# Robust PCA as Bilinear Decomposition with Outlier-Sparsity Regularization†

*Gonzalo Mateos and Georgios B. Giannakis (contact author)\**

## Abstract

Principal component analysis (PCA) is widely used for dimensionality reduction, with well-documented merits in various applications involving high-dimensional data, including computer vision, preference measurement, and bioinformatics. In this context, the fresh look advocated here permeates benefits from variable selection and compressive sampling, to robustify PCA against outliers. A least-trimmed squares estimator of a low-rank bilinear factor analysis model is shown closely related to that obtained from an $\ell_0$-(pseudo)norm-regularized criterion encouraging *sparsity* in a matrix explicitly modeling the outliers. This connection suggests robust PCA schemes based on convex relaxation, which lead naturally to a family of robust estimators encompassing Huber's optimal M-class as a special case. Outliers are identified by tuning a regularization parameter, which amounts to controlling sparsity of the outlier matrix along the whole *robustification* path of (group) least-absolute shrinkage and selection operator (Lasso) solutions. Beyond its neat ties to robust statistics, the developed outlier-aware PCA framework is versatile to accommodate novel and scalable algorithms to: i) track the low-rank signal subspace robustly, as new data are acquired in real time; and ii) determine principal components robustly in (possibly) infinite-dimensional feature spaces. Synthetic and real data tests corroborate the effectiveness of the proposed robust PCA schemes, when used to identify aberrant responses in personality assessment surveys, as well as unveil communities in social networks, and intruders from video surveillance data.

## Index Terms

Robust statistics, principal component analysis, outlier rejection, sparsity, (group) Lasso.

## EDICS Category: MLR-LEAR

∗ The authors are with the Dept. of Electrical and Computer Engineering, University of Minnesota, 200 Union Street SE, Minneapolis, MN 55455. Tel/fax: (612)626-7781/625-2002; Emails: {mate0058,georgios}@umn.edu

## I. INTRODUCTION

Principal component analysis (PCA) is the workhorse of high-dimensional data analysis and dimensionality reduction, with numerous applications in statistics, engineering, and the biobehavioral sciences; see, e.g., [22]. Nowadays ubiquitous e-commerce sites, the Web, and urban traffic surveillance systems generate massive volumes of data. As a result, the problem of extracting the most informative, yet low-dimensional structure from high-dimensional datasets is of paramount importance [17]. To this end, PCA provides least-squares (LS) optimal linear approximants in $\mathbb{R}^q$ to a data set in $\mathbb{R}^p$, for $q \leq p$. The desired linear subspace is obtained from the $q$ dominant eigenvectors of the sample data covariance matrix [22].

Data obeying postulated low-rank models include also outliers, which are samples not adhering to those nominal models. Unfortunately, LS is known to be very sensitive to outliers [19], [32], and this undesirable property is inherited by PCA as well [22]. Early efforts to robustify PCA have relied on robust estimates of the data covariance matrix; see, e.g., [4]. Related approaches are driven from statistical physics [39], and also from M-estimators [8]. Recently, polynomial-time algorithms with remarkable performance guarantees have emerged for low-rank matrix recovery in the presence of sparse – but otherwise arbitrarily large – errors [5], [7]. This pertains to an 'idealized robust' PCA setup, since those entries not affected by outliers are assumed error free. Stability in reconstructing the low-rank and sparse matrix components in the presence of 'dense' noise have been reported in [38], [42]. A hierarchical Bayesian model was proposed to tackle the aforementioned low-rank plus sparse matrix decomposition problem in [9].

In the present paper, a robust PCA approach is pursued requiring minimal assumptions on the outlier model. A natural least-trimmed squares (LTS) PCA estimator is first shown closely related to an estimator obtained from an $\ell_0$-(pseudo)norm-regularized criterion, adopted to fit a low-rank bilinear factor analysis model that explicitly incorporates an unknown *sparse* vector of outliers per datum (Section II). As in compressive sampling [35], efficient (approximate) solvers are obtained in Section III, by surrogating the $\ell_0$-norm of the outlier matrix with its closest convex approximant. This leads naturally to an M-type PCA estimator, which subsumes Huber's optimal choice as a special case [13]. Unlike Huber's formulation though, results here are not confined to an outlier contamination model. A tunable parameter controls the sparsity of the estimated matrix, and the number of outliers as a byproduct. Hence, effective data-driven methods to select this parameter are of paramount importance, and systematic approaches are pursued by efficiently exploring the entire *robustifaction* (a.k.a. homotopy) path of (group-) Lasso solutions [17], [41]. In this sense, the method here capitalizes on but *is not limited to* sparse settings where outliers are sporadic, since one can examine all sparsity levels along the robustification path. The outlier-aware generative data model and its sparsity-controlling estimator are quite general, since minor modifications

discussed in Section III-C enable robustifiying linear regression [14], dictionary learning [24], [34], and K-means clustering as well [12], [17]. Section IV deals with further modifications for bias reduction through nonconvex regularization, and automatic determination of the reduced dimension $q$.

Beyond its neat ties to robust statistics, the developed outlier-aware PCA framework is versatile to accommodate scalable *robust* algorithms to: i) track the low-rank signal subspace, as new data are acquired in real time (Section V); and ii) determine principal components in (possibly) infinite-dimensional feature spaces, thus robustifying kernel PCA [33], and spectral clustering as well [17, p. 544] (Section VI). The vast literature on *non-robust* subspace tracking algorithms includes [24], [40], and [2]; see also [18] for a first-order algorithm that is robust to outliers and incomplete data. Relative to [18], the online robust (OR-) PCA algorithm of this paper is a second-order method, which minimizes an outlier-aware exponentially-weighted LS estimator of the low-rank factor analysis model. Since the outlier and subspace estimation tasks decouple nicely in OR-PCA, one can readily devise a first-order counterpart when minimal computational loads are at a premium. In terms of performance, online algorithms are known to be markedly faster than their batch alternatives [2], [18], e.g., in the timely context of low-rank matrix completion [29], [30]. While the focus here is not on incomplete data records, extensions to account for missing data are immediate and will be reported elsewhere.

In Section VII, numerical tests with synthetic and real data corroborate the effectiveness of the proposed robust PCA schemes, when used to identify aberrant responses from a questionnaire designed to measure the Big-Five dimensions of personality traits [21], as well as unveil communities in a (social) network of college football teams [15], and intruders from video surveillance data [8]. Concluding remarks are given in Section VIII, while a few technical details are deferred to the Appendix.

*Notation:* Bold uppercase (lowercase) letters will denote matrices (column vectors). Operators $(\cdot)'$, $\text{tr}(\cdot)$, $\text{med}(\cdot)$, and $\odot$ will denote transposition, matrix trace, median, and Hadamard product, respectively. Vector $\text{diag}(\mathbf{M})$ collects the diagonal entries of $\mathbf{M}$, whereas the diagonal matrix $\text{diag}(\mathbf{v})$ has the entries of $\mathbf{v}$ on its diagonal. The $\ell_p$-norm of $\mathbf{x} \in \mathbb{R}^n$ is $\|\mathbf{x}\|_p := \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$ for $p \geq 1$; and $\|\mathbf{M}\|_F := \sqrt{\text{tr}(\mathbf{M}\mathbf{M}')}$ is the matrix Frobenious norm. The $n \times n$ identity matrix will be represented by $\mathbf{I}_n$, while $\mathbf{0}_n$ will denote the $n \times 1$ vector of all zeros, and $\mathbf{0}_{n \times m} := \mathbf{0}_n \mathbf{0}'_m$. Similar notation will be adopted for vectors (matrices) of all ones. The $i$-th vector of the canonical basis in $\mathbb{R}^n$ will be denoted by $\mathbf{b}_{n,i}$, $i = 1, \ldots, n$.

## II. ROBUSTIFYING PCA

Consider the standard PCA formulation, in which a set of data $\mathcal{T}_y := \{\mathbf{y}_n\}_{n=1}^N$ in the $p$-dimensional Euclidean *input* space is given, and the goal is to find the best $q$-rank ($q \leq p$) linear approximation to the

data in $\mathcal{T}_y$; see e.g., [22]. Unless otherwise stated, it is assumed throughout that the value of $q$ is given. One approach to solving this problem, is to adopt a low-rank bilinear (factor analysis) model

$$\mathbf{y}_n = \mathbf{m} + \mathbf{U}\mathbf{s}_n + \mathbf{e}_n, \quad n = 1, \dots, N \tag{1}$$

where $\mathbf{m} \in \mathbb{R}^p$ is a location (mean) vector; matrix $\mathbf{U} \in \mathbb{R}^{p \times q}$ has orthonormal columns spanning the signal subspace; $\{\mathbf{s}_n\}_{n=1}^N$ are the so-termed *principal components*, and $\{\mathbf{e}_n\}_{n=1}^N$ are zero-mean i.i.d. random errors. The unknown variables in (1) can be collected in $\mathcal{V} := \{\mathbf{m}, \mathbf{U}, \{\mathbf{s}_n\}_{n=1}^N\}$, and they are estimated using the LS criterion as

$$\min_{\mathcal{V}} \sum_{n=1}^N \|\mathbf{y}_n - \mathbf{m} - \mathbf{U}\mathbf{s}_n\|_2^2, \quad \text{s. to} \quad \mathbf{U}'\mathbf{U} = \mathbf{I}_q. \tag{2}$$

PCA in (2) is a nonconvex optimization problem due to the bilinear terms $\mathbf{U}\mathbf{s}_n$, yet a global optimum $\hat{\mathcal{V}}$ can be shown to exist; see e.g., [40]. The resulting estimates are $\hat{\mathbf{m}} = \sum_{n=1}^N \mathbf{y}_n/N$ and $\hat{\mathbf{s}}_n = \hat{\mathbf{U}}'(\mathbf{y}_n - \hat{\mathbf{m}})$, $n = 1, \dots, N$; while $\hat{\mathbf{U}}$ is formed with columns equal to the $q$-dominant right singular vectors of the $N \times p$ data matrix $\mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_N]'$ [17, p. 535]. The principal components (entries of) $\mathbf{s}_n$ are the projections of the centered data points $\{\mathbf{y}_n - \hat{\mathbf{m}}\}_{n=1}^N$ onto the signal subspace. Equivalently, PCA can be formulated based on maximum variance, or, minimum reconstruction error criteria; see e.g., [22].

## A. Least-trimmed squares PCA

Given training data $\mathcal{T}_x := \{\mathbf{x}_n\}_{n=1}^N$ possibly contaminated with outliers, the goal here is to develop a robust estimator of $\mathcal{V}$ that requires minimal assumptions on the outlier model. Note that there is an explicit notational differentiation between: i) the data in $\mathcal{T}_y$ which adhere to the nominal model (1); and ii) the given data in $\mathcal{T}_x$ that may also contain outliers, i.e., those $\mathbf{x}_n$ not adhering to (1). Building on LTS regression [32], the desired robust estimate $\hat{\mathcal{V}}_{LTS} := \{\hat{\mathbf{m}}, \hat{\mathbf{U}}, \{\hat{\mathbf{s}}_n\}_{n=1}^N\}$ for a prescribed $\nu < N$ can be obtained via the following LTS PCA estimator [cf. (2)]

$$\hat{\mathcal{V}}_{LTS} := \arg \min_{\mathcal{V}} \sum_{n=1}^\nu r_{[n]}^2(\mathcal{V}), \quad \text{s. to} \quad \mathbf{U}'\mathbf{U} = \mathbf{I}_q \tag{3}$$

where $r_{[n]}^2(\mathcal{V})$ is the $n$-th order statistic among the squared residual norms $r_1^2(\mathcal{V}), \dots, r_N^2(\mathcal{V})$, and $r_n(\mathcal{V}) := \|\mathbf{x}_n - \mathbf{m} - \mathbf{U}\mathbf{s}_n\|_2$. The so-termed *coverage* $\nu$ determines the breakdown point of the LTS PCA estimator [32], since the $N - \nu$ largest residuals are absent from the estimation criterion in (3). Beyond this universal outlier-rejection property, the LTS-based estimation offers an attractive alternative to robust linear regression due to its high breakdown point and desirable analytical properties, namely $\sqrt{N}$-consistency and asymptotic normality under mild assumptions [32].

**Remark 1 (Robust estimation of the mean):** In most applications of PCA, data in $\mathcal{T}_y$ are typically assumed zero mean. This is without loss of generality, since nonzero-mean training data can always be rendered zero mean, by subtracting the sample mean $\sum_{n=1}^{N} \mathbf{y}_n / N$ from each $\mathbf{y}_n$. In modeling zero-mean data, the known vector $\mathbf{m}$ in (1) can obviously be neglected. When outliers are present however, data in $\mathcal{T}_x$ are not necessarily zero mean, and it is unwise to center them using the non-robust sample mean estimator which has a breakdown point equal to zero [32]. Towards robustifying PCA, a more sensible approach is to estimate $\mathbf{m}$ robustly, and jointly with $\mathbf{U}$ and the principal components $\{\mathbf{s}_n\}_{n=1}^{N}$.

Because (3) is a nonconvex optimization problem, a nontrivial issue pertains to the existence of the proposed LTS PCA estimator, i.e., whether or not (3) attains a minimum. Fortunately, the answer is in the affirmative as asserted next.

**Property 1:** *The LTS PCA estimator is well defined, since* (3) *has (at least) one solution.*

Existence of $\hat{\mathcal{V}}_{LTS}$ can be readily established as follows: i) for each subset of $\mathcal{T}$ with cardinality $\nu$ (there are $\binom{N}{\nu}$ such subsets), solve the corresponding PCA problem to obtain a unique candidate estimator per subset; and ii) pick $\hat{\mathcal{V}}_{LTS}$ as the one among all $\binom{N}{\nu}$ candidates with the minimum cost.

Albeit conceptually simple, the solution procedure outlined under Property 1 is combinatorially complex, and thus intractable except for small sample sizes $N$. Algorithms to obtain approximate LTS solutions in large-scale linear regression problems are available; see e.g., [32].

### B. $\ell_0$-norm regularization for robustness

Instead of discarding large residuals, the alternative approach here explicitly accounts for outliers in the low-rank data model (1). This becomes possible through the vector variables $\{\mathbf{o}_n\}_{n=1}^{N}$ one per training datum $\mathbf{x}_n$, which take the value $\mathbf{o}_n \neq \mathbf{0}_p$ whenever datum $n$ is an outlier, and $\mathbf{o}_n = \mathbf{0}_p$ otherwise. Thus, the novel outlier-aware factor analysis model is

$$\mathbf{x}_n = \mathbf{y}_n + \mathbf{o}_n = \mathbf{m} + \mathbf{U}\mathbf{s}_n + \mathbf{e}_n + \mathbf{o}_n, \qquad n = 1, \ldots, N \tag{4}$$

where $\mathbf{o}_n$ can be deterministic or random with unspecified distribution. In the *under-determined* linear system of equations (4), both $\mathcal{V}$ as well as the $N \times p$ matrix $\mathbf{O} := [\mathbf{o}_1, \ldots, \mathbf{o}_N]'$ are unknown. The percentage of outliers dictates the degree of *sparsity* (number of zero rows) in $\mathbf{O}$. Sparsity control will prove instrumental in efficiently estimating $\mathbf{O}$, rejecting outliers as a byproduct, and consequently arriving at a *robust* estimator of $\mathcal{V}$. To this end, a natural criterion for controlling outlier sparsity is to seek the estimator [cf. (2)]

$$\{\hat{\mathcal{V}}, \hat{\mathbf{O}}\} = \arg \min_{\mathcal{V}, \mathbf{O}} \|\mathbf{X} - \mathbf{1}_N \mathbf{m}' - \mathbf{S}\mathbf{U}' - \mathbf{O}\|_F^2 + \lambda_0 \|\mathbf{O}\|_0, \quad \text{s. to } \mathbf{U}'\mathbf{U} = \mathbf{I}_q \tag{5}$$

where $\mathbf{X} := [\mathbf{x}_1, \ldots, \mathbf{x}_N]' \in \mathbb{R}^{N \times p}$, $\mathbf{S} := [\mathbf{s}_1, \ldots, \mathbf{s}_N]' \in \mathbb{R}^{N \times q}$, and $\|\mathbf{O}\|_0$ denotes the nonconvex $\ell_0$-norm that is equal to the number of nonzero rows of $\mathbf{O}$. Vector (group) sparsity in the rows $\hat{\mathbf{o}}_n$ of $\hat{\mathbf{O}}$ can be directly controlled by tuning the parameter $\lambda_0 \geq 0$.

As with compressive sampling and sparse modeling schemes that rely on the $\ell_0$-norm [35], the robust PCA problem (5) is NP-hard [26]. In addition, the sparsity-controlling estimator (5) is intimately related to LTS PCA, as asserted next.

**Proposition 1:** *If $\{\hat{\mathcal{V}}, \hat{\mathbf{O}}\}$ minimizes* (5) *with $\lambda_0$ chosen such that $\|\hat{\mathbf{O}}\|_0 = N - \nu$, then $\hat{\mathcal{V}}_{LTS} = \hat{\mathcal{V}}$.*

*Proof:* Given $\lambda_0$ such that $\|\hat{\mathbf{O}}\|_0 = N - \nu$, the goal is to characterize $\hat{\mathcal{V}}$ as well as the positions and values of the nonzero rows of $\hat{\mathbf{O}}$. Note that because $\|\hat{\mathbf{O}}\|_0 = N - \nu$, the last term in the cost of (5) is constant, hence inconsequential to the minimization. Upon defining $\hat{\mathbf{r}}_n := \mathbf{x}_n - \hat{\mathbf{m}} - \hat{\mathbf{U}}\hat{\mathbf{s}}_n$, it is not hard to see from the optimality conditions that the rows of $\hat{\mathbf{O}}$ satisfy

$$\hat{\mathbf{o}}_n = \begin{cases} \mathbf{0}_p, & \|\hat{\mathbf{r}}_n\|_2 \leq \sqrt{\lambda_0} \\ \hat{\mathbf{r}}_n, & \|\hat{\mathbf{r}}_n\|_2 > \sqrt{\lambda_0} \end{cases}, \quad n = 1, \ldots, N. \tag{6}$$

This is intuitive, since for those nonzero $\hat{\mathbf{o}}_n$ the best thing to do in terms of minimizing the overall cost is to set $\hat{\mathbf{o}}_n = \hat{\mathbf{r}}_n$, and thus null the corresponding squared-residual terms in (5). In conclusion, for the chosen value of $\lambda_0$ it holds that $N - \nu$ squared residuals effectively do not contribute to the cost in (5).

To determine $\hat{\mathcal{V}}$ and the row support of $\hat{\mathbf{O}}$, one alternative is to exhaustively test all $\binom{N}{N-\nu} = \binom{N}{\nu}$ admissible row-support combinations. For each one of these combinations (indexed by $j$), let $\mathcal{S}_j \subset \{1, \ldots, N\}$ be the index set describing the row support of $\hat{\mathbf{O}}^{(j)}$, i.e., $\hat{\mathbf{o}}_n^{(j)} \neq \mathbf{0}_p$ if and only if $n \in \mathcal{S}_j$; and $|\mathcal{S}_j| = N - \nu$. By virtue of (6), the corresponding candidate $\hat{\mathcal{V}}^{(j)}$ solves $\min_{\mathcal{V}} \sum_{n \in \mathcal{S}_j} r_n^2(\mathcal{V})$ subject to $\mathbf{U}'\mathbf{U} = \mathbf{I}_q$, while $\hat{\mathcal{V}}$ is the one among all $\{\hat{\mathcal{V}}^{(j)}\}$ that yields the least cost. Recognizing the aforementioned solution procedure as the one for LTS PCA outlined under Property 1, it follows that $\hat{\mathcal{V}}_{LTS} = \hat{\mathcal{V}}$. ∎

The importance of Proposition 1 is threefold. First, it formally justifies model (4) and its estimator (5) for robust PCA, in light of the well documented merits of LTS [32]. Second, it further solidifies the connection between sparsity-aware learning and robust estimation. Third, problem (5) lends itself naturally to efficient (approximate) solvers based on convex relaxation, the subject dealt with next.

## III. SPARSITY-CONTROLLING OUTLIER REJECTION

Recall that the row-wise $\ell_2$-norm sum $\|\mathbf{B}\|_{2,r} := \sum_{n=1}^N \|\mathbf{b}_n\|_2$ of matrix $\mathbf{B} := [\mathbf{b}_1, \ldots, \mathbf{b}_N]' \in \mathbb{R}^{N \times p}$ is the closest convex approximation of $\|\mathbf{B}\|_0$. This property motivates relaxing problem (5) to

$$\min_{\mathcal{V}, \mathbf{O}} \|\mathbf{X} - \mathbf{1}_N \mathbf{m}' - \mathbf{S}\mathbf{U}' - \mathbf{O}\|_F^2 + \lambda_2 \|\mathbf{O}\|_{2,r}, \quad \text{s. to } \mathbf{U}'\mathbf{U} = \mathbf{I}_q. \tag{7}$$

The nondifferentiable $\ell_2$-norm regularization term encourages row-wise (vector) sparsity on the estimator of $\mathbf{O}$, a property that has been exploited in diverse problems in engineering, statistics, and machine learning [17]. A noteworthy representative is the group Lasso [41], a popular tool for joint estimation and selection of grouped variables in linear regression.

It is pertinent to ponder on whether problem (7) still has the potential of providing robust estimates $\hat{\mathcal{V}}$ in the presence of outliers. The answer is positive, since it is shown in the Appendix that (7) is equivalent to an M-type estimator

$$\min_{\mathcal{V}} \sum_{n=1}^{N} \rho_v(\mathbf{x}_n - \mathbf{m} - \mathbf{U}\mathbf{s}_n), \quad \text{s. to } \mathbf{U}'\mathbf{U} = \mathbf{I}_q \tag{8}$$

where $\rho_v : \mathbb{R}^p \to \mathbb{R}$ is a vector extension to Huber's convex loss function [19]; see also [23], and

$$\rho_v(\mathbf{r}) := \begin{cases} \|\mathbf{r}\|_2^2, & \|\mathbf{r}\|_2 \leq \lambda_2/2 \\ \lambda_2\|\mathbf{r}\|_2 - \lambda_2^2/4, & \|\mathbf{r}\|_2 > \lambda_2/2 \end{cases}. \tag{9}$$

M-type estimators (including Huber's) adopt a fortiori an $\epsilon$-contaminated probability distribution for the outliers, and rely on minimizing the *asymptotic* variance of the resultant estimator for the least favorable distribution of the $\epsilon$-contaminated class (asymptotic min-max approach) [19]. The assumed degree of contamination specifies the tuning parameter $\lambda_2$ in (9) (and thus the threshold for deciding the outliers in M-estimators). In contrast, the present approach is universal in the sense that it is not confined to any assumed class of outlier distributions, and can afford a data-driven selection of the tuning parameter. In a nutshell, M-estimators can be viewed as a special case of the present formulation only for a specific choice of $\lambda_2$, which is not obtained via a data-driven approach, but from distributional assumptions instead.

All in all, the sparsity-controlling role of the tuning parameter $\lambda_2 \geq 0$ in (7) is central, since model (4) and the equivalence of (7) with (8) suggest that $\lambda_2$ is a robustness-controlling constant. Data-driven approaches to select $\lambda_2$ are described in detail under Section III-B. Before dwelling into algorithmic issues to solve (7), a couple of remarks are in order.

**Remark 2 ($\ell_1$-norm regularization for entry-wise outliers):** In computer vision applications where robust PCA schemes are particularly attractive, one may not wish to discard the entire (vectorized) images $\mathbf{x}_n$, but only specific pixels deemed as outliers [8]. This can be accomplished by replacing $\|\mathbf{O}\|_{2,r}$ in (7) with $\|\mathbf{O}\|_1 := \sum_{n=1}^{N} \|\mathbf{o}_n\|_1$, a Lasso-type regularization that encourages entry-wise sparsity in $\hat{\mathbf{O}}$.

**Remark 3 (Outlier rejection):** From the equivalence between problems (7) and (8), it follows that those data points $\mathbf{x}_n$ deemed as containing outliers ($\hat{\mathbf{o}}_n \neq \mathbf{0}_p$) are not completely discarded from the estimation process. Instead, their effect is downweighted as per Huber's loss function [cf. (9)]. Nevertheless, explicitly accounting for the outliers in $\hat{\mathbf{O}}$ provides the means of identifying and removing the contaminated data altogether, and thus possibly re-running PCA on the outlier-free data.

*A. Solving the relaxed problem*

To optimize (7) iteratively for a given value of $\lambda_2$, an alternating minimization (AM) algorithm is adopted which cyclically updates $\mathbf{m}(k) \rightarrow \mathbf{S}(k) \rightarrow \mathbf{U}(k) \rightarrow \mathbf{O}(k)$ per iteration $k = 1, 2, \ldots$. AM algorithms are also known as block-coordinate-descent methods in the optimization parlance; see e.g., [3], [36]. To update each of the variable groups, (7) is minimized while fixing the rest of the variables to their most up-to-date values. While the overall problem (7) is not jointly convex with respect to (w.r.t.) $\{\mathbf{S}, \mathbf{U}, \mathbf{O}, \mathbf{m}\}$, fixing all but one of the variable groups yields subproblems that are efficiently solved, and attain a unique solution.

Towards deriving the updates at iteration $k$ and arriving at the desired algorithm, note first that the mean update is $\mathbf{m}(k) = (\mathbf{X} - \mathbf{O}(k))'\mathbf{1}_N/N$. Next, form the centered and outlier-compensated data matrix $\mathbf{X}_o(k) := \mathbf{X} - \mathbf{1}_N \mathbf{m}(k)' - \mathbf{O}(k-1)$. The principal components are readily given by

$$\mathbf{S}(k) = \arg\min_{\mathbf{S}} \|\mathbf{X}_o(k) - \mathbf{S}\mathbf{U}(k-1)'\|_F^2 = \mathbf{X}_o(k)\mathbf{U}(k-1).$$

Continuing the cycle, $\mathbf{U}(k)$ solves

$$\min_{\mathbf{U}} \|\mathbf{X}_o(k) - \mathbf{S}(k)\mathbf{U}'\|_F^2, \quad \text{s. to } \mathbf{U}'\mathbf{U} = \mathbf{I}_q$$

a constrained LS problem also known as reduced-rank *Procrustes rotation* [43]. The minimizer is given in analytical form in terms of the left and right singular vectors of $\mathbf{X}_o'(k)\mathbf{S}(k)$ [43, Thm. 4]. In detail, one computes the SVD of $\mathbf{X}_o'(k)\mathbf{S}(k) = \mathbf{L}(k)\mathbf{D}(k)\mathbf{R}'(k)$ and updates $\mathbf{U}(k) = \mathbf{L}(k)\mathbf{R}'(k)$. Next, the minimization of (7) w.r.t. $\mathbf{O}$ is an orthonormal group Lasso problem. As such, it decouples across rows $\mathbf{o}_n$ giving rise to $N$ $\ell_2$-norm regularized subproblems, namely

$$\mathbf{o}_n(k) = \arg\min_{\mathbf{o}} \|\mathbf{r}_n(k) - \mathbf{o}\|_2^2 + \lambda_2 \|\mathbf{o}\|_2, \quad n = 1, \ldots, N$$

where $\mathbf{r}_n(k) := \mathbf{x}_n - \mathbf{m}(k) - \mathbf{U}(k)\mathbf{s}_n(k)$. The respective solutions are given by (see e.g., [27])

$$\mathbf{o}_n(k) = \frac{\mathbf{r}_n(k)(\|\mathbf{r}_n(k)\|_2 - \lambda_2/2)_+}{\|\mathbf{r}_n(k)\|_2}, \quad n = 1, \ldots, N \tag{10}$$

where $(\cdot)_+ := \max(\cdot, 0)$. For notational convenience, these $N$ parallel vector soft-thresholded updates are denoted as $\mathbf{O}(k) = \mathcal{S}[\mathbf{X} - \mathbf{1}_N \mathbf{m}'(k-1) - \mathbf{S}(k)\mathbf{U}'(k), (\lambda_2/2)\mathbf{I}_N]$ under Algorithm 1, where the thresholding operator $\mathcal{S}$ sets the entire outlier vector $\mathbf{o}_n(k)$ to zero whenever $\|\mathbf{r}_n(k)\|_2$ does not exceed $\lambda_2/2$, in par with the group sparsifying property of group Lasso. Interestingly, this is the same rule used to decide if datum $\mathbf{x}_n$ is deemed an outlier, in the equivalent formulation (8) which involves Huber's loss function. Whenever an $\ell_1$-norm regularizer is adopted as discussed in Remark 2, the only difference is that updates (10) boil down to soft-thresholding the scalar entries of $\mathbf{r}_n(k)$.

---

**Algorithm 1** : Batch robust PCA solver

Set $\mathbf{U}(0) = \mathbf{I}_p(:, 1:q)$ and $\mathbf{O}(0) = \mathbf{0}_{N \times p}$.

**for** $k = 1, 2, \ldots$ **do**

    Update $\mathbf{m}(k) = (\mathbf{X} - \mathbf{O}(k-1))' \mathbf{1}_N / N$.

    Form $\mathbf{X}_o(k) = \mathbf{X} - \mathbf{1}_N \mathbf{m}'(k) - \mathbf{O}(k-1)$.

    Update $\mathbf{S}(k) = \mathbf{X}_o(k) \mathbf{U}(k-1)$.

    Obtain $\mathbf{L}(k)\mathbf{D}(k)\mathbf{R}(k)' = \text{svd}[\mathbf{X}_o'(k)\mathbf{S}(k)]$ and update $\mathbf{U}(k) = \mathbf{L}(k)\mathbf{R}'(k)$.

    Update $\mathbf{O}(k) = \mathcal{S}\left[\mathbf{X} - \mathbf{1}_N \mathbf{m}'(k) - \mathbf{S}(k)\mathbf{U}'(k), (\lambda_2/2)\mathbf{I}_N\right]$.

**end for**

---

The entire AM solver is tabulated under Algorithm 1, indicating also the recommended initialization. Algorithm 1 is conceptually interesting, since it explicitly reveals the intertwining between the outlier identification process, and the PCA low-rank model fitting based on the outlier compensated data $\mathbf{X}_o(k)$.

The AM solver is also computationally efficient. Computing the $N \times q$ matrix $\mathbf{S}(k) = \mathbf{X}_o(k)\mathbf{U}(k-1)$ requires $Npq$ operations per iteration, and equally costly is to obtain $\mathbf{X}_o'(k)\mathbf{S}(k) \in \mathbb{R}^{p \times q}$. The cost of computing the SVD of $\mathbf{X}_o'(k)\mathbf{S}(k)$ is of order $\mathcal{O}(pq^2)$, while the rest of the operations including the row-wise soft-thresholdings to yield $\mathbf{O}(k)$ are linear in both $N$ and $p$. In summary, the total cost of Algorithm 1 is roughly $k_{\max}\mathcal{O}(Np + pq^2)$, where $k_{\max}$ is the number of iterations required for convergence (typically $k_{\max} = 5$ to $10$ iterations suffice). Because $q \leq p$ is typically small, Algorithm 1 is attractive computationally both under the classic setting where $N > p$, and $p$ is not large; as well as in high-dimensional data settings where $p \gg N$, a situation typically arising e.g., in microarray data analysis.

Because each of the optimization problems in the per-iteration cycles has a unique minimizer, and the nondifferentiable regularization only affects one of the variable groups ($\mathbf{O}$), the general results of [36] apply to establish convergence of Algorithm 1 as follows.

**Proposition 2:** *As $k \to \infty$, the iterates generated by Algorithm 1 converge to a stationary point of* (7).

*B. Selection of $\lambda_2$: robustification paths*

Selecting $\lambda_2$ controls the number of outliers rejected. But this choice is challenging because existing techniques such as cross-validation are not effective when outliers are present [32]. To this end, systematic data-driven approaches were devised in [14], which e.g., require a rough estimate of the percentage of outliers, or, robust estimates $\hat{\sigma}_e^2$ of the nominal noise variance that can be obtained using median absolute deviation (MAD) schemes [19]. These approaches can be adapted to the robust PCA setting considered here, and leverage the *robustification paths* of (group-)Lasso solutions [cf. (7)], which are defined as the

solution paths corresponding to $\|\hat{\mathbf{o}}_n\|_2$, $n = 1, \ldots, N$, for all values of $\lambda_2$. As $\lambda_2$ decreases, more vectors $\hat{\mathbf{o}}_n$ enter the model signifying that more of the training data are deemed to contain outliers.

Consider then a grid of $G_\lambda$ values of $\lambda_2$ in the interval $[\lambda_{\min}, \lambda_{\max}]$, evenly spaced on a logarithmic scale. Typically, $\lambda_{\max}$ is chosen as the minimum $\lambda_2$ value such that $\hat{\mathbf{O}} \neq \mathbf{0}_{N \times p}$, while $\lambda_{\min} = \epsilon \lambda_{\max}$ with $\epsilon = 10^{-4}$, say. Because Algorithm 1 converges quite fast, (7) can be efficiently solved over the grid of $G_\lambda$ values for $\lambda_2$. In the order of hundreds of grid points can be easily handled by initializing each instance of Algorithm 1 (per value of $\lambda_2$) using *warm starts* [17]. This means that multiple instances of (7) are solved for a sequence of decreasing $\lambda_2$ values, and the initialization of Algorithm 1 per grid point corresponds to the solution obtained for the immediately preceding value of $\lambda_2$ in the grid. For sufficiently close values of $\lambda_2$, one expects that the respective solutions will also be close (the row support of $\hat{\mathbf{O}}$ will most likely not change), and hence Algorithm 1 will converge after few iterations.

Based on the $G_\lambda$ samples of the robustification paths and the prior knowledge available on the outlier model (4), a couple of alternatives are also possible for selecting the 'best' value of $\lambda_2$ in the grid. A comprehensive survey of options can be found in [14].

*Number of outliers is known:* By direct inspection of the robustification paths one can determine the range of values for $\lambda_2$, such that the number of nonzero rows in $\hat{\mathbf{O}}$ equals the known number of outliers sought. Zooming-in to the interval of interest, and after discarding the identified outliers, $K$-fold cross-validation methods can be applied to determine the 'best' $\lambda_2^*$.

*Nominal noise covariance matrix is known:* Given $\boldsymbol{\Sigma}_e := E[\mathbf{e}_n \mathbf{e}_n']$, one can proceed as follows. Consider the estimates $\hat{\mathcal{V}}_g$ obtained using (7) after sampling the robustification path for each point $\{\lambda_{2,g}\}_{g=1}^G$. Next, pre-whiten those residuals corresponding to training data not deemed as containing outliers; i.e., form $\hat{\mathcal{R}}_g := \{\bar{\mathbf{r}}_{n,g} = \boldsymbol{\Sigma}_e^{-1/2}(\mathbf{x}_n - \hat{\mathbf{m}}_g - \hat{\mathbf{U}}_g \hat{\mathbf{s}}_{n,g}) : n \text{ s. to } \hat{\mathbf{o}}_n = \mathbf{0}\}$, and find the sample covariance matrices $\{\hat{\boldsymbol{\Sigma}}_{\bar{r},g}\}_{g=1}^G$. The winner $\lambda_2^* := \lambda_{2,g^*}$ corresponds to the grid point minimizing an absolute variance deviation criterion, namely $g^* := \arg\min_g |\text{tr}[\hat{\boldsymbol{\Sigma}}_{\bar{r},g}] - p|$.

### C. Connections with robust linear regression, dictionary learning, and clustering

Previous efforts towards robustifying linear regression have pointed out the equivalence between M-type estimators and $\ell_1$-norm regularized regression [13], and capitalized on this neat connection under a Bayesian framework [20]. However, they have not recognized the link to LTS via convex relaxation of the $\ell_0$-norm in (5). The treatment here goes beyond linear regression by considering the PCA framework, which entails a more challenging bilinear factor analysis model. Linear regression is subsumed as a special case, when matrix $\mathbf{U}$ is not necessarily tall but *assumed known*, while $\mathbf{s}_n = \mathbf{s}$, $\forall n = 1, \ldots, N$.

As an alternative to PCA, it is possible to device dimensionality reduction schemes when the data admit a sparse representation over a perhaps *unknown* basis. Such sparse representations comprise only a few elements (atoms) of the overcomplete basis (a.k.a. dictionary) to reconstruct the original data record. Thus, each datum is represented by a coefficient vector whose effective dimensionality (number of nonzero coefficients) is smaller than that of the original data vector. Recently, the *dictionary learning* paradigm offers techniques to design a dictionary over which the data assume a sparse representation; see e.g., [34] for a tutorial treatment. Dictionary learning schemes are flexible, in the sense that they utilize training data to learn an appropriate overcomplete basis customized for the data at hand [24], [34].

However, as in PCA the criteria adopted typically rely on a squared-error loss function as a measure of fit, which is known to be very sensitive to outliers [19], [32]. Interestingly, one can conceivably think of robustifying dictionary learning via minor modifications to the framework described so far. For instance, with the same matrix notation used in e.g., (5), one seeks to minimize

$$\min_{\mathcal{V},\mathbf{O}} \|\mathbf{X} - \mathbf{S}\mathbf{U}' - \mathbf{O}\|_F^2 + \lambda_1 \|\mathbf{S}\|_1 + \lambda_2 \|\mathbf{O}\|_{2,r}. \tag{11}$$

Different from the low-rank outlier-aware model adopted for PCA [cf. (4)], here the dictionary $\mathbf{U} \in \mathbb{R}^{p \times q}$ is fat $(q \gg p)$, with column vectors that are no longer orthogonal but still constrained to have unit $\ell_2$-norm. (This constraint is left implicit in (11) for simplicity.) Moreover, one seeks a sparse vector $\mathbf{s}_n$ to represent each datum $\mathbf{x}_n$, in terms of a few atoms of the learnt dictionary $\hat{\mathbf{U}}$. This is why (11) includes an additional sparsity-promoting $\ell_1$-norm regularization on $\mathbf{S}$, that is not present in (7). Sparsity is thus present both in the representation coefficients $\mathbf{S}$, as well as in the outliers $\mathbf{O}$.

Finally, it is shown here that a generative data model for K-means clustering [17] can share striking similarities with the bilinear model (1). Consequently, the sparsity-controlling estimator (7) can be adapted to robustify the K-means clustering task too [12]. Consider for instance that the data in $\mathcal{T}_x$ come from $q$ clusters, each of which is represented by a centroid $\mathbf{u}_i \in \mathbb{R}^p$, $i = 1, \ldots, q$. Moreover, for each input vector $\mathbf{x}_n$, K-means introduces the unknown membership variables $s_{ni} \in \{0,1\}$, $i = 1, \ldots, q$, where $s_{ni} = 1$ whenever $\mathbf{x}_n$ comes from cluster $i$, and $s_{ni} = 0$ otherwise. Typically, the membership variables are also constrained to satisfy $\sum_{n=1}^N s_{ni} > 0 \; \forall \, i$ (no empty clusters), and $\sum_{i=1}^q s_{ni} = 1 \; \forall \, n$ (single cluster membership). Upon defining $\mathbf{U} := [\mathbf{u}_1, \ldots, \mathbf{u}_q] \in \mathbb{R}^{p \times q}$ and the membership vectors $\mathbf{s}_n := [s_{n1}, \ldots, s_{nq}]' \in \mathbb{R}^q$, a pertinent model for hard K-means clustering assumes that input vectors can be expressed as $\mathbf{x}_n = \mathbf{U}\mathbf{s}_n + \mathbf{e}_n + \mathbf{o}_n$, where $\mathbf{e}_n$ and $\mathbf{o}_n$ are as in (4). Because the aforementioned constraints imply $\|\mathbf{s}_n\|_0 = \|\mathbf{s}_n\|_1 = 1 \; \forall \, n$, if $\mathbf{x}_n$ belongs to cluster $i$, then $s_{ni} = 1$ and in the absence of outliers one effectively has $\mathbf{x}_n = \mathbf{u}_i + \mathbf{e}_n$. Based on this data model, a natural approach towards

robustifying K-means clustering solves [12]

$$\min_{\mathcal{V},\mathbf{O}} \|\mathbf{X} - \mathbf{SU}' - \mathbf{O}\|_F^2 + \lambda_2\|\mathbf{O}\|_{2,r}, \quad \text{s. to } s_{ni} \in \{0,1\}, \sum_{n=1}^{N} s_{ni} > 0, \sum_{i=1}^{q} s_{ni} = 1. \quad (12)$$

Recall that in the robust PCA estimator (7), the subspace matrix is required to be orthonormal and the principal components are unrestrained. In the clustering context however, the centroid columns of $\mathbf{U}$ are free optimization variables, whereas the cluster membership variables adhere to the constraints in (12). Suitable relaxations to tackle the NP-hard problem (12) have been investigated in [12].

## IV. FURTHER ALGORITHMIC ISSUES

### A. Bias reduction through nonconvex regularization

Instead of substituting $\|\mathbf{O}\|_0$ in (5) by its closest convex approximation, namely $\|\mathbf{O}\|_{2,r}$, letting the surrogate function to be nonconvex can yield tighter approximations, and improve the statistical properties of the estimator. In rank minimization problems for instance, the logarithm of the determinant of the unknown matrix has been proposed as a smooth surrogate to the rank [11]; an alternative to the convex nuclear norm in e.g., [29]. Nonconvex penalties such as the smoothly clipped absolute deviation (SCAD) have been also adopted to reduce bias [10], present in uniformly weighted $\ell_1$-norm regularized estimators such as (7) [17, p. 92]. In the context of sparse signal reconstruction, the $\ell_0$-norm of a vector was surrogated in [6] by the logarithm of the geometric mean of its elements; see also [28].

Building on this last idea, consider approximating (5) by the *nonconvex* formulation

$$\min_{\mathcal{V},\mathbf{O}} \|\mathbf{X} - \mathbf{1}_N\mathbf{m}' - \mathbf{SU}' - \mathbf{O}\|_F^2 + \lambda_0 \sum_{n=1}^{N} \log(\|\mathbf{o}_n\|_2 + \delta), \quad \text{s. to } \mathbf{U}'\mathbf{U} = \mathbf{I}_q \quad (13)$$

where the small positive constant $\delta$ is introduced to avoid numerical instability. Since the surrogate term in (13) is concave, the overall minimization problem is nonconvex and admittedly more complex to solve than (7). Local methods based on iterative linearization of $\log(\|\mathbf{o}_n\|_2 + \delta)$ around the current iterate $\mathbf{o}_n(k)$, are adopted to minimize (13). Skipping details that can be found in [23], application of the majorization-minimization technique to (13) leads to an iteratively-reweighted version of (7), whereby $\lambda_2 \leftarrow \lambda_0 w_n(k)$ is used for updating $\mathbf{o}_n(k)$ in Algorithm 1. Specifically, per $k = 1, 2, \ldots$ one updates

$$\mathbf{O}(k) = \mathcal{S}\left[\mathbf{X} - \mathbf{1}_N\mathbf{m}'(k-1) - \mathbf{S}(k)\mathbf{U}'(k), (\lambda_0/2)\text{diag}(w_1(k), \ldots, w_N(k))\right]$$

where the weights are given by $w_n(k) = (\|\mathbf{o}_n(k-1)\|_2 + \delta)^{-1}$, $n = 1, \ldots, N$. Note that the thresholds vary both across rows (indexed by $n$), and across iterations. If the value of $\|\mathbf{o}_n(k-1)\|_2$ is small, then in the next iteration the regularization term $\lambda_0 w_n(k)\|\mathbf{o}_n\|_2$ has a large weight, thus promoting shrinkage

of that entire row vector to zero. If $\|\mathbf{o}_n(k-1)\|_2$ is large, the cost in the next iteration downweighs the regularization, and places more importance to the LS component of the fit.

All in all, the idea is to start from the solution of (7) for the 'best' $\lambda_2$, which is obtained using Algorithm 1. This initial estimate is refined after runnning a few iterations of the iteratively-reweighted counterpart to Algorithm 1. Extensive numerical tests suggest that even a couple iterations of this second stage refinement suffices to yield improved estimates $\hat{\mathcal{V}}$, in comparison to those obtained from (7). The improvements can be leveraged to bias reduction – and its positive effect with regards to outlier support estimation – also achieved by similar *weighted* norm regularizers proposed for linear regression [17, p. 92].

### B. Automatic rank determination: from nuclear- to Frobenius-norm regularization

Recall that $q \leq p$ is the dimensionality of the subspace where the outlier-free data (1) are assumed to live in, or equivalently, $q = \text{rank}[\mathbf{Y}]$ in the absence of noise. So far, $q$ was assumed known and fixed. This is reasonable in e.g., compression/quantization, where a target distortion-rate tradeoff dictates the maximum $q$. In other cases, the physics of the problem may render $q$ known. This is indeed the case in array processing for direction-of-arrival estimation, where $q$ is the dimensionality of the so-termed *signal subspace*, and is given by the number of plane waves impinging on a uniform linear array; see e.g., [40].

Other applications however, call for signal processing tools that can determine the 'best' $q$, as well as robustly estimate the underlying low-dimensional subspace $\mathbf{U}$ from data $\mathbf{X}$. Noteworthy representatives for this last kind of problems include unveiling traffic volume anomalies in large-scale networks [25], and automatic intrusion detection from video surveillance frames [5], [8], just to name a few. A related approach in this context is (stable) principal components pursuit (PCP) [38], [42], which solves

$$\min_{\mathbf{L},\mathbf{O}} \|\mathbf{X} - \mathbf{L} - \mathbf{O}\|_F^2 + \lambda_* \|\mathbf{L}\|_* + \lambda_2 \|\mathbf{O}\|_{2,r} \tag{14}$$

with the objective of reconstructing the low-rank matrix $\mathbf{L} \in \mathbb{R}^{N \times p}$, as well as the sparse matrix of outliers $\mathbf{O}$ in the presence of dense noise with known variance.[1] Note that $\|\mathbf{L}\|_*$ denotes the matrix nuclear norm, defined as the sum of the singular values of $\mathbf{L}$. The same way that the $\ell_2$-norm regularization promotes sparsity in the rows of $\hat{\mathbf{O}}$, the nuclear norm encourages a low-rank $\hat{\mathbf{L}}$ since it effects sparsity in the vector of singular values of $\mathbf{L}$. Upon solving the convex optimization problem (14), it is possible to obtain $\hat{\mathbf{L}} = \hat{\mathbf{S}}\hat{\mathbf{U}}'$ using the SVD. Interestingly, (14) does not fix (or require the knowledge of) $\text{rank}[\mathbf{L}]$ a fortiori, but controls it through the tuning parameter $\lambda_*$. Adopting a Bayesian framework, a similar problem was considered in [9].

---

[1] Actually, [42] considers entrywise outliers and adopts an $\ell_1$-norm regularization on $\mathbf{O}$.

Instead of assuming that $q$ is known, suppose that only an upper bound $\bar{q}$ is given. Then, the class of feasible noise-free low-rank matrix components of $\mathbf{Y}$ in (1) admit a factorization $\mathbf{L} = \mathbf{SU}'$, where $\mathbf{S}$ and $\mathbf{U}$ are $N \times \bar{q}$ and $p \times \bar{q}$ matrices, respectively. Building on the ideas used in the context of finding minimum rank solutions of linear matrix equations [29], a novel alternative approach to robustifying PCA is to solve

$$\min_{\mathbf{U,S,O}} \|\mathbf{X} - \mathbf{SU}' - \mathbf{O}\|_F^2 + \frac{\lambda_*}{2}(\|\mathbf{U}\|_F^2 + \|\mathbf{S}\|_F^2) + \lambda_2\|\mathbf{O}\|_{2,r}. \tag{15}$$

Different from (14) and (7), a Frobenius-norm regularization on both $\mathbf{U}$ and $\mathbf{S}$ is adopted to control the dimensionality of the estimated subspace $\hat{\mathbf{U}}$. Relative to (7), $\mathbf{U}$ in (15) is not constrained to be orthonormal. It is certainly possible to include the mean vector $\mathbf{m}$ in the cost of (15), as well as an $\ell_1$-norm regularization for entrywise outliers. The main motivation behind choosing the Frobenius-norm regularization comes from the equivalence of (14) with (15), as asserted in the ensuing result which adapts [29, Lemma 5.1] to the problem formulation considered here.

**Lemma 1:** *If $\{\hat{\mathbf{L}}, \hat{\mathbf{O}}\}$ minimizes (14) and rank$[\hat{\mathbf{L}}] \leq \bar{q}$, then (14) and (15) are equivalent.*

*Proof:* Because rank$[\hat{\mathbf{L}}] \leq \bar{q}$, the relevant feasible subset of (14) can be re-parametrized as $\{\mathbf{SU}', \mathbf{O}\}$, where $\mathbf{S}$ and $\mathbf{U}$ are $N \times \bar{q}$ and $p \times \bar{q}$ matrices, respectively. For every triplet $\{\mathbf{U}, \mathbf{S}, \mathbf{O}\}$ the objective of (15) is no smaller than the one of (14), since it holds that [29]

$$\|\mathbf{L}\|_* = \min_{\mathbf{U,S}} \frac{1}{2}(\|\mathbf{U}\|_F^2 + \|\mathbf{S}\|_F^2), \quad \text{s. to } \mathbf{L} = \mathbf{SU}'. \tag{16}$$

One can show that the gap between the objectives of (14) and (15) vanishes at $\mathbf{O}^* := \hat{\mathbf{O}}$, $\mathbf{S}^* := \mathbf{U}_L\mathbf{\Sigma}^{1/2}$, and $\mathbf{U}^* := \mathbf{V}_L\mathbf{\Sigma}^{1/2}$; where $\hat{\mathbf{L}} = \mathbf{U}_L\mathbf{\Sigma}\mathbf{V}'_L$ is the SVD of $\hat{\mathbf{L}}$. Therefore, from the previous arguments it follows that (14) and (15) attain the same global minimum objective, which completes the proof. ∎

Even though problem (15) is nonconvex, the number of optimization variables is reduced from $2Np$ to $Np + (N + p)\bar{q}$, which becomes significant when $\bar{q}$ is in the order of a few dozens and both $N$ and $p$ are large. Also note that the dominant $Np$-term in the variable count of (15) is due to $\mathbf{O}$, which is sparse and can be efficiently handled. While the factorization $\mathbf{L} = \mathbf{SU}'$ could have also been introduced in (14) to reduce the number of unknowns, the cost in (15) is separable and much simpler to optimize using e.g., an AM solver comprising the iterations tabulated as Algorithm 2. The decomposability of the Frobenius-norm regularizer has been recently exploited for parallel processing across multiple processors when solving large-scale matrix completion problems [30], or to unveil network anomalies [25].

Because (15) is a nonconvex optimization problem, most solvers one can think of will at most provide convergence guarantees to a stationary point that may not be globally optimum. Nevertheless, simulation results in Section VII demonstrate that Algorithm 2 is effective in providing good solutions most of

---

**Algorithm 2** : Batch robust PCA solver with controllable rank

Set $\mathbf{O}(0) = \mathbf{0}_{N \times p}$, and randomly initialize $\mathbf{S}(0)$.

**for** $k = 1, 2, \ldots$ **do**

    Update $\mathbf{m}(k) = [\mathbf{X} - \mathbf{O}(k-1)]'\mathbf{1}_N/N$.

    Form $\mathbf{X}_o(k) = \mathbf{X} - \mathbf{1}_N\mathbf{m}'(k) - \mathbf{O}(k-1)$.

    Update $\mathbf{U}(k) = \mathbf{X}_o(k)'\mathbf{S}(k-1)[\mathbf{S}'(k-1)\mathbf{S}(k-1) + (\lambda_*/2)\mathbf{I}_{\bar{q}}]^{-1}$.

    Update $\mathbf{S}(k) = \mathbf{X}_o(k)\mathbf{U}(k)[\mathbf{U}'(k)\mathbf{U}(k) + (\lambda_*/2)\mathbf{I}_{\bar{q}}]^{-1}$.

    Update $\mathbf{O}(k) = \mathcal{S}\left[\mathbf{X} - \mathbf{S}(k)\mathbf{U}'(k), \lambda_2/2\right]$.

**end for**

---

the time, which is somehow expected since there is quite a bit of structure in (15). Formally, the next proposition adapted from [25, Prop. 1] provides a sufficient condition under which Algorithm 2 yields an optimal solution of (14). For a proof of a slightly more general result, see [25].

**Proposition 3:** *If* $\{\bar{\mathbf{U}}, \bar{\mathbf{S}}, \bar{\mathbf{O}}\}$ *is a stationary point of* (15) *and* $\|\mathbf{X} - \bar{\mathbf{S}}\bar{\mathbf{U}}' - \bar{\mathbf{O}}\|_2 \leq \lambda_*/2$, *then* $\{\hat{\mathbf{L}} := \bar{\mathbf{S}}\bar{\mathbf{U}}', \hat{\mathbf{O}} := \bar{\mathbf{O}}\}$ *is the optimal solution of* (14).

## V. ROBUST SUBSPACE TRACKING

E-commerce and Internet-based retailing sites, the World Wide Web, and video surveillance systems generate huge volumes of data, which far outweigh the ability of modern computers to analyze them in real time. Furthermore, data are generated sequentially in time, which motivates updating previously obtained learning results rather than re-computing new ones from scratch each time a new datum becomes available. This calls for low-complexity real-time (adaptive) algorithms for robust subspace tracking.

One possible adaptive counterpart to (7) is the exponentially-weighted LS (EWLS) estimator found by

$$\min_{\{\mathcal{V}, \mathbf{O}\}} \sum_{n=1}^{N} \beta^{N-n} \left[\|\mathbf{x}_n - \mathbf{m} - \mathbf{U}\mathbf{s}_n - \mathbf{o}_n\|_2^2 + \lambda_2\|\mathbf{o}_n\|_2\right] \tag{17}$$

where $\beta \in (0, 1]$ is a forgetting factor. In this context, $n$ should be understood as a temporal variable, indexing the instants of data acquisition. Note that in forming the EWLS estimator (17) at time $N$, the entire history of data $\{\mathbf{x}_n\}_{n=1}^{N}$ is incorporated in the real-time estimation process. Whenever $\beta < 1$, past data are exponentially discarded thus enabling operation in nonstationary environments. Adaptive estimation of sparse signals has been considered in e.g., [1] and [24].

Towards deriving a real-time, computationally efficient, and recursive (approximate) solver of (17), an AM scheme will be adopted in which iterations $k$ coincide with the time scale $n = 1, 2, \ldots$ of data

acquisition. Per time instant $n$, a new datum $\mathbf{x}_n$ is drawn and the corresponding pair of decision variables $\{\mathbf{s}(n), \mathbf{o}(n)\}$ are updated via

$$\{\mathbf{s}(n), \mathbf{o}(n)\} := \arg \min_{\{\mathbf{s}, \mathbf{o}\}} \|\mathbf{x}_n - \mathbf{m}(n-1) - \mathbf{U}(n-1)\mathbf{s} - \mathbf{o}\|_2^2 + \lambda_2 \|\mathbf{o}\|_2. \tag{18}$$

As per (18), only $\mathbf{o}(n)$ is updated at time $n$, rather than the whole (growing with time) matrix $\mathbf{O}$ that minimization of (17) would dictate; see also [24] for a similar approximation.

Because (18) is a smooth optimization problem w.r.t. $\mathbf{s}$, from the first-order optimality condition the principal component update is $\mathbf{s}(n) = \mathbf{U}'(n-1)[\mathbf{x}_n - \mathbf{m}(n-1) - \mathbf{o}(n)]$. Interestingly, this resembles the projection approximation adopted in [40], and can only be evaluated after $\mathbf{o}(n)$ is obtained. To this end, plug $\mathbf{s}(n)$ in (18) to obtain $\mathbf{o}(n)$ via a particular instance of the group Lasso estimator

$$\mathbf{o}(n) = \arg \min_{\mathbf{o}} \|[\mathbf{I}_p - \mathbf{U}(n-1)\mathbf{U}'(n-1)](\mathbf{x}_n - \mathbf{m}(n-1) - \mathbf{o})\|_2^2 + \lambda_2 \|\mathbf{o}\|_2 \tag{19}$$

with a single group of size equal to $p$. The cost in (19) is non-differentiable at the origin, and different from e.g., ridge regression, it does not admit a closed-form solution. Upon defining

$$\mathbf{H}(n) := 2[\mathbf{I}_p - \mathbf{U}(n-1)\mathbf{U}'(n-1)]'[\mathbf{I}_p - \mathbf{U}(n-1)\mathbf{U}'(n-1)] \in \mathbb{R}^{p \times p} \tag{20}$$

$$\mathbf{g}(n) := -\mathbf{H}(n)[\mathbf{x}_n - \mathbf{m}(n-1)] \in \mathbb{R}^p \tag{21}$$

one can recognize (19) as the multidimensional shrinkage-thresholding operator $\mathcal{T}_{\mathbf{H}(n), \lambda_2}(\mathbf{g}(n))$ introduced in [27]. In particular, as per [27, Corollary 2] it follows that

$$\mathbf{o}(n) = \mathcal{T}_{\mathbf{H}(n), \lambda_2}(\mathbf{g}(n)) = \begin{cases} -(\mathbf{H}(n) + \gamma \mathbf{I}_p)^{-1}\mathbf{g}(n), & \text{if } \|\mathbf{g}(n)\|_2 > \lambda_2 \\ \mathbf{0}_p, & \text{otherwise} \end{cases} \tag{22}$$

where parameter $\gamma := \lambda_2^2/(2\eta)$ is such that $\eta > 0$ solves the scalar optimization

$$\min_{\eta > 0} \left(1 - \mathbf{g}'(n) \left(2\eta \mathbf{H}(n) + \lambda_2^2\right)^{-1} \mathbf{g}(n)\right) \eta. \tag{23}$$

Remarkably, one can easily determine if $\mathbf{o}(n) = \mathbf{0}_p$, by forming $\mathbf{g}(n)$ and checking whether $\|\mathbf{g}(n)\|_2 \leq \lambda_2$. This will be the computational burden incurred to solve (19) for most $n$, since outliers are typically sporadic and one would expect to obtain $\mathbf{o}(n) = \mathbf{0}_p$ most of the time. When datum $\mathbf{x}_n$ is deemed an outlier, $\|\mathbf{g}(n)\|_2 > \lambda_2$, and one needs to carry out the extra line search in (23) to determine $\mathbf{o}(n)$ as per (22); further details can be found in in [27]. Whenever an $\ell_1$-norm outlier regularization is adopted, the resulting counterpart of (19) can be solved using e.g., coordinate descent [1], or, the Lasso variant of least-angle regression (LARS) [24].

Moving on, the subspace update is given by

$$\mathbf{U}(n) = \arg \min_{\mathbf{U}} \sum_{i=1}^{n} \beta^{n-i} \|\mathbf{x}_i - \mathbf{m}(i-1) - \mathbf{U}\mathbf{s}(i) - \mathbf{o}(i)\|_2^2$$

---

**Algorithm 3** : Online robust (OR-)PCA

---

\* Batch initialization phase

Determine $\lambda_2$ and $\mathbf{U}(n_0)$ from $\{\mathbf{x}_n\}_{n=1}^{n_0}$, as in Section III-B. Initialize $\mathbf{P}(n_0) = 10^3 \mathbf{I}_p$ and $\mathbf{s}(n_0) = \mathbf{0}_q$.

\* Online phase

**for** $n = n_0 + 1, n_0 + 2, \ldots$ **do**

    Form $\mathbf{H}(n)$ and $\mathbf{g}(n)$ using (20) and (21).

    Update $\mathbf{o}(n) = \mathcal{T}_{\mathbf{H}(n), \lambda_2}(\mathbf{g}(n))$ via (22).

    Update $\mathbf{s}(n) = \mathbf{U}'(n-1)[\mathbf{x}_n - \mathbf{o}(n)]$.

    \* RLS subspace update

    Update $\mathbf{k}(n) = \mathbf{P}(n-1)\mathbf{s}(n)/[\beta + \mathbf{s}'(n)\mathbf{P}(n-1)\mathbf{s}(n)]$.

    Update $\mathbf{P}(n) = (1/\beta)[\mathbf{P}(n-1) - \mathbf{k}(n)(\mathbf{P}(n-1)\mathbf{s}(n))']$.

    Update $\mathbf{U}(n) = \mathbf{U}(n-1) + [\mathbf{x}_n - \mathbf{U}(n-1)\mathbf{s}(n) - \mathbf{o}(n)]\mathbf{k}'(n)$.

**end for**

---

and can be efficiently obtained from $\mathbf{U}(n-1)$, via a recursive LS update leveraging the matrix inversion lemma; see e.g., [40]. Note that the orthonormality constraint on $\mathbf{U}$ is not enforced here, yet the deviation from orthonormality is typically small as observed in [40]. Still, if orthonormal principal directions are required, an extra orthonormalization step can be carried out per iteration, or, once at the end of the process. Finally, $\mathbf{m}(n)$ is obtained recursively as the exponentially-weighted average of the outlier-compensated data $\{\mathbf{x}_i - \mathbf{o}(i)\}_{i=1}^n$. The resulting online robust (OR-)PCA algorithm and its initialization are summarized under Algorithm 3, where $\mathbf{m}$ and its update have been omitted for brevity.

For the batch case where all data in $\mathcal{T}_x$ are available for joint processing, two data-driven criteria to select $\lambda_2$ have been outlined in Section III-B. However, none of these sparsity-controlling mechanisms can be run in real-time, and selecting $\lambda_2$ for subspace tracking via OR-PCA is challenging. One possibility to circumvent this problem is to select $\lambda_2$ once during a short initialization (batch) phase of OR-PCA, and retain its value for the subsequent time instants. Specifically, the initialization phase of OR-PCA entails solving (7) using Algorithm 1, with a typically small batch of data $\{\mathbf{x}_n\}_{n=1}^{n_0}$. At time $n_0$, the criteria in Section III-B are adopted to find the 'best' $\lambda_2$, and thus obtain the subspace estimate $\hat{\mathbf{U}}(n_0)$ required to initialize the OR-PCA iterations.

Convergence analysis of OR-PCA algorithm is beyond the scope of the present paper, and is only confirmed via simulations. The numerical tests in Section VII also show that in the presence of outliers, the novel adaptive algorithm outperforms existing non-robust alternatives for subspace tracking.

## VI. ROBUSTIFYING KERNEL PCA

Kernel (K)PCA is a generalization to (linear) PCA, seeking principal components in a *feature space* nonlinearly related to the *input space* where the data in $\mathcal{T}_x$ live [33]. KPCA has been shown effective in performing nonlinear feature extraction for pattern recognition [33]. In addition, connections between KPCA and spectral clustering [17, p. 548] motivate well the novel KPCA method developed in this section, to robustly identify cohesive subgroups (communities) from social network data.

Consider a nonlinear function $\phi : \mathbb{R}^p \to \mathcal{H}$, that maps elements from the input space $\mathbb{R}^p$ to a feature space $\mathcal{H}$ of arbitrarily large – possibly infinite – dimensionality. Given transformed data $\mathcal{T}_{\mathcal{H}} := \{\phi(\mathbf{x}_n)\}_{n=1}^N$, the proposed approach to robust KPCA fits the model

$$\phi(\mathbf{x}_n) = \mathbf{m} + \mathbf{U}\mathbf{s}_n + \mathbf{e}_n + \mathbf{o}_n, \quad n = 1, \ldots, N \tag{24}$$

by solving ($\boldsymbol{\Phi} := [\phi(\mathbf{x}_1), \ldots, \phi(\mathbf{x}_N)]$)

$$\min_{\mathbf{U}, \mathbf{S}, \mathbf{O}} \|\boldsymbol{\Phi}' - \mathbf{1}_N \mathbf{m}' - \mathbf{S}\mathbf{U}' - \mathbf{O}\|_F^2 + \frac{\lambda_*}{2}(\|\mathbf{U}\|_F^2 + \|\mathbf{S}\|_F^2) + \lambda_2 \|\mathbf{O}\|_{2,r}. \tag{25}$$

It is certainly possible to adopt the criterion (7) as well, but (25) is chosen here for simplicity in exposition. Except for the principal components' matrix $\mathbf{S} \in \mathbb{R}^{N \times \bar{q}}$, both the data and the unknowns in (25) are now vectors/matrices of generally infinite dimension. In principle, this challenges the optimization task since it is impossible to store, or, perform updates of such quantities directly. For these reasons, assuming zero-mean data $\phi(\mathbf{x}_n)$, or, the possibility of mean compensation for that matter, cannot be taken for granted here [cf. Remark 1]. Thus, it is important to explicitly consider the estimation of $\mathbf{m}$.

Interestingly, this hurdle can be overcome by endowing $\mathcal{H}$ with the structure of a reproducing kernel Hilbert space (RKHS), where inner products between any two members of $\mathcal{H}$ boil down to evaluations of the reproducing kernel $K_{\mathcal{H}} : \mathbb{R}^p \times \mathbb{R}^p \to \mathbb{R}$, i.e., $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle_{\mathcal{H}} = K_{\mathcal{H}}(\mathbf{x}_i, \mathbf{x}_j)$. Specifically, it is possible to form the kernel matrix $\mathbf{K} := \boldsymbol{\Phi}'\boldsymbol{\Phi} \in \mathbb{R}^{N \times N}$, without directly working with the vectors in $\mathcal{H}$. This so-termed *kernel trick* is the crux of most kernel methods in machine learning [17], including kernel PCA [33]. The problem of selecting $K_{\mathcal{H}}$ (and $\phi$ indirectly) will not be considered here.

Building on these ideas, it is shown in the sequel that Algorithm 2 can be *kernelized*, to solve (25) at affordable computational complexity and memory storage requirements that do not depend on the dimensionality of $\mathcal{H}$.

**Proposition 4:** *For $k \geq 1$, the sequence of iterates generated by Algorithm 2 when applied to solve (25) can be written as $\mathbf{m}(k) = \boldsymbol{\Phi}\boldsymbol{\mu}(k)$, $\mathbf{U}(k) = \boldsymbol{\Phi}\boldsymbol{\Upsilon}(k)$, and $\mathbf{O}'(k) = \boldsymbol{\Phi}\boldsymbol{\Omega}(k)$. The quantities $\boldsymbol{\mu}(k) \in \mathbb{R}^N$, $\boldsymbol{\Upsilon}(k) \in \mathbb{R}^{N \times \bar{q}}$, and $\boldsymbol{\Omega}(k) \in \mathbb{R}^{N \times N}$ are recursively updated as in Algorithm 4, without the need of operating with vectors in $\mathcal{H}$.*

*Proof:* The proof relies on an inductive argument. Suppose that at iteration $k-1$, there exists a matrix $\boldsymbol{\Omega}(k-1) \in \mathbb{R}^{N \times N}$ such that the outliers can be expressed as $\mathbf{O}'(k-1) = \boldsymbol{\Phi}\boldsymbol{\Omega}(k-1)$. From Algorithm 2, the update for the mean vector is $\mathbf{m}(k) = [\boldsymbol{\Phi}' - \mathbf{O}(k-1)]'\mathbf{1}_N/N = [\boldsymbol{\Phi} - \boldsymbol{\Phi}\boldsymbol{\Omega}(k-1)]\mathbf{1}_N/N = \boldsymbol{\Phi}\boldsymbol{\mu}(k)$ where $\boldsymbol{\mu}(k) := [\mathbf{I}_n - \boldsymbol{\Omega}(k-1)]\mathbf{1}_N/N$. Likewise, $\mathbf{X}_o(k) = \boldsymbol{\Phi}' - \mathbf{1}_N\boldsymbol{\mu}'(k)\boldsymbol{\Phi}' - \boldsymbol{\Omega}'(k-1)\boldsymbol{\Phi}'$ so that one can write the subspace update as $\mathbf{U}(k) = \boldsymbol{\Phi}\boldsymbol{\Upsilon}(k)$, upon defining

$$\boldsymbol{\Upsilon}(k) := [\mathbf{I}_N - \boldsymbol{\mu}(k)\mathbf{1}_N' - \boldsymbol{\Omega}(k-1)]\mathbf{S}(k-1)[\mathbf{S}'(k-1)\mathbf{S}(k-1) + (\lambda_*/2)\mathbf{I}_{\bar{q}}]^{-1}.$$

With regards to the principal components, it follows that (cf. Algorithm 2)

$$\mathbf{S}(k) = [\mathbf{I}_N - \mathbf{1}_N\boldsymbol{\mu}'(k) - \boldsymbol{\Omega}'(k-1)]\boldsymbol{\Phi}'\boldsymbol{\Phi}\boldsymbol{\Upsilon}(k)[\boldsymbol{\Upsilon}(k)'\boldsymbol{\Phi}'\boldsymbol{\Phi}\boldsymbol{\Upsilon}(k) + (\lambda_*/2)\mathbf{I}_{\bar{q}}]^{-1}$$

$$= [\mathbf{I}_N - \mathbf{1}_N\boldsymbol{\mu}'(k) - \boldsymbol{\Omega}'(k-1)]\mathbf{K}\boldsymbol{\Upsilon}(k)[\boldsymbol{\Upsilon}(k)'\mathbf{K}\boldsymbol{\Upsilon}(k) + (\lambda_*/2)\mathbf{I}_{\bar{q}}]^{-1} \qquad (26)$$

which is expressible in terms of the kernel matrix $\mathbf{K} := \boldsymbol{\Phi}'\boldsymbol{\Phi}$. Finally, the columns $\mathbf{o}_n(k)$ are given by the vector soft-thresholding operation (10), where the residuals are

$$\mathbf{r}_n(k) = \phi(\mathbf{x}_n) - \mathbf{m}(k) - \mathbf{U}(k)\mathbf{s}_n(k) = \boldsymbol{\Phi}[\mathbf{b}_{N,n} - \boldsymbol{\mu}(k) - \boldsymbol{\Upsilon}(k)\mathbf{s}_n(k)] := \boldsymbol{\Phi}\boldsymbol{\rho}_n(k).$$

Upon stacking all columns $\mathbf{o}_n(k)$, $n = 1, \ldots, N$, one readily obtains [cf. (10)]

$$\mathbf{O}'(k) = \boldsymbol{\Phi}[\mathbf{I}_N - \boldsymbol{\mu}(k)\mathbf{1}_N' - \boldsymbol{\Upsilon}(k)\mathbf{S}'(k)]\boldsymbol{\Lambda}(k) \qquad (27)$$

where $\boldsymbol{\Lambda}(k) := \mathrm{diag}((\|\mathbf{r}_1(k)\|_2 - \lambda_2/2)_+/\|\mathbf{r}_1(k)\|_2, \ldots, (\|\mathbf{r}_N(k)\|_2 - \lambda_2/2)_+/\|\mathbf{r}_N(k)\|_2)$. Interestingly, the diagonal elements of $\boldsymbol{\Lambda}(k)$ can be computed using the kernel matrix, since $\|\mathbf{r}_n(k)\|_2 = \sqrt{\boldsymbol{\rho}_n'(k)\mathbf{K}\boldsymbol{\rho}_n(k)}$, $n = 1, \ldots, N$. From (27) it is apparent that one can write $\mathbf{O}'(k) = \boldsymbol{\Phi}\boldsymbol{\Omega}(k)$, after defining

$$\boldsymbol{\Omega}(k) := [\mathbf{I}_N - \boldsymbol{\mu}(k)\mathbf{1}_N' - \boldsymbol{\Upsilon}(k)\mathbf{S}'(k)]\boldsymbol{\Lambda}(k).$$

The proof is concluded by noting that for $k = 0$, Algorithm 2 is initialized with $\mathbf{O}'(0) = \mathbf{0}_{p \times N}$. One can thus satisfy the inductive base case $\mathbf{O}'(0) = \boldsymbol{\Phi}\boldsymbol{\Omega}(0)$, by letting $\boldsymbol{\Omega}(0) = \mathbf{0}_{N \times N}$. ■

In order to run the novel robust KPCA algorithm (tabulated as Algorithm 4), one does not have to store or process the quantities $\mathbf{m}(k)$, $\mathbf{U}(k)$, and $\mathbf{O}(k)$. As per Proposition 4, the iterations of the provably convergent AM solver in Section IV-B can be equivalently carried out by cycling through *finite-dimensional* 'sufficient statistics' $\boldsymbol{\mu}(k) \to \boldsymbol{\Upsilon}(k) \to \mathbf{S}(k) \to \boldsymbol{\Omega}(k)$. In other words, the iterations of the robust kernel PCA algorithm are devoid of algebraic operations among vectors in $\mathcal{H}$. Recall that the size of matrix $\mathbf{S}$ is independent of the dimensionality of $\mathcal{H}$. Nevertheless, its update in Algorithm 2 cannot be carried out verbatim in the high-dimensional setting here, and is instead kernelized to yield the update rule (26).

Because $\mathbf{O}'(k) = \boldsymbol{\Phi}\boldsymbol{\Omega}(k)$ and upon convergence of the algorithm, the outlier vector norms are computable in terms of $\mathbf{K}$, i.e., $[\|\mathbf{o}_1(\infty)\|_2^2, \ldots, \|\mathbf{o}_N(\infty)\|_2^2]' = \mathrm{diag}[\boldsymbol{\Omega}'(\infty)\mathbf{K}\boldsymbol{\Omega}(\infty)]$. These are critical to

---

**Algorithm 4** : Robust KPCA solver

---

Initialize $\boldsymbol{\Omega}(0) = \mathbf{0}_{N \times N}$, $\mathbf{S}(0)$ randomly, and form $\mathbf{K} = \boldsymbol{\Phi}'\boldsymbol{\Phi}$.

**for** $k = 1, 2, \ldots$ **do**

    Update $\boldsymbol{\mu}(k) = [\mathbf{I}_n - \boldsymbol{\Omega}(k-1)]\mathbf{1}_N / N$.

    Form $\boldsymbol{\Phi}_o(k) = \mathbf{I}_N - \boldsymbol{\mu}(k)\mathbf{1}'_N - \boldsymbol{\Omega}(k-1)$.

    Update $\boldsymbol{\Upsilon}(k) = \boldsymbol{\Phi}_o(k)\mathbf{S}(k-1)[\mathbf{S}'(k-1)\mathbf{S}(k-1) + (\lambda_*/2)\mathbf{I}_{\bar{q}}]^{-1}$.

    Update $\mathbf{S}(k) = \boldsymbol{\Phi}'_o(k)\mathbf{K}\boldsymbol{\Upsilon}(k)[\boldsymbol{\Upsilon}(k)'\mathbf{K}\boldsymbol{\Upsilon}(k) + (\lambda_*/2)\mathbf{I}_{\bar{q}}]^{-1}$.

    Form $\boldsymbol{\rho}_n(k) = \mathbf{b}_{N,n} - \boldsymbol{\mu}(k) - \boldsymbol{\Upsilon}(\boldsymbol{k})\mathbf{s}_n(k)$, $n = 1, \ldots, N$, and update $\boldsymbol{\Lambda}(k)$.

    Update $\boldsymbol{\Omega}(k) = [\mathbf{I}_N - \boldsymbol{\mu}(k)\mathbf{1}'_N - \boldsymbol{\Upsilon}(\boldsymbol{k})\mathbf{S}'(k)]\boldsymbol{\Lambda}(k)$.

**end for**

---

determine the robustification paths needed to carry out the outlier sparsity control methods in Section III-B. Moreover, the principal component corresponding to any given new data point $\mathbf{x}$ is obtained through the projection $\mathbf{s} = \mathbf{U}(\infty)'[\boldsymbol{\phi}(\mathbf{x}) - \mathbf{m}(\infty)] = \boldsymbol{\Upsilon}'(\infty)\boldsymbol{\Phi}'\boldsymbol{\phi}(\mathbf{x}) - \boldsymbol{\Upsilon}'(\infty)\mathbf{K}\boldsymbol{\mu}(\infty)$, which is again computable after $N$ evaluations the kernel function $K_{\mathcal{H}}$.

## VII. NUMERICAL TESTS

### A. Synthetic data tests

To corroborate the effectiveness of the proposed robust methods, experiments with computer generated data are carried out first. These are important since they provide a 'ground truth', against which performance can be assessed by evaluating suitable figures of merit.

**Outlier-sparsity control.** To generate the data (4), a similar setting as in [42, Sec. V] is considered here with $N = p$ and $\mathbf{m} = \mathbf{0}_p$. For $n = 1, \ldots, N$, the errors are $\mathbf{e}_n \sim \mathcal{N}(\mathbf{0}_p, \sigma_e^2 \mathbf{I}_p)$ (multivariate normal distribution) and i.i.d. The entries of $\mathbf{U}$ and $\{\mathbf{s}_n\}_{n=1}^N$ are i.i.d. zero-mean Gaussian distributed, with variance $\sigma_{U,s}^2 = 10\sigma_e/\sqrt{N}$. Outliers are generated as $\mathbf{o}_n = \mathbf{p}_n \odot \mathbf{q}_n$, where the entries of $\mathbf{p}_n$ are i.i.d. Bernoulli distributed with parameter $\rho_p$, and $\mathbf{q}_n$ has i.i.d. entries drawn from a uniform distribution supported on $[-5, 5]$. The chosen values of the parameters are $N = p = 200$, $q = 20$, $\rho_p = 0.01$, and varying noise levels $\sigma_e^2 = \{0.01, 0.05, 0.1, 0.25, 0.5\}$.

In this setup, the ability to recover the low-rank component of the data $\mathbf{L} := \mathbf{S}\mathbf{U}'$ is tested for the sparsity-controlling robust PCA method of this paper [cf. (7)], stable PCP (14), and (non-robust) PCA. The $\ell_1$-norm regularized counterparts of (7) and (14) are adopted to deal with entry-wise outliers. Both values of $q$ and $\sigma_e^2$ are assumed known to obtain $\hat{\mathbf{L}} := \hat{\mathbf{S}}\hat{\mathbf{U}}'$ and $\hat{\mathbf{O}}$ via (7). This way, $\lambda_2$ is chosen using the sparsity-controlling algorithm of Section III-B, searching over a grid where $G_\lambda = 200$, $\lambda_{\min} = 10^{-2}\lambda_{\max}$,

and $\lambda_{\max} = 20$. In addition, the solutions of (7) are refined by running two iterations of the iteratively reweighted algorithm in Section IV-A, where $\delta = 10^{-5}$. Regarding SPCP, only the knowledge of $\sigma_e^2$ is required to select the tuning parameters $\lambda_* = 2\sqrt{2N\sigma_e^2}$ and $\lambda_2 = 2\sqrt{2\sigma_e^2}$ in (14), as suggested in [42]. Finally, the best rank $q$ approximation to the data $\mathbf{X}$ is obtained using standard PCA.

The results are summarized in Table I, which shows the estimation errors $\bar{\text{err}} := \|\mathbf{L} - \hat{\mathbf{L}}\|_F / N$ attained by the aforementioned schemes, averaged over 15 runs of the experiment. The 'best' tuning parameters $\lambda_2^*$ used in (7) are also shown. Both robust schemes attain an error which is approximately an order of magnitude smaller than PCA. With the additional knowledge of the true data rank $q$, the sparsity-controlling algorithm of this paper outperforms stable PCP in terms of $\bar{\text{err}}$. This numerical test is used to validate Proposition 3 as well. For the same values of the tuning parameters chosen for (14) and the rank upper-bound set to $\bar{q} = 2q$, Algorithm 2 is run to obtain the solution $\{\bar{\mathbf{U}}, \bar{\mathbf{S}}, \bar{\mathbf{O}}\}$ of the nonconvex problem (15). The average (across realizations and values of $\sigma_e^2$) errors obtained are $\|\hat{\mathbf{L}} - \bar{\mathbf{S}}\bar{\mathbf{U}}'\|_F / N = 0.15 \times 10^{-6}$ and $\|\hat{\mathbf{O}} - \bar{\mathbf{O}}\|_F / N = 0.78 \times 10^{-7}$, where $\{\hat{\mathbf{L}}, \hat{\mathbf{O}}\}$ is the solution of stable PCP [cf. (14)]. Thus, the solutions are identical for all practical purposes.

**Identification of invalid survey protocols.** Robust PCA is tested here to identify invalid or otherwise aberrant item response (questionnaire) data in surveys, that is, to flag and hold in abeyance data that may negatively influence (i.e., bias) subsequent data summaries and statistical analyses. In recent years, item response theory (IRT) has become the dominant paradigm for constructing and evaluating questionnaires in the biobehavioral and health sciences and in high-stakes testing (e.g., in the development of college admission tests); see e.g., [37]. IRT entails a class of nonlinear models characterizing an individual's item response behavior by one or more latent traits, and one or more item parameters. An increasingly popular IRT model for survey data is the 2-parameter logistic IRT model (2PLM) [31]. 2PLM characterizes the probability of a keyed (endorsed) response $y_{nm}$, as a nonlinear function of a weighted difference between a person parameter $\theta_n$ and an item parameter $b_m$

$$\Pr(y_{nm} = 1 | \theta_n) = \frac{e^{1.7a_m(\theta_n - b_m)}}{1 + e^{1.7a_m(\theta_n - b_m)}} \tag{28}$$

where $\theta_n$ is a latent trait value for individual $n$; $a_m$ is an item discrimination parameter (similar to a factor loading) for item $m$; and $b_m$ is an item difficulty (or extremity) parameter for item $m$.

Binary item responses ('agree/disagree' response format) were generated for $N = 1,000$ hypothetical subjects who were administered $p = 200$ items (questions). The 2PLM function (28) was used to generate the underlying item response probabilities, which were converted into binary item responses as follows: a response was coded 1 if $\Pr(y_{nm} | \theta_n) \geq \mathcal{U}(0,1)$, and coded 0 otherwise, where $\mathcal{U}[0,1]$ denotes a uniform

random deviate over $[0, 1]$. Model parameters were randomly drawn as $\{a_m\}_{m=1}^{200} \sim \mathcal{U}[1, 1.5]$, $\{b_m\}_{m=1}^{200} \sim$ $\mathcal{U}[-2, 2]$, and $\{\boldsymbol{\theta}_l\}_{l=1}^{200} \sim \mathcal{N}(\mathbf{0}_5, \mathbf{I}_5)$. Each of the 200 items loaded on one of $q = 5$ latent factors. To simulate random responding – a prevalent form of aberrancy in e.g., web-collected data – rows 101-120 of the item response matrix $\mathbf{Y}$ were modified by (re)drawing each of the entries from a Bernoulli distribution with parameter 0.5, thus yielding the corrupted matrix $\mathbf{X}$.

Robust PCA in (7) was adopted to identify invalid survey data, with $q = 5$, and $\lambda_2$ chosen such that $\|\hat{\mathbf{O}}\|_0 = 150$, a safe overestimate of the number of outliers. Results of this study are summarized in Fig. 1, which displays the 100 largest outliers ($\|\hat{\mathbf{o}}_n\|_2$) from the robust PCA analysis of the $N = 1,000$ simulated response vectors. When the outliers are plotted against their ranks, there is an unmistakable break between the 20th and 21st ordered value indicating that the method correctly identified the *number* of aberrant response patterns in $\mathbf{X}$. Perhaps more impressively, the method also correctly identified rows 101-to-120 as containing the invalid data.

**Online robust subspace estimation.** A simulated test is carried out here to corroborate the convergence and effectiveness of the OR-PCA algorithm in Section V. For $N = 2,000$, $p = 150$, and $q = 5$, nominal data in $\mathcal{T}_y$ are generated according to the stationary model (1), where $\mathbf{e}_n \sim \mathcal{N}(\mathbf{0}_p, 10^{-3}\mathbf{I}_p)$. Vectors $\mathbf{x}_{1001}, \ldots, \mathbf{x}_{1005}$ are outliers, uniformly i.i.d. over $[-0.5, 0.5]$. The results depicted in Fig. 2 are obtained after averaging over 50 runs. Fig. 2 (left) depicts the time evolution of the angle between the learnt subspace (spanned by the columns of) $\hat{\mathbf{U}}(n)$ and the true subspace $\mathbf{U}$ generating $\mathcal{T}_y$, where $\lambda_2 = 1.65$ and $\beta = 0.99$. The convergent trend of Algorithm 3 to $\mathbf{U}$ is apparent; and markedly outperforms the non-robust subspace tracking method in [40], and the first-order GROUSE algorithm in [2]. Note that even though $\mathbf{U}$ is time-invariant, it is meaningful to select $0 \ll \beta < 1$ to quickly 'forget' and recover from the outliers. A similar trend can be observed in Fig. 2 (right), which depicts the time evolution of the reconstruction error $\|\mathbf{y}_n - \hat{\mathbf{U}}(n)\hat{\mathbf{U}}(n)'\mathbf{y}_n\|_2^2/p$.

**Robust spectral clustering.** The following simulated test demonstrates that robust KPCA in Section VI can be effectively used to robustify spectral clustering (cf. the connection between both non-robust methods in e.g., [17, p. 548]). Adopting the data setting from [17, p. 546]), $N = 450$ points in $\mathbb{R}^2$ are generated from three circular concentric clusters, with respective radii of 1, 2.8, and 5. The points are uniformly distributed in angle, and additive noise $\mathbf{e}_n \sim \mathcal{N}(\mathbf{0}_2, 0.15\mathbf{I}_2)$ is added to each datum. Five outliers $\{\mathbf{x}_n\}_{n=451}^{455}$ uniformly distributed in the square $[-7, 7]^2$ complete the training data $\mathcal{T}_x$; see Fig. 3 (left). To unveil the cluster structure from the data, Algorithm 4 is run using the Gaussian radial kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2/c)$, with $c = 10$. The sparsity-controlling parameter is set to $\lambda_2 = 1.85$ so that $\|\hat{\mathbf{O}}\|_0 = 5$, while $\lambda_* = 1$, and $\bar{q} = 2$. Upon convergence, the vector of estimated outlier norms is

$[\|\mathbf{o}_1(\infty)\|_2^2, \ldots, \|\mathbf{o}_{N+5}(\infty)\|_2^2]' = [0, \ldots, 0, 10^{-4}, 1.3 \times 10^{-3}, 1.5 \times 10^{-2}, 10^{-2}, 1.7 \times 10^{-2}]'$, which shows that the outliers are correctly identified. Estimates of the (rotated) first two dominant eigenvectors of the kernel matrix $\mathbf{K}$ are obtained as the columns of $\hat{\mathbf{\Upsilon}}$, and are depicted in Fig. 3 (right). After removing the rows of $\hat{\mathbf{\Upsilon}}$ corresponding to the outliers [black points in Fig. 3 (right)], e.g., K-means clustering of the remaining points in Fig. 3 (right) will easily reveal the three clusters sought. From Fig. 3 (right) it is apparent that a non-robust KPCA method will incorrectly assign the outliers to the outer (green) cluster.

## B. Real data tests

**Video surveillance.** To validate the proposed approach to robust PCA, Algorithm 1 was tested to perform background modeling from a sequence of video frames; an approach that has found widespread applicability for intrusion detection in video surveillance systems. The experiments were carried out using the dataset studied in [8], which consists of $N = 520$ images ($p = 120 \times 160$) acquired from a static camera during two days. The illumination changes considerably over the two day span, while approximately $40\%$ of the training images contain people in various locations. For $q = 10$, both standard PCA and the robust PCA of Section III were applied to build a low-rank background model of the environment captured by the camera. For robust PCA, $\ell_1$-norm regularization on $\mathbf{O}$ was adopted to identify outliers at a pixel level. The outlier sparsity-controlling parameter was chosen as $\lambda_2 = 9.69 \times 10^{-4}$, whereas a single iteration of the reweighted scheme in Section IV-A was run to reduce the bias in $\hat{\mathbf{O}}$.

Results are shown in Fig. 1, for three representative images. The first column comprises the original frames from the training set, while the second column shows the corresponding PCA image reconstructions. The presence of undesirable 'ghostly' artifacts is apparent, since PCA is unable to completely separate the people from the background. The third column illustrates the robust PCA reconstructions, which recover the illumination changes while successfully subtracting the people. The fourth column shows the reshaped outlier vectors $\hat{\mathbf{o}}_n$, which mostly capture the people and abrupt changes in illumination.

**Robust measurement of the Big Five personality factors.** The 'Big Five' are five factors ($q = 5$) of personality traits, namely extraversion, agreeableness, conscientiousness, neuroticism, and openness; see e.g., [21]. The Big Five inventory (BFI) on the other hand, is a brief questionnaire (44 items in total) tailored to measure the Big Five dimensions. Subjects taking the questionnaire are asked to rate in a scale from 1 (disagree strongly) to 5 (agree strongly), items of the form 'I see myself as someone who is talkative'. Each item consists of a short phrase correlating (positively or negatively) with one factor; see e.g., [21, pp. 157-58] for a copy of the BFI and scoring instructions.

Robust PCA is used to identify aberrant responses from real BFI data comprising the Eugene-Springfield

community sample [16]. The rows of $\mathbf{X}$ contain the $p = 44$ item responses for each one of the $N = 437$ subjects under study. For $q = 5$, (7) is solved over grid of $G_\lambda = 200$ values of $\lambda_2$, where $\lambda_{\min} = 10^{-2}\lambda_{\max}$, and $\lambda_{\max} = 20$. The first plot of Fig. 5 (left) shows the evolution of $\hat{\mathbf{O}}$'s row support as a function of $\lambda_2$ with black pixels along the $n$th row indicating that $\|\hat{\mathbf{o}}_n\|_2 = 0$, and white ones reflecting that the responses from subject $n$ are deemed as outliers for the given $\lambda_2$. For example subjects $n = 418$ and $204$ are strong outlier candidates due to random responding, since they enter the model ($\|\hat{\mathbf{o}}_n\|_2 > 0$) for relatively large values of $\lambda_2$. The responses of e.g., subjects $n = 63$ (all items rated '3') and $249$ (41 items rated '3' and 3 items rated '4') are also undesirable, but are well modeled by (1) and are only deemed as outliers when $\lambda_2$ is quite small. These two observations are corroborated by the second plot of Fig. 5 (left), which shows the robust PCA results on a corrupted dataset, obtained from $\mathbf{X}$ by overwriting: (i) rows $151 - 160$ with random item responses drawn from a uniform distribution over $\{1, 2, 3, 4, 5\}$; and (ii) rows $301 - 310$ with constant item responses of value 3.

For $\lambda_2 = 5.6107$ corresponding to $\|\hat{\mathbf{O}}\|_0 = 100$, Fig. 5 (right) depicts the norm of the 40 largest outliers. Following the methodology outlined in Section VII-A, 8 subjects including $n = 418$ and $204$ are declared as outliers by robust PCA. As a means of validating these results, the following procedure is adopted. Based on the BFI scoring key [21], a list of all pairs of items hypothesized to yield positively correlated responses is formed. For each $n$, one counts the 'inconsistencies' defined as the number of times that subject $n$'s ratings for these pairs differ in more than four, in absolute value. Interestingly, after rank-ordering all subjects in terms of this inconsistency score, one finds that $n = 418$ ranks highest with a count of 17, $n = 204$ ranks second (10), and overall the eight outliers found rank in the top twenty.

**Unveiling communities in social networks.** Next, robust KPCA is used to identify communities and outliers in a network of $N = 115$ college football teams, by capitalizing on the connection between KPCA and spectral clustering [17, p. 548]. Nodes in the network graph represent teams belonging to eleven conferences (plus five independent teams), whereas (unweighted) edges joining pairs of nodes indicate that both teams played against each other during the Fall 2000 Division I season [15]. The kernel matrix used to run robust KPCA is $\mathbf{K} = \zeta\mathbf{I}_N + \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$, where $\mathbf{A}$ and $\mathbf{D}$ denote the graph adjacency and degree matrices, respectively; while $\zeta > 0$ is chosen to render $\mathbf{K}$ positive semi-definite. The tuning parameters are chosen as $\lambda_2 = 1.297$ so that $\|\hat{\mathbf{O}}\|_0 = 10$, while $\lambda_* = 1$, and $\bar{q} = 3$. Fig. 6 (left) shows the entries of $\mathbf{K}$, where rows and columns are permuted to reveal the clustering structure found by robust KPCA (after removing the outliers); see also Fig. 6 (right). The quality of the clustering is assessed through the adjusted rand index (ARI) after excluding outliers [12], which yielded the value 0.8967. Four of the teams deemed as outliers are Connecticut, Central Florida, Navy, and Notre Dame,

which are indeed teams not belonging to any major conference. The community structure of traditional powerhouse conferences such as Big Ten, Big 12, ACC, Big East, and SEC was identified exactly.

## VIII. CONCLUDING SUMMARY

Outlier-robust PCA methods were developed in this paper, to obtain low-dimensional representations of (corrupted) data. Bringing together the seemingly unrelated fields of robust statistics and sparse regression, the novel robust PCA framework was found rooted at the crossroads of outlier-resilient estimation, learning via (group-) Lasso and kernel methods, and real-time adaptive signal processing. Social network analysis, video surveillance, and psychometrics, were highlighted as relevant application domains.

## APPENDIX

Towards establishing the equivalence between problems (7) and (8), consider the pair $\{\hat{\mathcal{V}}, \hat{\mathbf{O}}\}$ that solves (7). Assume that $\hat{\mathcal{V}}$ is given, and the goal is to determine $\hat{\mathbf{O}}$. Upon defining the residuals $\hat{\mathbf{r}}_n := \mathbf{x}_n - \hat{\mathbf{m}} - \hat{\mathbf{U}}\hat{\mathbf{s}}_n$ and from the row-wise decomposability of $\|\cdot\|_{2,r}$, the rows of $\hat{\mathbf{O}}$ are separately given by

$$\hat{\mathbf{o}}_n := \arg \min_{\mathbf{o}_n \in \mathbb{R}^p} \left[ \|\hat{\mathbf{r}}_n - \mathbf{o}_n\|_2^2 + \lambda_2 \|\mathbf{o}_n\|_2 \right], \quad n = 1, \ldots, N. \tag{29}$$

For each $n = 1, \ldots, N$, because (29) is nondifferentiable at the origin one should consider two cases: i) if $\hat{\mathbf{o}}_n = \mathbf{0}_p$, it follows that the minimum cost in (29) is $\|\hat{\mathbf{r}}_n\|_2^2$; otherwise, ii) if $\|\hat{\mathbf{o}}_n\|_2 > 0$, the first-order condition for optimality gives $\hat{\mathbf{o}}_n = \hat{\mathbf{r}}_n - (\lambda_2/2)\hat{\mathbf{r}}_n/\|\hat{\mathbf{r}}_n\|_2$ provided $\|\hat{\mathbf{r}}_n\|_2 > \lambda_2/2$, and the minimum cost is $\lambda_2\|\hat{\mathbf{r}}_n\|_2 - \lambda_2^2/4$. Compactly, the solution of (29) is given by $\hat{\mathbf{o}}_n = \hat{\mathbf{r}}_n(\|\hat{\mathbf{r}}_n\|_2 - \lambda_2/2)_+/\|\hat{\mathbf{r}}_n\|_2$ , while the minimum cost in (29) after minimizing w.r.t. $\mathbf{o}_n$ is $\rho_v(\hat{\mathbf{r}}_n)$ [cf. (9) and the argument following (29)]. The conclusion is that $\hat{\mathcal{V}}$ is the minimizer of (8), in addition to being the solution of (7) by definition.

## References

[1] D. Angelosante, J. A. Bazerque, and G. B. Giannakis, "Online adaptive estimation of sparse signals: Where RLS meets the $\ell_1$-norm," *IEEE Trans. Signal Process*, vol. 58, pp. 3436–3447, Jul. 2010.

[2] L. Balzano, R. Nowak, and B. Recht, "Online identification and tracking of subspaces from highly incomplete information," in *Proc. of 48th Allerton Conference*, Monticello, IL, Sep./Oct. 2010, pp. 704–711.

[3] D. P. Bertsekas, *Nonlinear Programming*, 2nd ed. Athena-Scientific, 1999.

[4] N. A. Campbell, "Robust procedures in multivariate analysis i: Robust covariance estimation," *Applied Stat.*, vol. 29, pp. 231–237, 1980.

[5] E. J. Candes, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis?" *Journal of the ACM*, vol. 58, Article No. 11, Mar. 2011.

[6] E. J. Candes, M. B. Wakin, and S. Boyd, "Enhancing sparsity by reweighted $\ell_1$ minimzation," *Journal of Fourier Analysis and Applications*, vol. 14, pp. 877–905, Dec. 2008.

[7] V. Chandrasekaran, S. Sanghavi, P. A. Parillo, and A. S. Willsky, "Rank-sparsity incoherence for matrix decomposition," *SIAM Journal on Optimization*, vol. 21, pp. 572–596, 2011.

[8] F. de la Torre and M. J. Black, "A framework for robust subspace learning," *Int. Jrnl. of Computer Vision*, vol. 54, pp. 1183–209, 2003.

[9] X. Ding, L. He, and L. Carin, "Bayesian robust principal component analysis," *IEEE Trans. Image Process.*, vol. 20, 2011.

[10] J. Fan and R. Li, "Variable selection via nonconcave penalized likelihood and its oracle properties," *J. Amer. Stat. Assoc.*, vol. 96, pp. 1348–1360, 2001.

[11] M. Fazel, H. Hindi, and S. Boyd, "Log-det heuristic for matrix rank minimization with applications to Hankel and Euclidean distance matrices," in *Proc. of the American Control Conf.*, Denver, CO, Jun. 2003, pp. 2156–2162.

[12] P. Forero, V. Kekatos, and G. B. Giannakis, "Outlier-aware robust clustering," in *Proc. of Intl. Conf. on Acoustics, Speech and Signal Processing*, Prague, Czech Republic, May 2011, pp. 2244–2247.

[13] J. J. Fuchs, "An inverse problem approach to robust regression," in *Proc. of Intl. Conf. on Acoustics, Speech and Signal Processing*, Phoeniz, AZ, Mar. 1999, pp. 180–188.

[14] G. B. Giannakis, G. Mateos, S. Farahmand, V. Kekatos, and H. Zhu, "USPACOR: Universal sparsity-controlling outlier rejection," in *Proc. of Intl. Conf. on Acoust., Speech and Signal Proc.*, Prague, Czech Republic, May 2011, pp. 1952–1955.

[15] M. Girvan and M. E. J. Newman, "Community structure in social and biological networks," *Proc. Natl. Acad. Sci. USA*, vol. 99, pp. 7821–7826, 2002.

[16] L. R. Goldberg, "The Eugene-Springfield community sample: Information available from the research participants," Oregon Research Institue, Tech. Rep. vol. 48, no. 1, 2008.

[17] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, 2nd ed. Springer, 2009.

[18] J. He, L. Balzano, and J. C. S. Lui, "Online robust subspace tracking from partial information," 2011, see also arXiv:1109.3827v2 [cs.IT].

[19] P. J. Huber and E. Ronchetti, *Robust Statistics*. New York: Wiley, 2009.

[20] Y. Jin and B. D. Rao, "Algorithms for robust linear regression by exploiting the connection to sparse signal recovery," in *Proc. of Intl. Conf. on Acoustics, Speech and Signal Processing*, Dallas, TX, Mar. 2010, pp. 3830–3833.

[21] O. P. John, L. P. Naumann, and C. J. Soto, "Paradigm shift to the integrative big-five trait taxonomy: History, measurement, and conceptual issues," in *Handbook of personality: Theory and research*, O. P. John, R. W. Robins, and L. A. Pervin, Eds. New York, NY: Guilford Press, 2008.

[22] I. T. Jolliffe, *Principal Component Analysis*.   New York: Springer, 2002.

[23] V. Kekatos and G. B. Giannakis, "From sparse signals to sparse residuals for robust sensing," *IEEE Trans. on Signal Processing*, vol. 59, pp. 3355–3368, Jul. 2011.

[24] J. Mairal, J. Bach, J. Ponce, and G. Sapiro, "Online learning for matrix factorization and sparse coding," *Jrnl. of Machine Learning Research*, vol. 11, pp. 19–60, Jan. 2010.

[25] M. Mardani, G. Mateos, and G. B. Giannakis, "Unveiling network anomalies in large-scale networks via sparsity and low rank," in *Proc. of 44th Asilomar Conf. on Signals, Systems, and Computers*, Pacific Grove, CA, Nov. 2011.

[26] B. K. Natarajan, "Sparse approximate solutions to linear systems," *SIAM J. Comput.*, vol. 24, pp. 227–234, 1995.

[27] A. T. Puig, A. Wiesel, and A. O. Hero, "Multidimensional shrinkage-thresholding operator and group LASSO penalties," *IEEE Signal Process. Letters*, vol. 18, pp. 363–366, Jun. 2011.

[28] I. Ramirez, F. Lecumberry, and G. Sapiro, "Universal priors for sparse modeling," in *Proc. of 3rd Intl. Workshop on Comp. Advances in Multi-Sensor Adapt. Process.*, Aruba, Dutch Antilles, Dec. 2009, pp. 197–200.

[29] B. Recht, M. Fazel, and P. A. Parrilo, "Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization," *SIAM Rev.*, vol. 52, pp. 471–501, 2010.

[30] B. Recht and C. Re, "Parallel stochastic gradient algorithms for large-scale matrix completion," 2011, (submitted). [Online]. Available: http://pages.cs.wisc.edu/~brecht/papers/11.Rec.Re.IPGM.pdf

[31] S. P. Reise and N. G. Waller, "Traitedness and the assessment of response pattern scalability," *Journal of Personality and Social Psychology*, vol. 65, pp. 143–151, 1993.

[32] P. J. Rousseeuw and A. M. Leroy, *Robust regression and outlier detection*.   New York: Wiley, 1987.

[33] B. Schlkopf, A. Smola, and K.-R. Mller, "Kernel principal component analysis," *Artificial Neural Networks: Lec. Notes in Computer Science*, vol. 1327, pp. 583–588, 1997.

[34] I. Tošić and P. Frossard, "Dictionary learning," *IEEE Signal Process. Mag.*, vol. 28, pp. 27–38, Mar. 2010.

[35] J. Tropp, "Just relax: Convex programming methods for identifying sparse signals," *IEEE Trans. on Information Theory*, vol. 51, pp. 1030–1051, Mar. 2006.

[36] P. Tseng, "Convergence of block coordinate descent method for nondifferentiable maximization," *J. Optim. Theory Appl.*, vol. 109, pp. 473–492, 2001.

[37] N. Waller and S. Reise, "Measuring psychopathology with non-standard IRT models: Fitting the four parameter model to the MMPI," in *New Directions in Psychological Measurement with Model-Based Approaches*, S. Embretson, Ed.   Washington, DC: Amer. Psych. Assoc., 2010.

[38] H. Xu, C. Caramanis, and S. Sanghavi, "Robust PCA via outlier pursuit," 2010, see also arXiv:1010.4237v2 [cs.LG].

[39] L. Xu and A. L. Yuille, "Robust principal component analysis by self-organizing rules based on statistical physics approach," *IEEE Trans. Neural Nets.*, vol. 6, pp. 131–143, Jan. 1995.

[40] B. Yang, "Projection approximation subspace tracking," *IEEE Trans. Sig. Proc.*, vol. 43, pp. 95–107, Jan. 1995.

[41] M. Yuan and Y. Lin, "Model selection and estimation in regression with grouped variables," *J. Royal. Statist. Soc B*, vol. 68, pp. 49–67, 2006.

[42] Z. Zhou, X. Li, J. Wright, E. Candes, and Y. Ma, "Stable principal component pursuit," in *Proc. of Intl. Symp. on Information Theory*, Austin, TX, Jun. 2010, pp. 1518–1522.

[43] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," *Jrnl. of Comp. and Graphical Statistics*, vol. 15, no. 2, pp. 265–286, 2006.

TABLE I

| $\sigma_e^2$ | $\lambda_2^*$ in (7) | e̅rr for (7) (refined) | e̅rr for (14) | e̅rr for PCA |
|---|---|---|---|---|
| 0.01 | 0.7142 | 0.0622 | 0.0682 | 0.4679 |
| 0.05 | 1.7207 | 0.1288 | 0.1519 | 1.0122 |
| 0.1 | 2.4348 | 0.1742 | 0.2150 | 1.4141 |
| 0.25 | 3.6084 | 0.2525 | 0.3403 | 2.2480 |
| 0.5 | 6.1442 | 0.3361 | 0.4783 | 3.1601 |



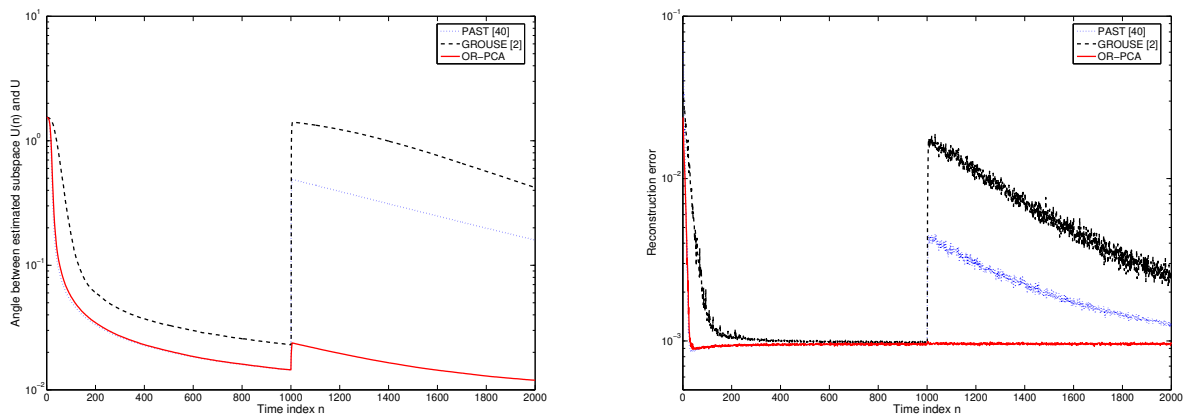Fig. 1.   Pseudo scree plot of outlier size ($\|\hat{\mathbf{o}}_n\|_2$); the 100 largest outliers are shown.



Fig. 2.   (Left) Time evolution of the angle between the learnt subspace $\mathbf{U}(n)$, and the true $\mathbf{U}$ used to generate the data ($\beta = 0.99$ and $\lambda_2 = 1.65$). Outlier contaminated data is introduced at time $n = 1001$. (Right) Time evolution of the reconstruction error.
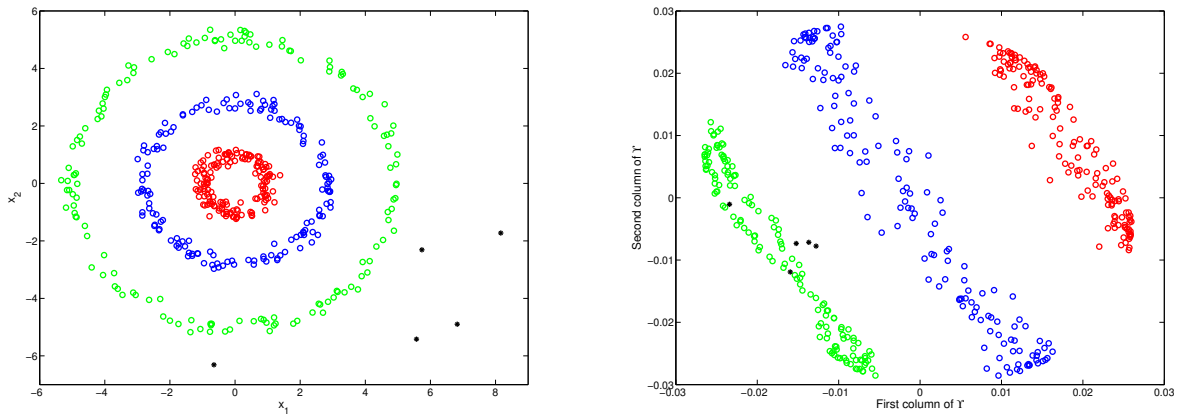
Fig. 3. (Left) Data in three concentric clusters, in addition to five outliers shown in black. (Right) Coordinates of the first two columns of $\Upsilon$, obtained by running Algorithm 4. The five outlying points are correctly identified, and thus can be discarded. Non-robust methods will assign them to the green cluster.
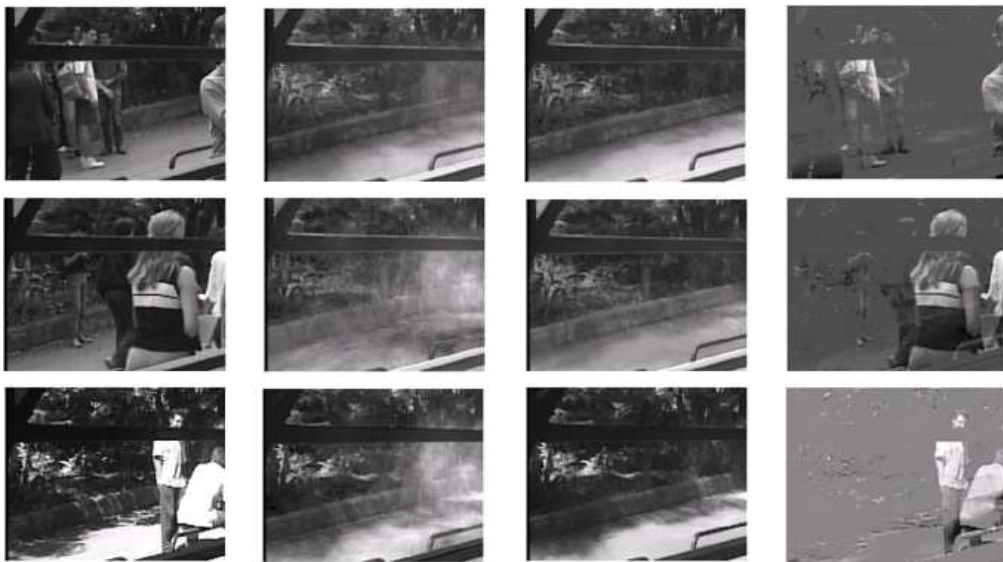


Fig. 4. Background modeling for video surveillance. First column: original frames. Second column: PCA reconstructions, where the presence of undesirable 'ghostly' artifacts is apparent, since PCA is not able to completely separate the people from the background. Third column: robust PCA reconstructions, which recover the illumination changes while successfully subtracting the people. Fourth column: outliers in $\hat{\mathbf{o}}$, which mostly capture the people and abrupt changes in illumination.
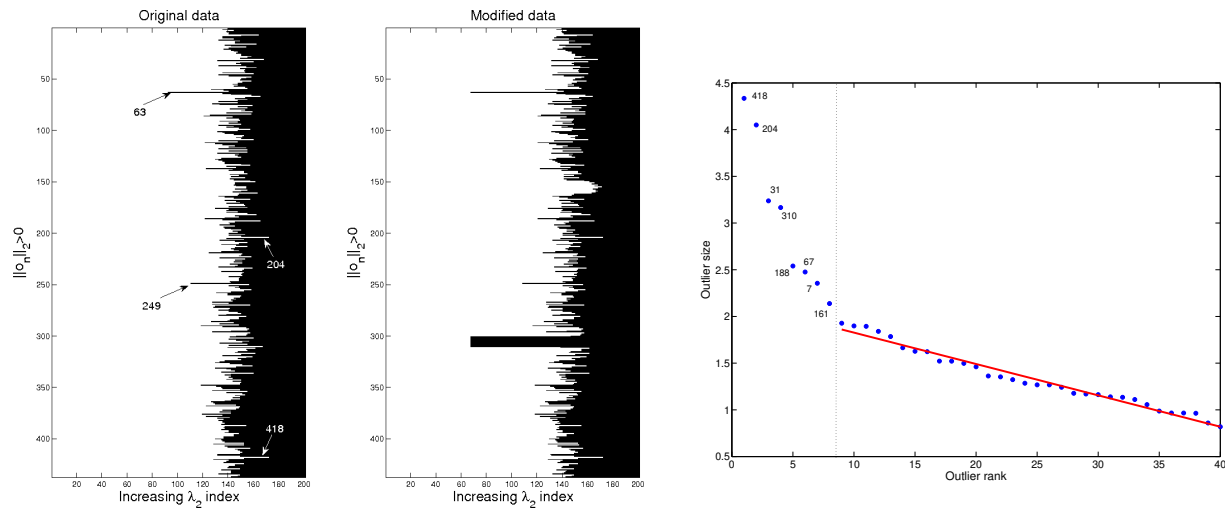
Fig. 5. (Left) Evolution of $\hat{\mathbf{O}}$'s row support as a function of $\lambda_2$ – black pixels along the $n$th row indicate that $\|\hat{\mathbf{o}}_n\|_2 = 0$, whereas white ones reflect that the responses from subject $n$ are deemed as outliers for given $\lambda_2$. The results for both the original and modified (introducing random and constant item responses) BFI datasets are shown. (Right) Pseudo scree plot of outlier size ($\|\hat{\mathbf{o}}_n\|_2$); the 40 largest outliers are shown. Robust PCA declares the largest 8 as aberrant responses.
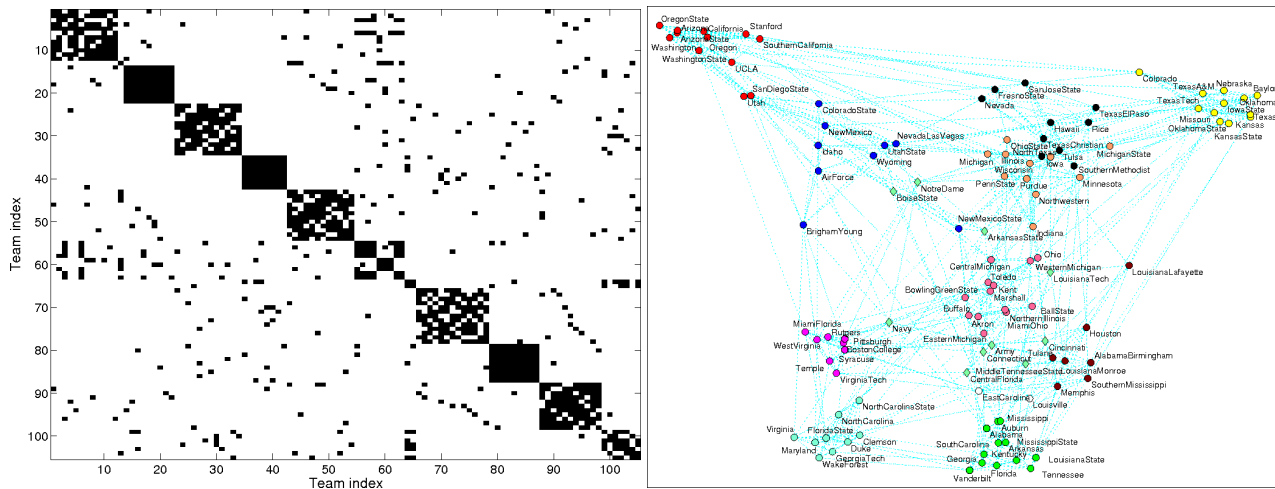


Fig. 6. (Left) Entries of $\mathbf{K}$ after removing the outliers, where rows and columns are permuted to reveal the clustering structure found by robust KPCA. (Right) Graph depiction of the clustered network. Teams belonging to the same estimated conference (cluster) are colored identically. The outliers are represented as diamond-shaped nodes.