

# A Simple and Fast Algorithm for $L_1$ -norm Kernel PCA

Cheolmin Kim, Diego Klabjan

**Abstract**—We present an algorithm for  $L_1$ -norm kernel PCA and provide a convergence analysis for it. While an optimal solution of  $L_2$ -norm kernel PCA can be obtained through matrix decomposition, finding that of  $L_1$ -norm kernel PCA is not trivial due to its non-convexity and non-smoothness. We provide a novel reformulation through which an equivalent, geometrically interpretable problem is obtained. Based on the geometric interpretation of the reformulated problem, we present a “fixed-point” type algorithm that iteratively computes a binary weight for each observation. As the algorithm requires only inner products of data vectors, it is computationally efficient and the kernel trick is applicable. In the convergence analysis, we show that the algorithm converges to a local optimal solution in a finite number of steps. Moreover, we provide a rate of convergence analysis, which has been never done for any  $L_1$ -norm PCA algorithm, proving that the sequence of objective values converges at a linear rate. In numerical experiments, we show that the algorithm is robust in the presence of entry-wise perturbations and computationally scalable, especially in a large-scale setting. Lastly, we introduce an application to outlier detection where the model based on the proposed algorithm outperforms the benchmark algorithms.

**Index Terms**—Principal Component Analysis,  $L_1$ -norm, Kernel, Outlier Detection.



## 1 INTRODUCTION

PRINCIPAL Component Analysis (PCA) is one of the most popular dimensionality reduction techniques [1]. Given a large set of possibly correlated features, it attempts to find a small set of features (*principal components*) that retain as much information as possible. To generate such new dimensions, it linearly transforms original features by multiplying *loading vectors* in a way that newly generated features are orthogonal and have the largest variance.

In traditional PCA, variance is measured using the  $L_2$ -norm. This has a nice property in that although the problem itself is non-convex, an optimal solution can be easily found through matrix factorization. With this property and easy interpretability, PCA is extensively used in a variety of applications. Nonetheless, it still has some limitations. First, since it generates a new dimension through a linear combination of features, it cannot capture non-linear relationships among features. Second, as it uses the  $L_2$ -norm for measuring variance, its outcome tends to be affected by influential outliers. In order to overcome these limitations, the following two approaches have been proposed.

**Kernel PCA** The idea of kernel PCA is to map original features into a high-dimensional feature space, and perform PCA in that high-dimensional feature space [2]. Using a non-linear mapping, it can capture non-linear relationships among features in an efficient way using the *kernel trick*. Using the trick, principal components can be computed with no explicit mapping.

**$L_1$ -norm PCA** To alleviate the effects of influential outliers,  $L_1$ -norm PCA uses the  $L_1$ -norm instead of the  $L_2$ -norm to measure variance. The  $L_1$ -norm is more advanta-

geous than the  $L_2$ -norm in presence of observations having large feature values since it is less influenced by them. Using this property, more robust results can be obtained by  $L_1$ -norm PCA in the presence of influential outliers.

In this paper, we combine the two approaches for the variance maximization version of  $L_1$ -norm PCA. In what follows, we always refer to the variance maximization version of  $L_1$ -norm PCA which is not the same as minimizing reconstruction error with respect to the  $L_1$ -norm. Compared to  $L_2$ -norm kernel PCA, the kernel version of  $L_1$ -norm PCA is a hard problem in that it is not only non-convex but also non-smooth. However, through a novel reformulation, we convert it to a geometrically interpretable problem where the objective is to minimize the  $L_2$ -norm of a vector subject to a linear constraint consisting of terms involving the  $L_1$ -norm. For the reformulated problem, we present a “fixed point” type algorithm that iteratively computes a weight of  $-1$  or  $1$  for each observation using the kernel matrix and previous weights. We show that the kernel trick is applicable to this algorithm. Moreover, we prove that the algorithm converges to a local optimal solution in a finite number of steps and the sequence of objective values converges at a linear rate. In numerical experiments, we computationally investigate the robustness of the algorithm and introduce an application to outlier detection. We also provide a runtime comparison to other robust kernel PCA algorithms and  $L_2$ -norm kernel PCA.

Our work has the following contributions.

1. We provide a novel reformulation of  $L_1$ -norm kernel PCA and present an iterative algorithm based on the geometric interpretation of the reformulated problem. This approach is not specific to  $L_1$ -norm kernel PCA but can be applied to a more general problem. Particularly, its application to  $L_2$ -norm PCA results in Power iteration [3].

• Cheolmin Kim and Diego Klabjan are with the Department of Industrial Engineering and Management Science, Northwestern University, Evanston, IL, 60208.  
E-mail: cheolmkim@u.northwestern.edu, d-klabjan@northwestern.edu

2. We not only prove convergence but also provide a rate of convergence analysis. Although many algorithms have been proposed for  $L_1$ -norm PCA, none of them provided a rate of convergence analysis. We stress that our analysis is for the kernel version which clearly covers  $L_1$ -norm PCA. Through a novel analysis, we show that the algorithm attains a linear rate of convergence.
3. We introduce a methodology based on  $L_1$ -norm kernel PCA for outlier detection and demonstrate that it outperforms the benchmark algorithms.

The paper is organized as follows. Section 2 reviews related works and points out how our work is different. Section 3 introduces a novel reformulation of  $L_1$ -norm kernel PCA and provides a geometric interpretation behind it. Based on the geometric interpretation, we present an iterative algorithm in Section 4. Section 5 provides a convergence analysis for it and the experimental results are followed in Section 6.

## 2 RELATED WORK

Extracting a low-rank representation from a large matrix is an important problem in machine learning and statistics. In a variety of contexts, many previous works [4], [5], [6], [7] have been proposed to address this problem. Recovering a low-rank matrix from a sampling of its entries is studied in [4]. Given that the number of sampled entries is sufficiently large, exact recovery is guaranteed with high probability by solving a simple convex optimization problem [4]. Assuming that a data matrix can be decomposed into the sum of a low-rank matrix  $L_0$  and a sparse matrix  $S_0$ , a convex program (known as *robust PCA*) that minimizes a weighted combination of the nuclear norm of  $L_0$  and the  $L_1$  norm of  $S_0$  is presented in [5]. Also, a variant of robust PCA that identifies outliers by additionally imposing a column-sparse structure on  $S_0$  is considered in [6]. Under some mild conditions, exact recovery is shown for both models [5], [6]. Moreover, exact recovery of mixture data is studied in [7], [8], [9], [10]. Utilizing a dictionary matrix, low-rank representation (LRR) [7] is shown to better handle mixture data than robust PCA. While matrix recovery is the main focus of these works, our work considers dimensionality reduction with emphasis on robustness, especially focusing on kernel PCA with the  $L_1$ -norm.

To reduce the number of features in a robust way, the  $L_1$ -norm has been involved in many PCA studies [11], [12], [13], [14], [15], [16], [17] and subspace estimation formulations [18], [19]. Finding a subspace onto which the  $L_1$  projections of data vectors have the smallest reconstruction error is studied in [11]. Based on the observation that the  $L_1$  projection occurs along a single unit direction, it finds an optimal subspace for each unit direction by solving  $d$  least absolute deviation regression problems, each having one dimension as a dependent variable while having the other dimensions as independent variables. Using linear programming, this approach can find a global optimal subspace in polynomial time [11].

Minimizing reconstruction error with respect to the  $L_1$ -norm is considered in [13], [14], [18]. While the PCA problem of minimizing  $\|M - XX^T M\|_1$  subject to  $X^T X = I$  is considered in [13], the subspace estimation problem of

minimizing  $E(U, V) = \|M - UV\|_1$  is studied in [18] where  $M$  is a data matrix. In order to solve the former problem, an iterative algorithm that computes a weight for each observation and applies  $L_2$ -norm PCA on the weighted data matrix is presented in [13]. On the other hand, the latter problem is solved using alternative convex minimization based on the observation that  $E(U, V)$  becomes a convex function once  $U$  or  $V$  is known. It alternatively optimizes one matrix at a time while keeping the other one fixed, repeating this process until convergence. Also, a subspace estimation formulation that minimizes reconstruction error with respect to the  $R_1$ -norm,  $\|M - UV\|_{R_1} = \sum_{i=1}^n \|x_i - Uv_i\|_2$  where  $x_i$  is the  $i^{th}$  column of  $M$  and  $v_i$  is that of  $V$ , is presented in [19]. Since this formulation minimizes the sum of distances with respect to the  $L_2$ -norm, it is different from  $L_2$ -norm PCA which minimizes the sum of squared distances with respect to the  $L_2$ -norm. Nonetheless, they share the same property that they have a unique global solution which is rotational invariant [19].

Maximizing variance with respect to the  $L_1$ -norm, which we refer to as  $L_1$ -norm PCA, is studied in [12], [15], [16], [17]. Our work also considers this formulation rather than the previous two since it has a favorable structure in that an optimal solution can be represented as a linear combination of data vectors with a weight of  $-1$  or  $1$ .  $L_1$ -norm PCA is shown to be NP-hard in [17] and [16]. Nevertheless, an algorithm finding a global optimal solution is proposed in [17]. Utilizing the auxiliary-unit-vector technique [20], it computes a global optimal solution with complexity  $O(n^{pr+p-1})$  where  $n$  is the number of observations,  $r$  is the rank of the data matrix, and  $p$  is the desired number of principal components. Assuming  $r$  and  $p$  are fixed, the runtime of this algorithm is polynomial in  $n$ . However, if  $n, p, r$  are large, it can be computationally prohibitive. Instead of finding a global optimal solution which is intractable in general, our work focuses on developing an efficient algorithm finding a local optimal solution for  $L_1$ -norm kernel PCA.

Recognizing the hardness of  $L_1$ -norm PCA, an approximation algorithm is presented in [16] based on the known Nesterov's theorem [21]. In this work,  $L_1$ -norm PCA is relaxed to a semi-definite programming (SDP) problem and alternatively, the SDP relaxation is considered. After solving the relaxed problem, it generates a random vector and uses randomized rounding to produce a feasible solution. This randomized algorithm is a  $\sqrt{2/\pi}$ -approximate algorithm in expectation. To achieve this approximation ratio with high probability, it performs randomized rounding multiple times and takes the one having the best objective value. Rather than providing an approximation guarantee by solving a relaxed problem, our work directly considers the kernel version of  $L_1$ -norm PCA and develops an efficient algorithm finding a local optimal solution.

Another approach utilizing a known mathematical programming model is introduced in [12] where the author proposes an iterative algorithm that solves a mixed integer programming problem in each iteration. Given an orthonormal matrix of loading vectors, it perturbs the matrix slightly in a way that the resulting matrix yields the largest objective value. After the perturbation, it uses singular value decomposition to recover orthogonality. The algorithm is completely different from the one proposed herein and the

sequence of objective values does not necessarily improve over iterations. Unlike it, our algorithm guarantees that the sequence of objective values keeps improving and converges at a linear rate.

A simple numerical algorithm finding a local optimal solution is proposed in [22]. In this work, an optimal solution is assumed to have a certain form, and weights involved in that form are updated in each iteration, improving the objective value. A similar algorithm and its extended version that finds multiple loading vectors at once are derived in [15] utilizing an optimization algorithm for general  $L_1$ -norm maximization problems. In the case of linear kernel, our algorithm uses the same framework as the one in [22] and [15]. However, while the algorithm in [22] is derived without any justification, we provide a geometric interpretation behind the algorithm, which is different from the derivation in [15]. Moreover, we provide a rate of convergence analysis and introduce a kernel version, which are not considered in [22] and [15].

On other hand, the kernel version of  $L_1$ -norm PCA has been rarely studied. Due to the difficulty of applying the kernel trick to  $L_1$ -norm kernel PCA, an alternative method named *nonlinear projection trick* is applied in [23]. Based on the finding that an optimal loading vector lies in the span of  $\Phi(A)^T U \Lambda^{-1/2}$  where  $\Phi(A)$  is a high-dimensionally mapped data matrix and  $U \Lambda U^T$  is the eigenvalue decomposition of the kernel matrix  $K$ , it alternatively considers  $L_1$ -norm PCA having  $U \Lambda^{1/2}$  in place of  $\Phi(A)$  and solves it using the algorithm in [22]. Another kernel extension of  $L_1$ -norm PCA is studied in [24]. In this work, a linear system involving a kernel matrix is solved in each iteration and the resulting solution is used to update the iterate. While the algorithms in [23] and [24] entail either eigenvalue decomposition or solving a linear system, our algorithm requires only a matrix-vector multiplication in each iteration, making it suitable in a large-scale setting.

### 3 KERNEL-BASED $L_1$ -NORM PCA FORMULATIONS

We consider  $L_1$ -norm PCA in a high-dimensional feature space  $F$ . Suppose we map data vectors  $a_i \in \mathbb{R}^d, i = 1, \dots, n$  into a feature space  $F$  by a possibly non-linear mapping  $\Phi : \mathbb{R}^d \rightarrow F$ . Assuming that each feature is standardized with a mean of 0 and standard deviation of 1 and that the kernel matrix  $K$  defined by  $K_{ij} = \Phi(a_i)^T \Phi(a_j)$  satisfies

- 1)  $K_{ii} > 0$  for  $1 \leq i \leq n$
- 2)  $|K_{ij}| < \infty$  for  $1 \leq i, j \leq n$ ,

the kernel version of  $L_1$ -norm PCA is formulated as

$$\begin{aligned} & \underset{x \in F}{\text{maximize}} && f(x) = \sum_{i=1}^n |\Phi(a_i)^T x| \\ & \text{subject to} && \|x\|_2 = 1. \end{aligned} \quad (1)$$

This formulation having  $\Phi(a_i)$  in place of  $a_i$  extends the variance maximization version of  $L_1$ -norm PCA in the obvious way and is also considered in [23], [24]. In this formulation, we only consider extracting the first loading vector. This assumption is justifiable since the subsequent loading vectors can be found by repeatedly solving (1). For example, once we obtain the first loading vector  $x^*$ ,

we can find the second loading vector by solving (1) with  $\Phi(a_i) - x^*(\Phi(a_i)^T x^*)$  in place of  $\Phi(a_i)$ .

Solving (1) is not trivial since it has a convex non-smooth objective function to maximize and a Euclidean unit ball constraint. In order to better understand the problem and set an algorithmic foundation, we reformulate (1) as

$$\begin{aligned} & \underset{x \in F}{\text{minimize}} && g(x) = \|x\|_2 \\ & \text{subject to} && \sum_{i=1}^n |\Phi(a_i)^T x| = 1. \end{aligned} \quad (2)$$

In order to prove the equivalence of (1) and (2), we argue that an optimal solution of one formulation can be derived from an optimal solution of the other formulation by means of some mapping. Two optimization problems are equivalent if there exists some mapping  $h$  such that if  $x^*$  is an optimal solution to one problem, then  $h(x^*)$  is an optimal solution to the other problem, and vice versa for a possible different mapping function [25].

**Proposition 1.** *Let  $x_1^*$  and  $y_2^*$  be an optimal solution to (1) and (2), respectively. Then,*

$$x_2^* = \frac{x_1^*}{\sum_{i=1}^n |\Phi(a_i)^T x_1^*|}$$

*is an optimal solution to (2), and*

$$y_1^* = \frac{y_2^*}{\|y_2^*\|_2}$$

*is an optimal solution to (1).*

*Proof.* It is easy to check that  $x_2^*$  is a feasible solution to (2). Suppose that  $x_2^*$  is not optimal to (2). Then, there exists some feasible  $z$  such that

$$\|z\|_2 < \|x_2^*\|_2.$$

As  $z$  is feasible to (2), we have

$$\sum_{i=1}^n |\Phi(a_i)^T z| = 1.$$

Let  $w = \frac{z}{\|z\|_2}$ . Then, we have

$$f(w) = \sum_{i=1}^n |\Phi(a_i)^T w| = \frac{\sum_{i=1}^n |\Phi(a_i)^T z|}{\|z\|_2} = \frac{1}{\|z\|_2}.$$

In the same way, we obtain

$$f(x_1^*) = \frac{1}{\|x_2^*\|_2}$$

since

$$x_1^* = \frac{x_2^*}{\|x_2^*\|_2} = \frac{\sum_{i=1}^n |\Phi(a_i)^T x_1^*|}{\sum_{i=1}^n |\Phi(a_i)^T x_1^*|} \frac{x_2^*}{\|x_2^*\|_2} = \frac{x_2^*}{\|x_2^*\|_2}.$$

This leads to

$$f(x_1^*) < f(w),$$

which contradicts the assumption that  $x_1^*$  is an optimal solution of (1). Therefore,  $x_2^*$  is optimal to (2)

On the other hand, it is obvious that  $y_1^*$  is feasible to (1). To derive a contradiction, suppose that  $y_1^*$  is not optimal to (1). Then, there exists some feasible  $w$  such that

$$\sum_{i=1}^n |\Phi(a_i)^T y_1^*| < \sum_{i=1}^n |\Phi(a_i)^T w|.$$

Let

$$z = \frac{w}{\sum_{i=1}^n |\Phi(a_i)^T w|}.$$

Then, we have

$$g(z) = \frac{\|w\|_2}{\sum_{i=1}^n |\Phi(a_i)^T w|} = \frac{1}{\sum_{i=1}^n |\Phi(a_i)^T w|}$$

since  $\|w\|_2 = 1$ . In the same way, we obtain

$$g(y_2^*) = \frac{1}{\sum_{i=1}^n |\Phi(a_i)^T y_1^*|}.$$

for

$$y_2^* = \frac{y_1^*}{\sum_{i=1}^n |\Phi(a_i)^T y_1^*|}$$

due to  $\|y_1^*\|_2 = 1$ . As a result, we have

$$g(y_2^*) > g(z),$$

contradicting the assumption that  $y_2^*$  is optimal to (2). Therefore,  $y_1^*$  is optimal to (1).  $\square$

To understand formulation (2), we first look at the constraint set,

$$\partial P = \left\{ x \mid \sum_{i=1}^n |\Phi(a_i)^T x| = 1 \right\}.$$

Geometrically, this constraint set is symmetric with respect to the origin and represents the boundary of polytope

$$P = \left\{ x \mid \sum_{i=1}^n |\Phi(a_i)^T x| \leq 1 \right\}.$$

It is easy to check that  $P$  is a polytope since it can be written as the intersection of a finite set of linear inequalities each having the form of  $\sum_{i=1}^n c_i \Phi(a_i)^T x \leq 1$  where  $c_i \in \{-1, 1\}$ . As the objective function measures the distance from the origin, formulation (2) can be understood as a problem of finding the closest point to the origin from the boundary of the polytope  $\partial P$ . The following proposition shows that an optimal solution  $x^*$  must be perpendicular to one of the faces of  $\partial P$ .

**Proposition 2.** *An optimal solution  $x^*$  is perpendicular to the face which it lies on.*

*Proof.* Let  $x^*$  be an optimal solution of (2) and define a face  $E$  such that

$$E = \left\{ x \mid \sum_{i=1}^n c_i^* \Phi(a_i)^T x = 1 \right\} \cap \partial P$$

where

$$c_i^* = \text{sgn}(\Phi(a_i)^T x^*)$$

for  $1 \leq i \leq n$ . If  $x^*$  is not perpendicular to face  $E$ , then

$$w = \frac{\sum_{i=1}^n \Phi(a_i) c_i^*}{\left\| \sum_{i=1}^n \Phi(a_i) c_i^* \right\|_2}$$

is the closest point to the origin from

$$\left\{ x \mid \sum_{i=1}^n c_i^* \Phi(a_i)^T x = 1 \right\}$$

having

$$\|w\|_2 < \|x^*\|_2. \quad (3)$$

Let

$$z = \frac{w}{\sum_{i=1}^n |\Phi(a_i)^T w|}.$$

Then,  $z$  is feasible to (2) and has the objective value of

$$\|z\|_2 = \frac{\|w\|_2}{\sum_{i=1}^n |\Phi(a_i)^T w|}. \quad (4)$$

From

$$\left\| \sum_{i=1}^n \Phi(a_i) c_i^* \right\|_2^2 = \sum_{i=1}^n \Phi(a_i)^T c_i^* \left( \sum_{j=1}^n \Phi(a_j) c_j^* \right),$$

we have

$$\sum_{i=1}^n |\Phi(a_i)^T \left( \sum_{j=1}^n \Phi(a_j) c_j^* \right)| - \left\| \sum_{i=1}^n \Phi(a_i) c_i^* \right\|_2^2 \geq 0$$

resulting in

$$\sum_{i=1}^n |\Phi(a_i)^T w| = \frac{\sum_{i=1}^n |\Phi(a_i)^T \left( \sum_{j=1}^n \Phi(a_j) c_j^* \right)|}{\left\| \sum_{i=1}^n \Phi(a_i) c_i^* \right\|_2} \geq 1. \quad (5)$$

As a result, by (3), (4), and (5), we have

$$\|z\|_2 \leq \|w\|_2 < \|x^*\|_2,$$

which contradicts the assumption that  $x^*$  is optimal to (2). Therefore,  $x^*$  must be perpendicular to  $E$ .  $\square$

Proposition 2 is important since it helps to characterize the form of an optimal solution  $x^*$ . From Proposition 2, we obtain the following corollary.

**Corollary 1.** *An optimal solution  $x^*$  of (2) has the form of*

$$x^* = \frac{y^*}{\sum_{i=1}^n |\Phi(a_i)^T y^*|}$$

for some  $y^*$  and  $c^*$  such that

$$y^* = \sum_{i=1}^n \Phi(a_i) c_i^*$$

and

$$c_i^* = \text{sgn}(\Phi(a_i)^T y^*),$$

for  $1 \leq i \leq n$ .

The characterization of an optimal loading vector using a sign vector is first proposed in [22] without any justification. However, we provide a derivation based on the geometry of  $\partial P$ , which is different from the one in [15] that uses the KKT conditions. Moreover, since we have

$$\|x^*\|_2 = \frac{\|y^*\|_2}{\sum_{i=1}^n |\Phi(a_i)^T y^*|} = \frac{1}{\left\| \sum_{i=1}^n \Phi(a_i) c_i^* \right\|_2} \quad (6)$$

due to

$$\sum_{i=1}^n |\Phi(a_i)^T y^*| = \sum_{i=1}^n c_i^* \Phi(a_i)^T y^* = \left\| \sum_{i=1}^n \Phi(a_i) c_i^* \right\|_2^2,$$

we can further show that an optimal solution of formulation (2) can be found from an optimal solution of the following binary problem,

$$\text{maximize}_{c \in \{-1,1\}^n} \left\| \sum_{i=1}^n \Phi(a_i) c_i \right\|_2^2. \quad (7)$$

**Proposition 3.** *Let  $c^*$  be an optimal solution of binary formulation (7). Then,*

$$y^* = \sum_{i=1}^n \Phi(a_i) c_i^*$$

satisfies

$$c_i^* = \text{sgn}(\Phi(a_i)^T y^*), \quad (8)$$

for  $1 \leq i \leq n$ . Moreover,

$$x^* = \frac{y^*}{\sum_{i=1}^n |\Phi(a_i)^T y^*|}$$

is an optimal solution of formulation (2).

*Proof.* To deduce a contradiction, let us assume that there exists some nonempty set  $J \subset \{1, \dots, n\}$  such that

$$c_j^* = -\text{sgn}(\Phi(a_j)^T y^*)$$

for  $j \in J$ . Since  $c^*$  is an optimal solution of (7), flipping the sign of  $c_j^*$  for  $j \in J$  must not improve the objective value of (7). However, for any  $j \in J$ , flipping the sign of  $c_j^*$  results in

$$\|y - 2\Phi(a_j) c_j^*\|_2^2 > \|y\|_2^2$$

since

$$\|y^* - 2\Phi(a_j) c_j^*\|_2^2 = \|y^*\|_2^2 + 4|y^{*T}(\Phi(a_j))| + 4\|\Phi(a_j)\|_2^2$$

and  $\|\Phi(a_j)\|_2^2 > 0$ . This contradicts the assumption that  $c^*$  is an optimal solution to (7). Therefore,  $y^*$  must satisfy

$$c_i^* = \text{sgn}(\Phi(a_i)^T y^*)$$

for  $1 \leq i \leq n$ . Since  $y^*$  and  $c^*$  satisfy (8) and  $c^*$  maximizes the objective value of (7),

$$x^* = \frac{y^*}{\sum_{i=1}^n |\Phi(a_i)^T y^*|}$$

is a minimizer of (2) due to Corollary 1 and (6).  $\square$

The following result has been shown in [17] for the linear kernel case but here we generalize it.

**Corollary 2.** *Formulation (2) is equivalent to formulation (7).*

*Proof.* Based on Corollary 1 and (6), we can formulate (2) as

$$\begin{aligned} & \text{maximize}_{c \in \{-1,1\}^n} && \left\| \sum_{i=1}^n \Phi(a_i) c_i \right\|_2^2 \\ & \text{subject to} && y = \sum_{i=1}^n \Phi(a_i) c_i \\ & && c_i = \text{sgn}(\Phi(a_i)^T y), \quad 1 \leq i \leq n. \end{aligned}$$

Since an optimal solution  $c^*$  to (7) satisfies the constraints of the above optimization problem by Proposition 3, the two formulations are essentially the same.  $\square$

It is interesting to note that we can reduce formulation (7) to the weighted max-cut problem since

$$\left\| \sum_{i=1}^n \Phi(a_i) c_i \right\|_2^2 = \sum_{i,j=1}^n K_{ij} + \sum_{i,j=1}^n (-2K_{ij}) \left( \frac{1 - c_i c_j}{2} \right). \quad (9)$$

Using the above reduction, we can alternatively consider the weighted max-cut problem on a complete graph with weight  $w_{ij} = -K_{ij}$ . Therefore, a popular approximation algorithm for the weighted max-cut problem [26] can be used to solve (7). However, due to the additional constant terms in (9), this does not imply a constant worst case approximation ratio algorithm for (7).

## 4 ALGORITHM

In this section, we develop an algorithm that finds a local optimal solution to (2) based on the findings in Section 3. Before giving details of the algorithm, we first provide the idea behind the algorithm.

The main idea of the algorithm is to move along the boundary of  $P$  so that the  $L_2$ -norm of  $x_k$  successively decreases. Figure 1 illustrates a step of the algorithm. Starting with an iterate  $x^k$ , we first identify the hyperplane  $h^k$  which the current iterate  $x^k$  lies on. After identifying the equation of  $h^k$ , we find the closest point to the origin from  $h^k$ , which we denote by  $z^k$ . After that, we obtain  $x^{k+1}$  by projecting  $z^k$  to the constraint set  $\partial P$ , which is done by multiplying an appropriate scalar between 0 and 1. We repeat this process until the sequence of iterates  $\{x^k\}$  converges.

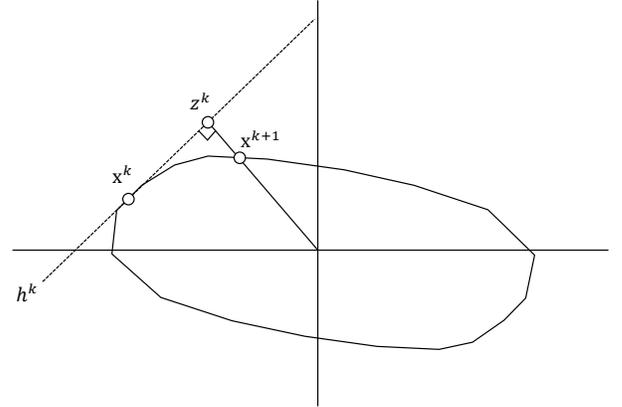


Fig. 1. Geometric derivation of the algorithm

Now, we develop an algorithm based on the above idea. Given  $x^k$ , we define the normal vector  $y^k$  at  $x^k$  by

$$y^k = \sum_{i=1}^n \Phi(a_i) c_i^k \quad (10)$$

using the sign vector  $c^k = [c_1^k, \dots, c_n^k]^T$  defined by

$$c_i^k = \text{sgn}(\Phi(a_i)^T x^k)$$

for  $1 \leq i \leq n$ . Using the normal vector  $y^k$  at  $x^k$ , we can find the equation of hyperplane  $h^k$  as

$$(y^k)^T(x - x^k) = 0. \quad (11)$$

The closest point  $z^k$  to the origin from  $h^k$  has the form of

$$z^k = sy^k. \quad (12)$$

Plugging (12) into (11), we have

$$s = \frac{(y^k)^T x^k}{(y^k)^T y^k}$$

resulting in

$$z^k = \frac{(y^k)^T x^k}{(y^k)^T y^k} y^k. \quad (13)$$

Projecting  $z^k$  to  $\partial P$ , we obtain

$$x^{k+1} = \frac{z^k}{\sum_{i=1}^n |\Phi(a_i)^T z^k|}. \quad (14)$$

Using

$$(y^k)^T x^k = \sum_{i=1}^n \Phi(a_i)^T x^k c_i^k = \sum_{i=1}^n |\Phi(a_i)^T x^k| = 1, \quad (15)$$

we can further write (13) as

$$z^k = \frac{y^k}{\|y^k\|_2^2} \quad (16)$$

leading to

$$x^{k+1} = \frac{y^k}{\sum_{i=1}^n |\Phi(a_i)^T y^k|}. \quad (17)$$

Also, from (10) and

$$\sum_{i=1}^n |\Phi(a_i)^T y^k| = \sum_{i=1}^n \Phi(a_i)^T y^k c_i^k = (c^k)^T K c^k,$$

we can represent  $x^{k+1}$  as a function of  $c^k$  by

$$x^{k+1} = \frac{\sum_{i=1}^n \Phi(a_i) c_i^k}{(c^k)^T K c^k}. \quad (18)$$

Since

$$c_i^{k+1} = \text{sgn}(\Phi(a_i)^T x^{k+1}) = \text{sgn}(K_i \cdot c^k),$$

we can update  $c_i^{k+1}$  using only  $K$  and  $c^k$  by

$$c^{k+1} = \text{sgn}(K c^k).$$

Moreover, from

$$\|x^{k+1} - x^k\|_2^2 = \frac{(c^k - c^{k+1})^T K (c^k - c^{k+1})}{(c^k)^T K c^k (c^{k+1})^T K c^{k+1}},$$

we can represent the termination criteria  $x^{k+1} = x^k$  by

$$(c^k - c^{k+1})^T K (c^k - c^{k+1}) = 0.$$

On the other hand, due to non-convexity of the problem, the algorithm can be stuck at a local optimum unless it is initialized close to a global optimum. In order to obtain a good initial iterate  $x^0$ , we consider each  $\Phi(a_j)$  and select the

one such that  $\Phi(a_j)/\|\Phi(a_j)\|_2$  yields the largest objective value for  $f$ , which is computed by

$$\frac{\sum_{i=1}^n |\Phi(a_i)^T \Phi(a_j)|}{\|\Phi(a_j)\|_2} = \frac{\sum_{i=1}^n |K_{ij}|}{\sqrt{K_{jj}}}. \quad (19)$$

Once we find the index  $j^*$  maximizing (19), we set

$$x^0 = \frac{\Phi(a_{j^*})}{\sum_{i=1}^n |\Phi(a_i)^T \Phi(a_{j^*})|}$$

resulting in

$$c_i^0 = \text{sgn}(\Phi(a_i)^T x^0) = \text{sgn}(\Phi(a_i)^T \Phi(a_{j^*})) = \text{sgn}(K_{ij^*}).$$

Since an optimal loading vector  $x^*$  must be located somewhere between  $\Phi(a_i)$  where  $1 \leq i \leq n$ , the above initialization scheme is likely to yield an initial iterate  $x^0$  close to the optimal loading vector  $x^*$ .

Summarizing all the above, we obtain Algorithm 1.

---

#### Algorithm 1 $L_1$ -norm Kernel PCA

---

**Input:** kernel matrix  $K$

find  $j^* = \arg \max_{1 \leq j \leq n} \sum_{i=1}^n |K_{ij}| / \sqrt{K_{jj}}$

initialize the sign vector  $c^0$  with  $c_i^0 = \text{sgn}(K_{ij^*})$

$k \leftarrow -1$

**repeat**

$k \leftarrow k + 1$

compute  $c^{k+1} = \text{sgn}(K c^k)$

**until**  $(c^k - c^{k+1})^T K (c^k - c^{k+1}) = 0$

**Output:** sign vector  $c^*$

---

Once we get the output  $c^*$  from Algorithm 1, we can compute principal scores with no explicit mapping. For example, the principal component of the  $i^{\text{th}}$  observation can be computed by

$$\begin{aligned} \frac{\Phi(a_i)^T x^*}{\|x^*\|_2} &= \frac{\sum_{j=1}^n \Phi(a_i)^T \Phi(a_j) c_j^*}{\sqrt{\sum_{i=1}^n \sum_{j=1}^n \Phi(a_i)^T \Phi(a_j) c_i^* c_j^*}} \\ &= \frac{K_i \cdot c^*}{\sqrt{(c^*)^T K c^*}}. \end{aligned}$$

Also, we can proceed to find more principal components with no explicit mapping. Noting that computing a loading vector and principal components requires only the kernel matrix, it suffices to update the kernel matrix each time a new loading vector is found. Fortunately, updating the kernel matrix can be done with no explicit mapping by

$$\begin{aligned} \tilde{K}_{ij} &= \left( \Phi(a_i) - \frac{\Phi(a_i)^T x^*}{\|x^*\|_2^2} x^* \right)^T \left( \Phi(a_j) - \frac{\Phi(a_j)^T x^*}{\|x^*\|_2^2} x^* \right) \\ &= \Phi(a_i)^T \Phi(a_j) - \frac{\Phi(a_i)^T x^* \Phi(a_j)^T x^*}{\|x^*\|_2^2} \\ &= K_{ij} - \frac{K_i \cdot c^* K_j \cdot c^*}{(c^*)^T K c^*}, \end{aligned}$$

which is equivalent to

$$\tilde{K} = K - \frac{(K c^*)(K c^*)^T}{(c^*)^T K c^*}$$

in a matrix form.

From  $y_k = \nabla f(x_k)$ , update rule (17) can be understood as projecting a gradient  $\nabla f(x_k)$  to the constraint set  $\partial P$  in

each iteration. In this sense, Algorithm 1 resembles Power iteration [3] for solving the eigenvalue problem, and interestingly, the application of our framework to the eigenvalue problem yields the same algorithm. The framework developed in this work such as reformulation, geometric interpretation and algorithm derivation is not specific to  $L_1$ -norm kernel PCA but can be extended to solve a more general problem. For example, our approach can be used to solve

$$\text{maximize } f(x) \quad \text{subject to } \|x\|_2 = 1$$

for any function  $f$  that is scale-invariant (*homogeneous* or *homothetic*). The application of our framework to this problem yields the following update rule

$$x^{k+1} \leftarrow \nabla f(x^k) / \|\nabla f(x^k)\|_2.$$

Compared to the other  $L_1$ -norm kernel PCA algorithms [23], [24] considering the same formulation (1), Algorithm 1 is much simple and computationally efficient as it involves just one matrix-vector multiplication in each iteration. In the case of L1-KPCA [24], a system of linear equations having the form of

$$K\eta = \sum_{j=1}^n c_j^k K_{\cdot j}$$

is repeatedly solved. Solving the above linear system is not only computationally costly but also numerically unstable since it is singular due to the presence of non-trivial solution  $c^k$ . On the other hand, KPCA-L1 [23] requires one matrix-vector multiplication but it does not directly consider the kernel matrix  $K$ . Instead, the eigenvalue decomposition of the kernel matrix  $K = U\Lambda U^T$  must be computed before starting to find each loading vector. Also,  $U\Lambda^{1/2}$  is involved in computation instead of the kernel matrix  $K$ . As Algorithm 1 entails neither solving a linear system nor computing the eigenvalue decomposition of  $K$ , it is computationally more efficient than the other algorithms.

When it comes to initialization, L1-KPCA [24] uses the optimal loading vector from  $L_2$ -norm kernel PCA. While KPCA-L1 [23] finds the data vector having the largest norm and uses its normalization for the initial iterate, Algorithm 1 finds the normalized data vector with the largest objective value for  $f$  and set it to be the initial iterate. As the initialization scheme of Algorithm 1 is based on the objective function  $f$  while the others are not, it is more likely to obtain a good initial iterate compared to the others.

## 5 CONVERGENCE ANALYSIS

In this section, we provide a convergence analysis of Algorithm 1. We first prove that the algorithm converges in a finite number of iterations, and then provide a rate of convergence analysis. Before proving the finite convergence of the algorithm, we first show that the sequence  $\{\|x^k\|_2\}$  generated by Algorithm 1 is non-increasing.

**Lemma 1.** *Let  $\{x_k\}$  and  $\{z_k\}$  be a sequence of vectors generated by Algorithm 1 and (16), respectively. Then, we have*

$$\|x^{k+1}\|_2 \leq \|z^k\|_2 \leq \|x^k\|_2.$$

Moreover, if  $\|x^k\|_2 = \|z^k\|_2$ , we have  $x^k = ry^k$  for some  $r \in \mathbb{R}$ .

*Proof.* The inequality  $\|z^k\|_2 \leq \|x^k\|_2$  follows from

$$\begin{aligned} \|x^k\|_2^2 - \|z^k\|_2^2 &= \|x^k\|_2^2 - \frac{1}{\|y^k\|_2^2} \\ &= \|x^k\|_2^2 - \frac{((y^k)^T x^k)^2}{\|y^k\|_2^2} \\ &= \frac{\|x^k\|_2^2 \|y^k\|_2^2 - ((y^k)^T x^k)^2}{\|y^k\|_2^2} \\ &\geq 0 \end{aligned}$$

where the second equality holds follows from (15) and the last inequality holds due to the Cauchy-Schwarz inequality. If  $\|x^k\|_2 = \|z^k\|_2$ , the Cauchy-Schwarz inequality becomes an equality resulting in

$$x^k = ry^k$$

for some  $r \in \mathbb{R}$ .

Next, from (14), we have

$$\|x^{k+1}\|_2^2 = \frac{\|z^k\|_2^2}{(\sum_{i=1}^n |\Phi(a_i)^T z^k|)^2}. \quad (20)$$

Using (13), we can represent the denominator as

$$\sum_{i=1}^n |\Phi(a_i)^T z^k| = \frac{\sum_{i=1}^n |\Phi(a_i)^T y^k|}{(y^k)^T y^k}.$$

From

$$\begin{aligned} \sum_{i=1}^n |\Phi(a_i)^T y^k| &= \sum_{i=1}^n |\Phi(a_i)^T (\sum_{j=1}^n \Phi(a_j) c_j^k)| \\ &= \sum_{i=1}^n |\sum_{j=1}^n \Phi(a_i)^T \Phi(a_j) c_i^k c_j^k| \end{aligned}$$

and

$$(y^k)^T y^k = \sum_{i=1}^n \sum_{j=1}^n \Phi(a_i)^T \Phi(a_j) c_i^k c_j^k,$$

we obtain

$$\sum_{i=1}^n |\Phi(a_i)^T z^k| = \frac{\sum_{i=1}^n |\sum_{j=1}^n \Phi(a_i)^T \Phi(a_j) c_i^k c_j^k|}{\sum_{i=1}^n \sum_{j=1}^n \Phi(a_i)^T \Phi(a_j) c_i^k c_j^k} \quad (21)$$

resulting in

$$\sum_{i=1}^n |\Phi(a_i)^T z^k| \geq 1. \quad (22)$$

By (20) and (22), we have

$$\|x^{k+1}\|_2^2 \leq \|z^k\|_2^2.$$

□

**Lemma 2.** *If*

$$\|x^k\|_2 = \|x^{k+1}\|_2,$$

*then, we have*

$$x^k = \frac{y^k}{\|y^k\|_2}, \quad y^k = \frac{x^k}{\|x^k\|_2},$$

*resulting in*

$$x^k = x^{k+1}.$$

*Proof.* Since  $\|x^k\|_2 = \|x^{k+1}\|_2$ , we have

$$\|z^k\|_2 = \|x^k\|_2, \quad x^k = ry^k$$

by Lemma 1 for some  $r \in \mathbb{R}$ . Using (15), we have

$$r = \frac{1}{\|y^k\|_2^2}$$

resulting in

$$x^k = \frac{y^k}{\|y^k\|_2^2}.$$

In the same way, we can show

$$y^k = \frac{x^k}{\|x^k\|_2^2}.$$

Since this implies  $z^k = x^k$  by (16), we finally have

$$x^{k+1} = \frac{z^k}{\sum_{i=1}^n |\Phi(a_i)^T z^k|} = \frac{x^k}{\sum_{i=1}^n |\Phi(a_i)^T x^k|} = x^k$$

where the first equality follows from (14) and the last equality holds from the feasibility of  $x^k$ .  $\square$

**Theorem 1.** *The sequence  $\{x^k\}$  converges in a finite number of steps.*

*Proof.* Suppose the sequence  $\{x^k\}$  does not converge. As an iterate  $x^k$  is solely determined by a sign vector  $c^k \in \{-1, +1\}^n$ , the number of possible vectors that  $x^k$  can take is finite. Therefore, if the sequence  $\{x^k\}$  does not converge, some vectors must appear more than once. Without loss of generality, let  $x^l = x^{l+m}$ . By Lemma 1, we have

$$\|x^{l+m}\|_2 = \|x^l\|_2 \geq \|x^{l+1}\|_2 \geq \dots \geq \|x^{l+m}\|_2$$

forcing us to have

$$\|x^l\|_2 = \|x^{l+1}\|_2 = \dots = \|x^{l+m}\|_2.$$

This implies

$$x^l = x^{l+1} = \dots = x^{l+m}$$

by Lemma 2, contradicting the assumption that the sequence  $\{x^k\}$  does not converge. Therefore, the sequence  $\{x^k\}$  generated by Algorithm 1 must converge in a finite number of steps.  $\square$

Next, we show that the sequence of  $\{\|x^k\|_2\}$  generated by Algorithm 1 converges at a linear rate. Although Theorem 1 shows that the algorithm converges in a finite number of steps, it may take an exponential number of steps to converge, due to the combinatorial structure of the problem, making it not appropriate in a large-scale setting. To make sure that this does not happen for Algorithm 1, we additionally prove linear convergence, which ensures that the optimality gap decreases no worse than a certain rate  $\rho < 1$ . Since this result implies that an  $\epsilon$ -optimal local solution can be attained after  $\mathcal{O}(1/(1-\rho) \log(1/\epsilon))$  iterations, we can obtain a near-optimal solution after a sufficient number of iterations without waiting for an exponential number of steps.

**Theorem 2.** *Let Algorithm 1 start from  $x^0$  and terminate with  $x^*$  at iteration  $k^*$ . Then, for some  $\rho < 1$ , we have*

$$\|x^k\|_2 - \|x^*\|_2 \leq \rho^k (\|x^0\|_2 - \|x^*\|_2)$$

for  $k < k^*$ .

*Proof.* From (14), we have

$$\|x^k\|_2 = \frac{\|z^{k-1}\|_2}{\sum_{i=1}^n |\Phi(a_i)^T z^{k-1}|}.$$

Since  $\|z^{k-1}\|_2 \leq \|x^{k-1}\|_2$  holds by Lemma 1, we obtain

$$\|x^k\|_2 \leq \frac{\|x^{k-1}\|_2}{\sum_{i=1}^n |\Phi(a_i)^T z^{k-1}|}. \quad (23)$$

Subtracting  $\|x^*\|_2$  to (23), we have

$$\begin{aligned} \|x^k\|_2 - \|x^*\|_2 &\leq \frac{\|x^{k-1}\|_2}{\sum_{i=1}^n |\Phi(a_i)^T z^{k-1}|} - \|x^*\|_2 \\ &\leq \frac{1}{\sum_{i=1}^n |\Phi(a_i)^T z^{k-1}|} (\|x^{k-1}\|_2 - \|x^*\|_2) \end{aligned} \quad (24)$$

where the last inequality follows from (22).

By induction on (24), we obtain

$$\|x^k\|_2 - \|x^*\|_2 \leq (\|x^0\|_2 - \|x^*\|_2) \prod_{l=1}^k \frac{1}{\sum_{i=1}^n |\Phi(a_i)^T z^{l-1}|}. \quad (25)$$

From (22), we know that

$$\sum_{i=1}^n |\Phi(a_i)^T z^{l-1}| \geq 1.$$

If

$$\sum_{i=1}^n |\Phi(a_i)^T z^{l-1}| = 1,$$

we have

$$\frac{\sum_{i=1}^n |\sum_{j=1}^n \Phi(a_i)^T \Phi(a_j) c_i^{l-1} c_j^{l-1}|}{\sum_{i=1}^n \sum_{j=1}^n |\Phi(a_i)^T \Phi(a_j) c_i^{l-1} c_j^{l-1}|} = 1$$

resulting in

$$c_i^{l-1} = \text{sgn}\left(\sum_{j=1}^n \Phi(a_i)^T \Phi(a_j) c_j^{l-1}\right).$$

Since this implies

$$c^l = \text{sgn}(Kc^{l-1}) = c^{l-1},$$

we have

$$x^l = x^{l+1}.$$

Therefore, as long as  $l < k^*$ , we must have

$$\sum_{i=1}^n |\Phi(a_i)^T z^{j-1}| > 1.$$

For  $c \in \{-1, 1\}^n$ , let

$$\rho(c) = \frac{\sum_{i=1}^n \sum_{j=1}^n \Phi(a_i)^T \Phi(a_j) c_i c_j}{\sum_{i=1}^n |\sum_{j=1}^n \Phi(a_i)^T \Phi(a_j) c_i c_j|}$$

and define

$$\rho = \max_{c \in \{-1, 1\}^n} \rho(c) \text{ subject to } \rho(c) < 1.$$

Then, for  $l < k^*$ , we have

$$\frac{1}{\sum_{i=1}^n |\Phi(a_i)^T \mathbf{z}^{j-1}|} = \rho(c^{j-1}) < \rho < 1.$$

By combining it with (25), we get the desired result.  $\square$

As shown in Theorem 2, no matter where the algorithm starts, the sequence of objective values of (2) converges at a linear rate. Now, we show that we can obtain a local optimal solution of (1) by scaling the output of Algorithm 1.

**Theorem 3.** *Let the output of Algorithm 1 be  $\mathbf{x}^*$ . Then,*

$$\bar{\mathbf{x}}^* = \frac{\mathbf{x}^*}{\|\mathbf{x}^*\|_2}$$

is a local optimal solution of (1).

*Proof.* It is easy to see that  $\bar{\mathbf{x}}^*$  is feasible. Since  $\mathbf{x}^*$  is the output of Algorithm 1,

$$\mathbf{y}^* = \frac{\mathbf{x}^*}{\|\mathbf{x}^*\|_2^2}$$

holds by Lemma 2. Next, consider

$$L(\lambda, \mathbf{x}) = \sum_{i=1}^n |\Phi(a_i)^T \mathbf{x}| - \lambda(\|\mathbf{x}\|_2^2 - 1).$$

From

$$\nabla_{\mathbf{x}} L(\lambda, \mathbf{x}) = \sum_{i=1}^n \text{sgn}(\Phi(a_i)^T \mathbf{x}) \Phi(a_i) - 2\lambda \mathbf{x},$$

we have

$$\begin{aligned} \nabla_{\mathbf{x}} L(\lambda, \bar{\mathbf{x}}^*) &= \sum_{i=1}^n \text{sgn}(\Phi(a_i)^T \bar{\mathbf{x}}^*) \Phi(a_i) - 2\lambda \bar{\mathbf{x}}^* \\ &= \sum_{i=1}^n \text{sgn}(\Phi(a_i)^T \mathbf{x}^*) \Phi(a_i) - 2\lambda \bar{\mathbf{x}}^* \\ &= \mathbf{y}^* - 2\lambda \bar{\mathbf{x}}^* \\ &= \frac{\mathbf{x}^*}{\|\mathbf{x}^*\|_2^2} - 2\lambda \bar{\mathbf{x}}^* \\ &= \left( \frac{1}{\|\mathbf{x}^*\|_2} - 2\lambda \right) \bar{\mathbf{x}}^*. \end{aligned}$$

Therefore, with

$$\lambda^* = \frac{1}{2\|\mathbf{x}^*\|_2},$$

we have

$$\nabla_{\mathbf{x}} L(\lambda^*, \bar{\mathbf{x}}^*) = 0,$$

meaning that  $(\lambda^*, \bar{\mathbf{x}}^*)$  satisfies the first-order necessary conditions. Moreover, from

$$\nabla_{\mathbf{xx}} L(\lambda^*, \bar{\mathbf{x}}^*) = -2\lambda^* I \prec 0,$$

the second-order sufficient condition is also satisfied. Since  $(\lambda^*, \bar{\mathbf{x}}^*)$  satisfies the first and second order conditions, from the theory of constrained optimization,  $\bar{\mathbf{x}}^*$  is a local optimal solution of (1).  $\square$

## 6 EXPERIMENTAL RESULTS

In this section, we assess the robustness and scalability of Algorithm 1 by running it on several tasks and compare it with other kernel PCA algorithms. First, we apply them on datasets having entry-wise perturbations and investigate how well each algorithm extracts principal components in a noisy setting. Next, we introduce their application to outlier detection and compare their performance with other popular outlier detection models. Lastly, we provide their runtime comparison.

In addition to Algorithm 1, the two other  $L_1$ -norm kernel PCA algorithms (KPCA-L1 [23], L1-KPCA [24]), the kernel version of  $R_1$ -norm PCA (R1-KPCA [19]) and  $L_2$ -norm kernel PCA (L2-KPCA [2]) are considered in the experiments. While  $R_1$ -norm PCA [19] is not originally designed to incorporate kernels, we include it as it is easy to develop a kernel variant. Other  $L_1$ -norm PCA algorithms were also considered but since it is not straightforward to develop a kernel version for them, they are disregarded.

### 6.1 Robust Extraction of PCs

To measure robustness, we first run the algorithms on datasets having entry-wise perturbations (noisy datasets) to obtain loading vectors. After that, we compute how much variation in the perturbation-excluded datasets (normal datasets) is explained by the loading vectors obtained from the noisy datasets. For this experiment, we prepare synthetic datasets having entry-wise perturbations so that loading vectors obtained by running  $L_2$ -norm kernel PCA on noisy and normal datasets are different from each other.

To generate synthetic datasets, we first construct a  $1000 \times 50$  data matrix with the rank of 10 following the data generation procedure in [13]. While the largest size in [13] is  $300 \times 50$ , we choose the size of  $1000 \times 50$  to consider larger datasets. To obtain entry-wise perturbations, we corrupt  $r\%$  of observations by adding some random noises. We refer to the resulting dataset as a noisy dataset and the noisy dataset without the entry-wise perturbations as a normal dataset. For each value of  $r \in \{5, 10, 15, 20, 25, 30\}$ , we generate 10 instances.

Let  $K$  denote a kernel matrix of a normal dataset and  $\mathbf{x}_1, \dots, \mathbf{x}_p$  be  $p$  loading vectors obtained by running  $L_2$ -norm kernel PCA on  $K$ . Also, let  $\tilde{K}$  be a kernel matrix of a noisy dataset and  $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_p$  be loading vectors obtained by running one of the kernel PCA algorithms (Algorithm 1, KPCA-L1, L1-KPCA, R1-KPCA, L2-KPCA) on  $\tilde{K}$ . Assuming that the normal dataset is standardized,

$$\sum_{j=1}^p \sum_{i=1}^n (\Phi(a_i)^T \tilde{\mathbf{x}}_j)^2 = \sum_{j=1}^p \tilde{\mathbf{x}}_j^T K \tilde{\mathbf{x}}_j \quad (26)$$

represents the amount of variation in the normal dataset explained by the  $p$  loading vectors  $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_p$  where  $n$  is the number of observations in the normal dataset. After dividing (26) by  $\sum_{j=1}^p \mathbf{x}_j^T K \mathbf{x}_j$ , which is the maximum amount of variation in the normal dataset that the  $p$  orthogonal vectors can explain, and multiplying by 100, we get the following measure:

$$(\text{Total Explained Variation}) \quad 100 \times \frac{\sum_{j=1}^p \tilde{\mathbf{x}}_j^T K \tilde{\mathbf{x}}_j}{\sum_{j=1}^p \mathbf{x}_j^T K \mathbf{x}_j}. \quad (27)$$

Metric (27) captures how well the loading vectors obtained from the noisy dataset explain variation in the normal dataset with respect to the  $L_2$ -norm. Therefore, it can be used to measure the robustness of each kernel PCA algorithm in the presence of entry-wise perturbations. For example, if one algorithm has a value close to one, then it is robust with respect to entry-wise perturbations. Using this metric, we compare the robustness of Algorithm 1 with that of KPCA-L1, L1-KPCA, R1-KPCA, and L2-KPCA. For each value of  $r$ , we compute (27) for the ten datasets with  $p = 4$  and average them. We arbitrarily choose  $p = 4$  since the result is consistent regardless of the choice of  $p$ . Figure 2 shows the results for the linear kernel and Figure 3 shows the results for the Gaussian kernel with the width parameter  $\sigma$  varying from 10 to 25.

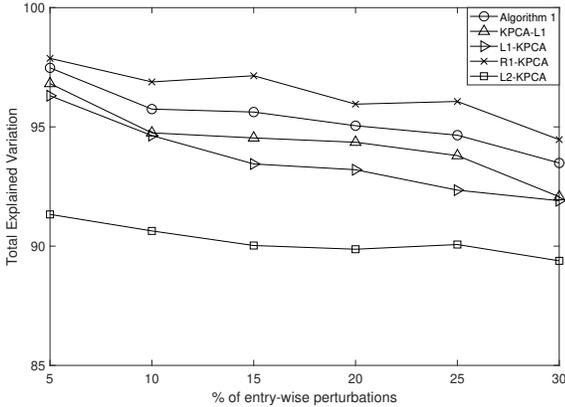


Fig. 2. Robust Extraction of PCs (Linear Kernel)

In the case of the linear kernel, R1-KPCA achieves the best performance for all values of  $r$  followed by the  $L_1$ -norm based kernel PCA algorithms and L2-KPCA. While the loading vectors from L2-KPCA explain about 90% of the variation, those from R1-KPCA, Algorithm 1, KPCA-L1, and L1-KPCA explain around 96%, 95%, 94%, and 93% of the variation, respectively. This demonstrates the robustness of the  $R_1$ -norm and  $L_1$ -norm based kernel PCA algorithms with respect to the presence of entry-wise perturbations. Among the three  $L_1$ -norm based kernel PCA algorithms, Algorithm 1 consistently outperforms KPCA-L1 and L1-KPCA by 1% and 2%, respectively. As the percentage of corrupted observations ( $r\%$ ) increases, the total explained variation tends to decrease for all of them but the gaps between them remain the same.

When the Gaussian kernel is used, the results are slightly different depending on the value of  $r$  and  $\sigma$ . If  $r$  and  $\sigma$  are small, the effects of entry-wise perturbations are relatively small so that all the algorithms give pretty similar results. However, if  $r$  or  $\sigma$  is large, the effects of entry-wise perturbations are pronounced in the kernel matrix, and therefore, the results are different depending on the robustness of the algorithms. As shown in Figure 3, the three  $L_1$ -norm kernel PCA algorithms and R1-KPCA outperform L2-KPCA as in the case of the linear kernel. However, while R1-KPCA achieves the best performance for the linear kernel, the  $L_1$ -norm based kernel PCA algorithms work better than R1-KPCA when the Gaussian kernel is used. Especially, Algo-

rithm 1 outperforms all the other algorithms if  $r$  exceeds 20. The superior performance of Algorithm 1 ranges from 1% to 5% in these cases.

## 6.2 Outlier Detection

$L_2$ -norm PCA has been shown to be effective for anomaly detection [27]. The idea is to extract loading vectors using datasets consisting of only normal samples and use these loading vectors to develop a detection model. Specifically, a boundary of normal samples is constructed from the loading vectors and the boundary is used to discriminate normal and abnormal samples.

We extend this principle to outlier detection, i.e. its unsupervised counterpart. In the outlier detection setting, sample labels are not given when the model is built. Therefore, it is not possible to build a detection model solely based on normal samples. Given this context, we run robust kernel PCA algorithms on the entire dataset (with outliers) and use the resulting loading vectors to characterize a boundary of normal samples. Since these loading vectors are less influenced by outliers as illustrated in Section 6.1, we expect that they would better construct a normal boundary. We compare the performance of Algorithm 1 based models to that of KPCA-L1, L1-KPCA, R1-KPCA, and L2-KPCA based models as well as two other popular outlier detection models [28] [29].

### 6.2.1 Toy Examples

We first illustrate the advantage of using robust kernel PCA for outlier detection using the following two-dimensional toy examples.

Figure 4 displays the distribution of normal samples and outliers. As the normal samples follow a linear pattern, we run the kernel PCA algorithms with the linear kernel and represent their first loading vectors in Figure 4. In the figure, the first loading vectors of the three  $L_1$ -norm based kernel PCA algorithms are represented using a single dashed line since they yield the same first loading vector in this example. In addition to the normal samples forming a linear pattern, there are some outliers scattered exhibiting two different patterns; the two triangle points are outliers due to their scale and the six square points are outliers since they do not follow the linear pattern. If the first loading vector exactly matches the linear pattern, outliers can be easily detected in the principal space; the triangle points can be detected due to large first principal components and the square points can be detected from large second principal components. However, due to the presence of outliers, it is impossible that the first loading vector exactly matches the linear pattern. Given this context, we use robust kernel PCA algorithms to obtain the first loading vector with lower deviation from the linear pattern.

Figure 5 displays the PCA results of the five kernel PCA algorithms. In the figure, the x-axis and the y-axis represents the first and the second principal component, respectively. As shown in the figure, the triangle outliers can be easily separated by the first principal component for any kernel PCA algorithm. However, while the square outliers can be discriminated by the second principal component of the  $L_1$ -norm based kernel PCA algorithms and R1-KPCA, there

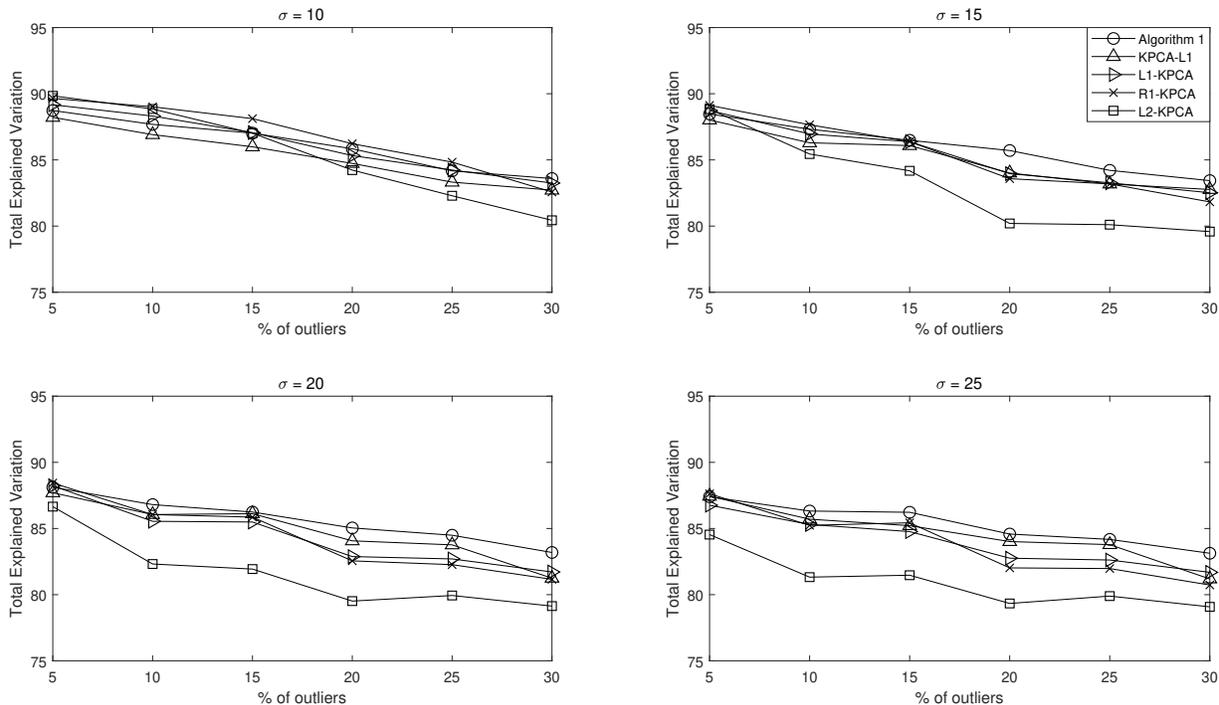


Fig. 3. Robust Extraction of PCs (Gaussian Kernel with  $\sigma$  ranging from 10 to 25)

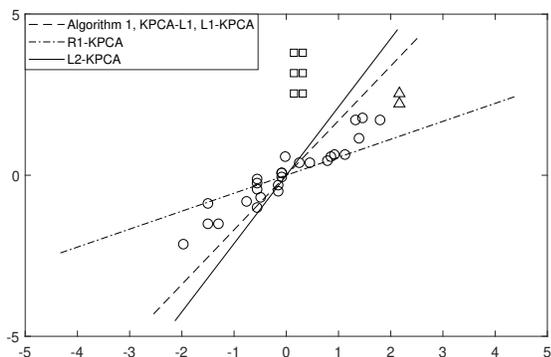


Fig. 4. The first toy example - original space

exists some overlap between the normal samples and the square outliers in the range of the second principal component of L2-KPCA. As seen in the figure, two outliers appear closer to the origin than some normal samples making the circular boundary of the normal samples include them. On the other hand, all the normal samples are clearly separated from the outliers in the principal space of the  $L_1$ -norm based kernel PCA algorithms and R1-KPCA, demonstrating the advantage of using robust kernel PCA in outlier detection. This result is consistent with the findings in Figure 2.

In order to see if the same result holds for the Gaussian kernel, we consider another example. As shown in Figure 6, the second example has a spiral pattern consisting of normal samples as well as two types of outliers. As in the previous example, it has both trivial outliers (the triangle

points) and more challenging outliers (the square points). In order to obtain nonlinear principal components, we run the five kernel PCA algorithms with the Gaussian kernel. As Figure 7 displays, only Algorithm 1 succeeds to exclude the square outliers from the boundary while the other kernel PCA algorithms include them within the boundary. This superior performance of Algorithm 1 with the Gaussian kernel is consistent with the results in Section 6.1 and attests the effectiveness of using it for outlier detection, especially with the Gaussian kernel.

### 6.2.2 Real-world Datasets

For outlier detection, we use datasets from the UCI Machine Learning Repository [30] and the ODDS Library [31], see Table 1.

TABLE 1  
Real-world Datasets for Outlier Detection

Data set	# samples	# features	# outliers
WBC	378	30	21 (7.6%)
Ionosphere	351	33	126 (36%)
BreastW	683	9	239 (35%)
Cardio	1831	21	176 (9.6%)
Musk	3062	166	97 (3.2%)
Mnist	7603	100	700 (9.2%)

In this experiment, we use a similar detection rule as the one in [27] where it is applied for anomaly detection. Let  $Y \in \mathbb{R}^{n \times p}$  denote  $p$  principal components and  $m_j$  and  $\lambda_j$  be the mean and variance of the  $j^{th}$  principal component, respectively. To detect outliers, we consider the following

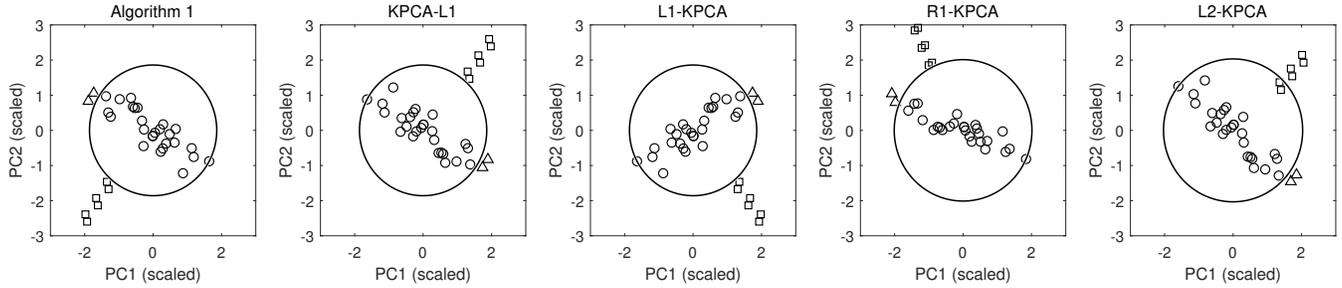


Fig. 5. The first toy example - principal space

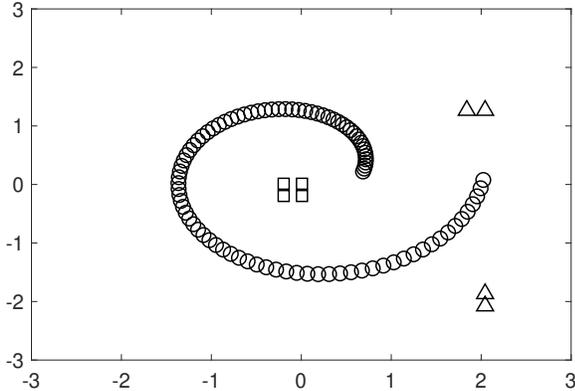


Fig. 6. The second toy example - original space

detection model, which classify the  $i^{th}$  sample as an outlier if

$$\sum_{\{j:\lambda_j \geq \alpha\}} \frac{(Y_{ij} - m_j)^2}{\lambda_j} > c. \quad (28)$$

The metric appearing on the left-hand side of (28) represents the squared Euclidean distance to the origin in the standardized principal space consisting of principal components whose variance is greater than or equal to  $\alpha$ . Therefore, our model can be understood as drawing a circular boundary (as illustrated in Figures 5 and 7) on this reduced standardized principal space. Since sample labels are unknown at the stage of building a model in the outlier detection setting, it is unclear how to choose an appropriate  $c$ . So, we compute precision and recall with varying  $c$  and evaluate the performance of each model using AUC under the precision-recall curve. We compare AUC of the Algorithm 1 based models to that of the KPCA-L1, L1-KPCA, R1-KPCA, and L2-KPCA based models as well as that of the two popular outlier detection models, Local Outlier Factor (LOF) [28] and Isolation Forest (iForest) [29].

Since principal components having small sample variance provide minor information, we only consider principal components whose sample variance is greater than or equal to some threshold value  $\alpha$ . We set  $\alpha$  be to the largest  $\bar{\alpha}$  such that

$$0.8 \times \sum_{j=1}^d \lambda_j \leq \sum_{\{j:\lambda_j \geq \bar{\alpha}\}} \lambda_j$$

holds where  $d$  is the number of features. For the choice of the kernel function, we consider both the linear kernel and the Gaussian kernel with the width parameter  $\sigma$  of the Gaussian kernel to be equal to  $d$ . On the other hand, we set the number of nearest neighbors to 10 in LOF, and the number of trees, the size of subsample, and the number of rounds to 100, 256, and 10, respectively in iForest since these parameter values are commonly used.

Table 2 displays the AUCs of the 12 different detection models. The numbers in bold present the highest AUC cases (there can be several similar top performances). If outliers are obvious, any kernel PCA based model works well as seen in the case of Breastw and Musk. However, if outliers are unclear, the Algorithm 1 based detection models tend to outperform the other detection models. Especially, the Algorithm 1 based model with the Gaussian kernel consistently achieves top AUC values. Compared to the kernel PCA based models, LOF and iForest do not work well. LOF never achieves the top performance and iForest is not competitive for high-dimensional datasets such as Must and MNIST although it yields the top AUC values for WBC and Breastw. As opposed to them, the Algorithm 1 based model with the Gaussian kernel consistently works well regardless of the size of the problem, demonstrating its effectiveness in outlier detection.

### 6.3 Runtime Comparison

Lastly, we compare the runtime of Algorithm 1 to that of KPCA-L1, L1-KPCA, R1-KPCA, and L2-KPCA. In order to obtain a runtime comparison, we run them on the six real-world datasets presented in Table 1 and measure the time taken to get all the principal components.

As shown in Table 3, the runtime largely varies across the algorithms. Among the  $L_1$ -norm based kernel PCA algorithms, Algorithm 1 has the smallest runtime for all datasets. Actually, it is much faster than the other two algorithms since it requires only one matrix-vector multiplication while the other algorithms entail either eigen-decomposition or solving a system of equations. R1-KPCA is also not as fast as Algorithm 1 since it involves QR-decomposition in each iteration to make loading vectors orthogonal. Among the robust kernel PCA algorithms, only Algorithm 1 is computationally comparable to L2-KPCA, making it the best choice for robust kernel PCA in a large-scale setting.

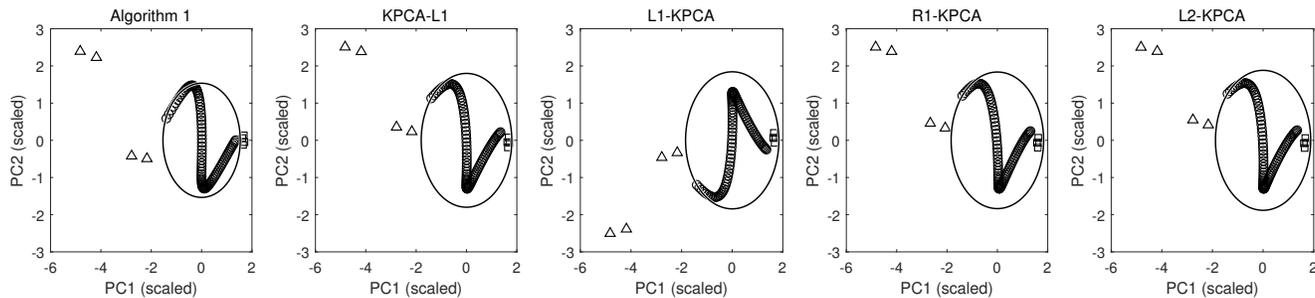


Fig. 7. The second toy example - principal space

TABLE 2  
AUC of the Outlier Detection Models

Datasets	AUC											
	Linear					Gaussian					LOF	iForest
	Algo 1	KPCA -L1	L1- KPCA	R1- KPCA	L2- KPCA	Algo 1	KPCA -L1	L1- KPCA	R1- KPCA	L2- KPCA		
WBC	0.5208	0.5288	0.5320	0.4658	0.4798	0.5292	<b>0.5340</b>	<b>0.5337</b>	0.5072	0.5224	0.3451	<b>0.5525</b>
Ionosphere	0.6625	<b>0.7319</b>	0.6834	<b>0.7642</b>	0.7057	<b>0.7238</b>	0.6806	0.6887	0.7041	0.6992	0.7032	0.7067
Breastw	0.9250	0.9125	0.9218	0.9269	0.9152	<b>0.9428</b>	0.9287	0.9354	<b>0.9521</b>	0.9309	0.3750	<b>0.9513</b>
Cardio	0.5790	0.5551	<b>0.5799</b>	0.4265	0.5066	<b>0.6096</b>	0.5752	<b>0.5963</b>	0.5116	0.4664	0.1921	0.5114
Musk	<b>0.9947</b>	<b>0.9947</b>	<b>0.9947</b>	0.8055	0.9358	<b>0.9947</b>	<b>0.9947</b>	<b>0.9947</b>	0.9916	<b>0.9947</b>	0.0925	0.7596
MNIST	<b>0.3985</b>	<b>0.4002</b>	N/A	N/A	0.3914	<b>0.3966</b>	0.3913	N/A	N/A	0.3639	0.1924	0.3380

\* N/A: The experiments can not be completed within the period of 24 hours.

TABLE 3  
Runtime Comparison

Datasets	Runtime (minutes)									
	Linear					Gaussian				
	Algo 1	KPCA -L1	L1- KPCA	R1- KPCA	L2- KPCA	Algo 1	KPCA -L1	L1- KPCA	R1- KPCA	L2- KPCA
WBC	0.0	0.0	0.1	0.2	0.0	0.0	0.0	0.0	0.2	0.0
Ionosphere	0.0	0.0	0.1	0.2	0.0	0.0	0.0	0.1	0.2	0.0
Breastw	0.0	0.0	0.1	0.2	0.0	0.0	0.0	0.1	0.1	0.0
Cardio	0.0	1.9	12.3	7.6	0.1	0.0	1.8	15.3	6.2	0.1
Musk	0.5	69.6	999.1	973.5	0.3	0.5	12.9	1018.0	968.3	0.1
MNIST	0.9	263.3	> 1440	> 1440	1.8	1.0	263.2	> 1440	> 1440	1.8

## 7 CONCLUSION

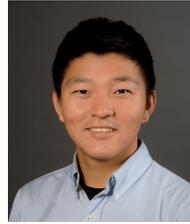
In this work, we present a simple algorithm for  $L_1$ -norm kernel PCA and provide its convergence analysis. In order to develop it, we first reformulate  $L_1$ -norm kernel PCA into a geometrically interpretable problem and derive a geometric interpretation behind it. Based on the geometric interpretation, we develop an algorithm to which the kernel trick is applicable. In the convergence analysis, we prove that the algorithm converges to a local optimal solution in a finite number of steps and the sequence of objective values converges at a linear rate.

The computational experiments demonstrate the robustness of the proposed algorithm in the presence of entry-wise perturbations and the runtime comparison shows that it takes much less time than the other robust kernel PCA algorithms. Also, its application to outlier detection outperforms all of the other benchmark algorithms. The model based on the proposed algorithm is not only better than that of the other kernel PCA based models but also outperforms LOF and iForest, especially when high-dimensional datasets are considered.

## REFERENCES

- [1] I. Jolliffe, *Principal Component Analysis*. Wiley Online Library, 2002.
- [2] B. Schölkopf, A. Smola, and K.-R. Müller, "Kernel Principal Component Analysis," in *International Conference on Artificial Neural Networks*, 1997, pp. 583–588.
- [3] G. H. Golub and C. F. Van Loan, *Matrix Computations*. JHU Press, 2012, vol. 3.
- [4] E. J. Candès and B. Recht, "Exact Matrix Completion via Convex Optimization," *Foundations of Computational Mathematics*, vol. 9, no. 6, p. 717, 2009.
- [5] E. J. Candès, X. Li, Y. Ma, and J. Wright, "Robust Principal Component Analysis?" *Journal of the ACM*, vol. 58, no. 3, p. 11, 2011.
- [6] H. Xu, C. Caramanis, and S. Sanghavi, "Robust PCA via outlier pursuit," in *Advances in Neural Information Processing Systems*, 2010, pp. 2496–2504.
- [7] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma, "Robust Recovery of Subspace Structures by Low-Rank Representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 171–184, 2013.
- [8] G. Liu and P. Li, "Recovery of Coherent Data via Low-Rank Dictionary Pursuit," in *Advances in Neural Information Processing Systems*, 2014, pp. 1206–1214.

- [9] G. Liu, H. Xu, J. Tang, Q. Liu, and S. Yan, "A Deterministic Analysis for LRR," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 3, pp. 417–430, 2016.
- [10] G. Liu, Q. Liu, and P. Li, "Blessing of Dimensionality: Recovering Mixture Data via Dictionary Pursuit," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 1, pp. 47–60, 2017.
- [11] J. P. Brooks, J. H. Dulá, and E. L. Boone, "A Pure  $L_1$ -norm Principal Component Analysis," *Computational Statistics & Data Analysis*, vol. 61, pp. 83–98, 2013.
- [12] Y. W. Park, "Optimization for Regression, PCA, and SVM: Optimality and Scalability," Ph.D. dissertation, Northwestern University, 2015.
- [13] Y. W. Park and D. Klabjan, "Iteratively Reweighted Least Squares Algorithms for  $L_1$ -Norm Principal Component Analysis," in *IEEE International Conference on Data Mining*, 2016, pp. 430–438.
- [14] —, "Three iteratively reweighted least squares algorithms for  $L_1$ -norm principal component analysis," *Knowledge and Information Systems*, vol. 54, no. 3, pp. 541–565, 2018.
- [15] F. Nie, H. Huang, C. Ding, D. Luo, and H. Wang, "Robust Principal Component Analysis with Non-Greedy  $L_1$ -Norm Maximization," in *Proceedings of the International Joint Conference on Artificial Intelligence*, vol. 22, no. 1, 2011, pp. 1433–1438.
- [16] M. McCoy and J. A. Tropp, "Two Proposals for Robust PCA using Semidefinite Programming," *Electronic Journal of Statistics*, vol. 5, pp. 1123–1160, 2011.
- [17] P. P. Markopoulos, G. N. Karystinos, and D. A. Pados, "Optimal Algorithms for  $L_1$ -subspace Signal Processing," *IEEE Transactions on Signal Processing*, vol. 62, no. 19, pp. 5046–5058, 2014.
- [18] Q. Ke and T. Kanade, "Robust  $L_1$  Norm Factorization in the Presence of Outliers and Missing Data by Alternative Convex Programming," in *Proceedings of the 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2005, pp. 739–746.
- [19] C. Ding, D. Zhou, X. He, and H. Zha, " $R_1$ -PCA: Rotational Invariant  $L_1$ -norm Principal Component Analysis for Robust Subspace Factorization," in *Proceedings of the 23rd International Conference on Machine Learning*. ACM, 2006, pp. 281–288.
- [20] G. N. Karystinos and A. P. Liavas, "Efficient Computation of the Binary Vector That Maximizes a Rank-Deficient Quadratic Form," *IEEE Transactions on Information Theory*, vol. 56, no. 7, pp. 3581–3593, 2010.
- [21] Y. Nesterov, "Semidefinite Relaxation and Nonconvex Quadratic Optimization," *Optimization Methods and Software*, vol. 9, no. 1-3, pp. 141–160, 1998.
- [22] N. Kwak, "Principal Component Analysis based on  $L_1$ -norm Maximization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, no. 9, pp. 1672–1680, 2008.
- [23] —, "Nonlinear projection trick in kernel methods: An alternative to the kernel trick," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 24, no. 12, pp. 2113–2119, 2013.
- [24] Y. Xiao, H. Wang, W. Xu, and J. Zhou, "L1 norm based KPCA for novelty detection," *Pattern Recognition*, vol. 46, no. 1, pp. 389–396, 2013.
- [25] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [26] M. X. Goemans and D. P. Williamson, "Improved Approximation Algorithms for Maximum Cut and Satisfiability Problems Using Semidefinite Programming," *Journal of the ACM (JACM)*, vol. 42, no. 6, pp. 1115–1145, 1995.
- [27] M.-L. Shyu, S.-C. Chen, K. Sarinapakorn, and L. Chang, "A Novel Anomaly Detection Scheme Based on Principal Component Classifier," in *IEEE International Conference on Data Mining*, 2003, pp. 172–179.
- [28] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander, "LOF: Identifying Density-Based Local Outliers," in *ACM SIGMOD Record*, vol. 29, no. 2, 2000, pp. 93–104.
- [29] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation Forest," in *IEEE International Conference on Data Mining*, 2008, pp. 413–422.
- [30] M. Lichman, "UCI Machine Learning Repository," 2013. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [31] S. Rayana, "ODDS Library," 2016. [Online]. Available: <http://odds.cs.stonybrook.edu>



**Cheolmin Kim** is a Ph.D student in Industrial Engineering and Management Science at Northwestern University. He received his B.S in Industrial Engineering and B.A in Economics from Seoul National University. His research interests are at the interaction of optimization and machine learning. He is interested in designing a new machine learning model, developing a training algorithm and analyzing the efficiency of the algorithm from an optimization perspective.



**Diego Klabjan** is a professor at Northwestern University, Department of Industrial Engineering and Management Sciences. He is also Founding Director, Master of Science in Analytics. After obtaining his doctorate from the School of Industrial and Systems Engineering of the Georgia Institute of Technology in 1999 in Algorithms, Combinatorics, and Optimization, in the same year he joined the University of Illinois at Urbana-Champaign. In 2007 he became an associate professor at Northwestern and in 2012 he was promoted to a full professor. His research is focused on machine learning, deep learning and analytics with concentration in finance, transportation, sport, and bioinformatics. Professor Klabjan has led projects with large companies such as Intel, Baxter, Allstate, AbbVie, FedEx Express, General Motors, United Continental, and many others, and he is also assisting numerous start-ups with their analytics needs. He is also a founder of Opex Analytics LLC.