. Stochastic Modelling and Applied Probability*. Springer New York, 1997.
- [5] Pedro A Forero, Vassilis Kekatos, and Georgios B Giannakis. Robust clustering using outlier-sparsity regularization. *Signal Processing, IEEE Transactions on*, 60(8):4163?4177, 2012.
- [6] Mar??a Teresa Gallegos and Gunter Ritter. A robust method for cluster analysis. *Annals of Statistics*, pages 347?380, 2005.
- [7] Luis Angel Garc??a-Escudero and Alfonso Gordaliza. Robustness properties of k-means and trimmed k-means. *Journal of the American Statistical Association*, 94(447):956?969, 1999.
- [8] Michael R. Garey and David S. Johnson.

Computers and Intractability: A Guide to the Theory of NP-Completeness. W. H. Freeman & Co., New York, NY, USA, 1979. [9] Frank R Hampel, Elvezio M Ronchetti, Peter J Rousseeuw, and Werner A Stahel. Robust statistics: the approach based on influence functions, volume 114. John Wiley & Sons, 2011. [10] Christian Hennig. trimcluster: Cluster analysis with trimming, 2012. R package version 0.1-2. [11] Tamás Linder. Learning-theoretic methods in vector quantization. In Principles of nonparametric learning, pages 163–210. Springer, 2002. [12] Stuart P Lloyd. Least squares quantization in pcm. Information Theory, IEEE Transactions on, 28(2):129–137, 1982. [13] Rahul Mazumder, Jerome H Friedman, and Trevor Hastie. Sparsenet: Coordinate descent with nonconvex penalties. Journal of the American Statistical Association, 2012. [14] Volodymyr Melnykov, Wei-Chen Chen, and Ranjan Maitra. MixSim: An R package for simulating data to study performance of clustering algorithms. Journal of Statistical Software, 51(12):1–25, 2012. [15] David Pollard. Strong consistency of k-means clustering. The Annals of Statistics, 9(1):135–140, 1981. [16] David Pollard. Convergence of stochastic processes. Springer Science & Business Media, 1984. [17] William M Rand. Objective criteria for the evaluation of clustering methods. Journal of the American Statistical association, 66(336):846–850, 1971. [18] G. Ritter. Robust Cluster Analysis and Variable Selection. Chapman & Hall/CRC Monographs on Statistics & Applied Probability. CRC Press, 2014. [19] R Tyrrell Rockafellar and Roger J-B Wets. Variational analysis, volume 317. Springer Science & Business Media, 2009. [20] Yiyuan She et al. Thresholding-based iterative selection procedures for model selection and shrinkage. Electronic Journal of statistics, 3:384–415, 2009. [21] Marc Teboulle. A unified continuous optimization framework for center-based clustering methods. The Journal of Machine Learning Research, 8:65–102, 2007. [22] Paul Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. Journal of optimization theory and applications, 109(3):475–494, 2001. [23] Sara Van De Geer. Empirical processes in m-estimation. June 13, 2003. Handout at New Directions in General Equilibrium Analysis (Cowles Workshop, Yale University). [24] Daniela M Witten. Penalized unsupervised learning with outliers. Statistics and its Interface, 6(2):211, 2013. [25] Stephen J Wright. Coordinate descent algorithms. Mathematical Programming, 151(1):3–34, 2015. [26] Yaoliang Yu, Xun Zheng, Micol Marchetti-Bowick, and Eric P Xing. Minimizing nonconvex nonseparable functions. In AISTATS, 2015.