

Provable Self-Representation Based Outlier Detection in a Union of Subspaces

Chong You, Daniel P. Robinson, René Vidal
Johns Hopkins University, Baltimore, MD, 21218, USA

Abstract

Many computer vision tasks involve processing large amounts of data contaminated by outliers, which need to be detected and rejected. While outlier detection methods based on robust statistics have existed for decades, only recently have methods based on sparse and low-rank representation been developed along with guarantees of correct outlier detection when the inliers lie in one or more low-dimensional subspaces. This paper proposes a new outlier detection method that combines tools from sparse representation with random walks on a graph. By exploiting the property that data points can be expressed as sparse linear combinations of each other, we obtain an asymmetric affinity matrix among data points, which we use to construct a weighted directed graph. By defining a suitable Markov Chain from this graph, we establish a connection between inliers/outliers and essential/inessential states of the Markov chain, which allows us to detect outliers by using random walks. We provide a theoretical analysis that justifies the correctness of our method under geometric and connectivity assumptions. Experimental results on image databases demonstrate its superiority with respect to state-of-the-art sparse and low-rank outlier detection methods.

1. Introduction

In many applications in computer vision, including motion estimation and segmentation [19] and face recognition [2], high-dimensional datasets can be well approximated by a union of low-dimensional subspaces. Such applications have motivated a lot of research on the problems of learning one or more subspaces from data, a.k.a. subspace learning and subspace clustering, respectively. In practice, datasets are often contaminated by points that do not lie in the subspaces, i.e. outliers. In such situations, it is often essential to detect and reject these outliers before any subsequent processing/analysis is performed.

Prior work. We address the problem of outlier detection in the setting when the inlier data are assumed to lie close to a union of unknown low-dimensional subspaces (low relative to the dimension of the ambient space). A traditional method for solving this problem is RANSAC [13], which is

based on randomly selecting a subset of points, fitting a subspace to them, and counting the number of points that are well fit by this subspace; this process is repeated for sufficiently many trials and the best fit is chosen. RANSAC is intrinsically combinatorial and the number of trials needed to find a good estimate of the subspace grows exponentially with the subspace dimension. Consequently, the methods of choice have been to robustly learn the subspaces by penalizing the sum of *unsquared* distances (in lieu of *squared* distances used in classical methods such as PCA) of points to the closest subspace [9, 22, 62, 61]. Such a penalty is robust to outliers because it reduces the contributions from large residuals arising from outliers. However, the optimization problem is usually nonconvex and a good initialization is extremely important for finding the optimal solution.

The groundbreaking work of Wright et al. [54] and Candès et al. [4] on using convex optimization techniques to solve the PCA problem with robustness to corrupted entries has led to many recent methods for PCA with robustness to outliers [55, 32, 25, 60, 21]. For example, Outlier Pursuit [55] uses the nuclear norm $\|\cdot\|_*$ to seek low-rank solutions by solving the problem $\min_L \|X - L\|_{2,1} + \lambda \|L\|_*$ for some $\lambda > 0$. A prominent advantage of convex optimization techniques is that they are guaranteed to correctly identify outliers under certain conditions. Very recently, several non-convex outlier detection methods have also been developed with guaranteed correctness [20, 6]. Nonetheless, these methods typically model a *unique* inlier subspace, e.g., by a low rank matrix L in Outlier Pursuit, and therefore cannot deal with multiple inlier subspaces since the union of multiple subspaces could be high-dimensional.

Another class of methods with theoretical guarantees for correctness utilizes the fact that outliers are expected to have low similarities with other data points. In [5, 1], a multi-way similarity is introduced that is defined from the polar curvature, which has the advantage of exploiting the subspace structure. However, the number of combinations in multi-way similarity can be prohibitively large. Some recent works have explored using inner products between data points for outlier detection [17, 40]. Although computationally very efficient, these methods require the inliers to be well distributed and densely sampled within the subspaces.

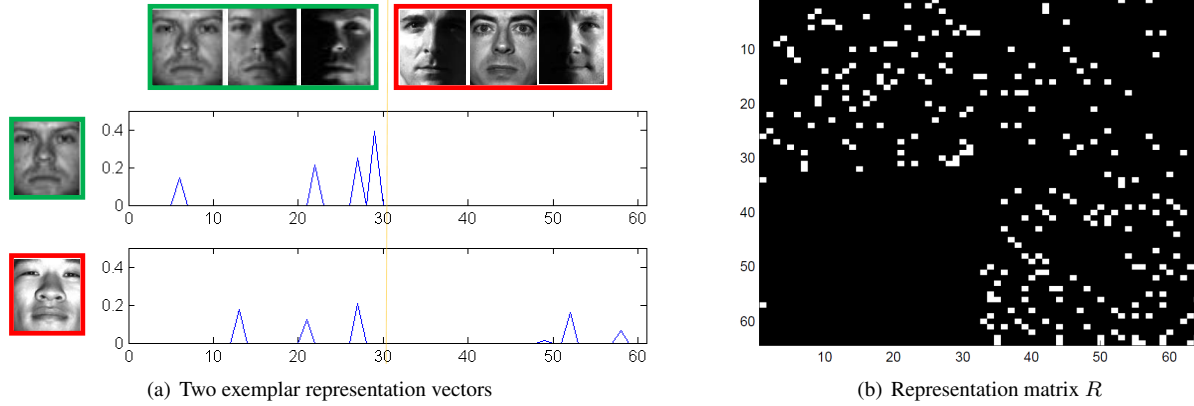


Figure 1. An illustration of a self-representation matrix R in the presence of outliers. The first 32 columns of the data matrix X correspond to 32 images of one individual under different illuminations from the Extended Yale B database, and the next 32 images are randomly chosen from all other individuals; three examples from each category are shown near the top of 1(a). We also show a typical representation vector for an inlier and an outlier image in 1(a), and the complete representation matrix R in 1(b), where white and black denote $r_{ij} \neq 0$ and $r_{ij} = 0$. Notice that inliers use only other inliers in their representation, while outliers use both inliers and outliers in their representations.

Overview of our method and contributions. In this work, we address the problem of outlier detection by using data self-representation. The proposed approach builds on the self-expressiveness property of data in a union of low-dimensional subspaces, originally introduced in [11], which states that a point in a subspace can always be expressed as a linear combination of other points in the subspace. In particular, if the columns of $X = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ lie in multiple subspaces, then for all $j = 1, \dots, N$, there exists a vector $\mathbf{r}_j \in \mathbb{R}^N$ such that $\mathbf{x}_j = X\mathbf{r}_j$ and the nonzero entries of \mathbf{r}_j correspond to points in the same subspace as \mathbf{x}_j . If the subspace dimensions are small, \mathbf{r}_j can be taken to be sparse and be computed by solving the ℓ_1 minimization problem

$$\min_{\mathbf{r}_j} \|\mathbf{r}_j\|_1 + \frac{\gamma}{2} \|\mathbf{x}_j - X\mathbf{r}_j\|_2^2 \quad \text{s.t. } r_{jj} = 0 \quad (1)$$

for some $\gamma > 0$. In [11], an *undirected* graph is constructed from $R = [\mathbf{r}_1, \dots, \mathbf{r}_N]$ in which each vertex corresponds to a data point, and vertices corresponding to \mathbf{x}_i and \mathbf{x}_j are connected if either r_{ij} or r_{ji} is nonzero. Such a graph can be used to segment the data into their respective subspaces by applying spectral clustering [48] to the graph’s Laplacian.

Consider now the case where X contains outliers to the subspaces. Figure 1 illustrates an example representation matrix R computed from (1) for data drawn from a **single** subspace (face images from one individual) plus outliers (other images). In this case, the representation R is such that inliers express themselves as linear combinations of a few other inliers, while outliers express themselves as linear combinations of both inliers and outliers. Motivated by this observation, we use a *directed* graph to model data relations: a directed edge from \mathbf{x}_j to \mathbf{x}_i indicates that \mathbf{x}_j uses \mathbf{x}_i in its representation (i.e. $r_{ij} \neq 0$). Then a random walk on the representation graph initialized at an outlier will

not return to the set of outliers since once the random walk reaches an inlier it cannot return to the outliers. Therefore, we design a random walk process and identify outliers as those whose probabilities tend to zero. Our work makes the following contributions with respect to the state of the art:

1. Our method can detect outliers using the probability distribution of a *random walk* on a graph constructed from *data self-representation*.
2. Our *data self-representation* allows our method to handle multiple inlier subspaces. Knowledge of the number of subspaces and their dimensions is not required, and the subspaces may have a nontrivial intersection.
3. Our method can explore contextual information by using a *random walk*, i.e., the “outlierness” of a particular point depends on the “outlierness” of its neighbors.
4. Our analysis shows that our method correctly identifies outliers under suitable assumptions on the data distribution and connectivity of the representation graph.
5. Experiments on real image databases illustrate the effectiveness of our method.

2. Related work

Outlier detection by self-representation. Prior work has explored using data self-representation as a tool for outlier detection in a union of subspaces. Specifically, motivated by the observation that outliers do not have *sparse* representations, [43, 8] declare a point \mathbf{x}_j as an outlier if $\|\mathbf{r}_j\|_1$ is above a threshold. However, this ℓ_1 -thresholding strategy is not robust to outliers that are close to each other since their representation vectors may have small ℓ_1 -norms. The LRR [28] solves for a low-rank self-representation matrix R in lieu of a sparse representation and penalizes the sum of unsquared self-representation errors $\|\mathbf{x}_j - X\mathbf{r}_j\|_2$, which

makes it more robust to outliers. However, LRR requires the subspaces to be independent and the sum of the union of subspaces to be low-dimensional [29].

Outlier detection by maximum consensus. In a diverse range of contexts such as maximum consensus [63, 7] and robust linear regression [33, 49], people have studied problems of the form

$$\min_{\mathbf{b}} \sum_{i=1}^N \mathbb{I}(|\mathbf{x}_i^\top \mathbf{b} - y_i| \geq \epsilon), \quad (2)$$

in which $\mathbb{I}(\cdot)$ is the indicator function. Note that if we set $y_i = 1$ for all i , then (2) can be interpreted as detecting outliers in data X where the inliers lie close to an *affine* hyperplane. A problem closely related to (2) is

$$\min_{\mathbf{b}} \sum_{i=1}^N \mathbb{I}(|\mathbf{x}_i^\top \mathbf{b}| \geq \epsilon) \text{ s.t. } \mathbf{b} \neq 0, \quad (3)$$

which appears in many applications (e.g. see [39]). In particular, (3) can be used to learn a *linear* hyperplane from data corrupted by outliers. To detect outliers in a general low-dimensional subspace, one can apply (2) and (3) recursively to find a basis for the orthogonal complement of the subspace [46]. However, such an approach is limited because there can be only one inlier subspace and the dimension of that subspace must be known in advance.

Outlier detection by random walk. Perhaps the most well-known random walk based algorithm is PageRank [3]. Originally introduced to determine the authority of website pages from web graphs, PageRank and its variants have been used in different contexts for ranking the centrality of the vertices in a graph. In particular, [34, 35] propose the OutRank, which ranks the “outlierness” of points in a dataset by applying PageRank to an undirected graph in which the weight of an edge is the cosine similarity or RBF similarity between the two connected data points. Then, points that have low centrality are regarded as outliers. The outliers returned by OutRank are those that have low similarity to other data points. Therefore, OutRank does not work if points in a subspace are not dense enough.

3. Outlier detection by self-representation

In this section, we present our data self-representation based outlier detection method. We first describe the data self-representation and its associated properties for inliers and outliers. We then design a random walk algorithm on the representation graph whose limiting behavior allows us to identify the sets of inliers and outliers.

3.1. Data self-representation

Given an unlabeled dataset $X = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ containing inliers and outliers, the first step of our algorithm is

to construct the data self-representation matrix denoted by $R = [\mathbf{r}_1, \dots, \mathbf{r}_N]$. As briefly discussed in the introduction (see also Figure 1), a self-representation matrix R computed from (1) is observed to have different properties for inliers and outliers. Specifically, inliers usually use only other inliers for self-representation, i.e. for an inlier \mathbf{x}_j , the representation is such that $r_{ij} \neq 0$ only if \mathbf{x}_i is also an inlier, where r_{ij} is the (i, j) -th entry of R . This property is expected to hold if the inliers lie in a union of low dimensional subspaces, as evidenced from the works [12, 43, 59, 52, 50]. As an intuitive explanation, if the inliers lie in a low dimensional subspace, then any inlier has a *sparse* representation using other points in this subspace. Thus such a representation can be found by using sparsity-inducing regularization as seen in (1). In contrast, outliers are generally randomly distributed in the ambient space, so that a self-representation usually contains both inliers and outliers.

Since the representation R computed from (1) is sparse, there are potentially connectivity issues in the representation graph, i.e. an inlier that is not well-connected to other inliers may be detected as an outlier, and an outlier that is not well connected may be detected as an inlier. To address the connectivity issue, we compute the data self-representation matrix R by the elastic net problem [64, 56]:

$$\min_{\mathbf{r}_j} \lambda \|\mathbf{r}_j\|_1 + \frac{1-\lambda}{2} \|\mathbf{r}_j\|_2^2 + \frac{\gamma}{2} \|\mathbf{x}_j - X\mathbf{r}_j\|_2^2 \text{ s.t. } r_{jj} = 0, \quad (4)$$

in which $\lambda \in [0, 1]$ controls the balance between sparseness (via ℓ_1 regularization) and connectivity (via ℓ_2 regularization). Specifically, if λ is chosen close to 1, we can still expect that the computed representation for an inlier will only use inliers. The ℓ_2 regularization has been introduced to promote more connections between data points, i.e. if $\lambda \in [0, 1)$, then one expects more nonzero entries in R . A detailed discussion of the representation computed from (4) and the connectivity issue is provided in Section 4.

3.2. Representation graph and random walk

We use a directed graph G , which we call a *representation graph*, to capture the behavior of inliers and outliers from the representation matrix R . The vertices of G correspond to the data points X , and the edges are given by the (weighted) adjacency matrix $A := |R|^\top \in \mathbb{R}^{N \times N}$ with the absolute value taken elementwise, i.e., the weight of the edge from \mathbf{x}_i to \mathbf{x}_j is given by $a_{ij} = |r_{ji}|$. In the representation graph, we expect that vertices corresponding to inliers will have edges that only lead to inliers, while vertices that are outliers will have edges that lead to both inliers and outliers. In other words, we do not expect to have any edges that lead from an inlier to an outlier.

Using the previous paragraph as motivation, we design a random walk procedure to identify the outliers. A random walk on the representation graph G is a discrete time

Markov chain, for which the transition probability from x_i at a given time to x_j at the next time is given by $p_{ij} := a_{ij}/d_i$ with $d_i := \sum_j a_{ij}$. By this definition, if the starting point of a random walk is an inlier then it will never escape the set of inliers as there is no edge going from any inlier to any outlier. In contrast, a random walk starting from an outlier will likely end up in an inlier state since once it enters any inlier it will never return to an outlier state. Thus, by using different data points to initialize random walks, outliers can be identified by observing the final probability distribution of the state of the random walks (see Figure 2).

If $P \in \mathbb{R}^{N \times N}$ is the transition matrix with entries p_{ij} , then P is related to the representation matrix R by

$$p_{ij} = |r_{ji}| / \|\mathbf{r}_i\|_1 \text{ for all } \{i, j\} \subset \{1, 2, \dots, N\}. \quad (5)$$

We define $\pi^{(t)} = [\pi_1^{(t)}, \dots, \pi_N^{(t)}]$ to be the state probability distribution at time t , then the state transition is given by $\pi^{(t+1)} = \pi^{(t)} P$. Thus, a t -step transition is $\pi^{(t)} = \pi^{(0)} P^t$ with $\pi^{(0)}$ the chosen initial state probability distribution.

3.3. Main algorithm: Outlier detection by R-graph

We propose to perform outlier detection by using random walks on the representation graph G . We set the initial probability distribution as $\pi^{(0)} = [1/N, \dots, 1/N]$, and then compute the t -step transition $\pi^{(t)} = \pi^{(0)} P^t$. This can be interpreted as initializing a random walk from each of the N data points, and then finding the sum of probability distributions of all random walks after t steps. It is expected that all random walks—starting from either an inlier or an outlier—will eventually have high probabilities for the inlier states and low probabilities for the outlier states.

We note that the $\pi^{(t)}$ defined as above need not converge, as shown by the 2-dimensional example $P = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$. Instead, we choose to use the T -step Cesàro mean, given by

$$\bar{\pi}^{(T)} = \frac{1}{T} \sum_{t=1}^T \pi^{(0)} P^t \equiv \frac{1}{T} \sum_{t=1}^T \pi^{(t)}, \quad (6)$$

which is the average of the first T t-step probability distributions (see Figure 2). The sequence $\{\bar{\pi}^{(T)}\}$ has the benefit that it always converges, and its limit is the same as that of $\pi^{(t)}$ whenever the latter exists. In the next section, we give a more detailed discussion of this choice, its properties for outlier detection, and its convergence behavior.

Our complete algorithm is stated as Algorithm 1.

4. Theoretical guarantees for correctness

Let us first formally define the problem of outlier detection when data is drawn from a union of subspaces.

Problem 4.1 (Outlier detection in a union of subspaces). *Given data $X = [x_1, \dots, x_N] \in \mathbb{R}^{D \times N}$ whose columns*

Algorithm 1 Outlier detection by representation graph

Input: Data $X = [x_1, \dots, x_N]$, #iterations T , threshold ϵ .

- 1: Use X to solve for $R = [r_1, \dots, r_N]$ using (4).
- 2: Compute P from R using (5).
- 3: Initialize $t = 0$, $\pi = [1/N, \dots, 1/N]$, and $\bar{\pi} = \mathbf{0}$.
- 4: **for** $t = 1, 2, \dots, T$ **do**
- 5: Compute $\pi \leftarrow \pi \cdot P$, and then set $\bar{\pi} \leftarrow \bar{\pi} + \pi$.
- 6: **end for**
- 7: $\bar{\pi} \leftarrow \bar{\pi}/T$.

Output: An indicator of outliers: x_j is an outlier if $\bar{\pi}_j \leq \epsilon$.

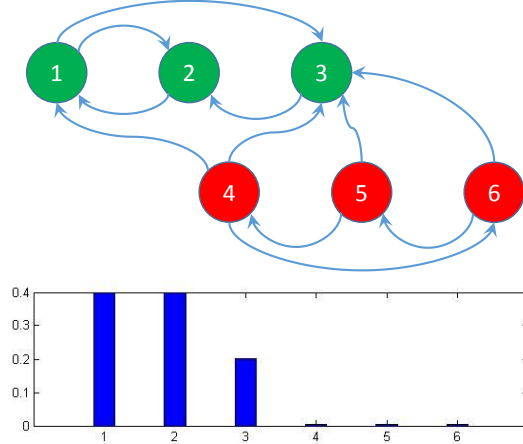


Figure 2. Illustration of random walks on a representation graph. Top: green balls represent inliers and red balls represent outliers, and arrows represent edges among nodes. Notice that there is no edge going from inliers to outliers. A random walk starting from any point will end up at only inlier points. Bottom: bar plot of $\bar{\pi}^{(100)}$ with the i th bar corresponding to the i th entry in $\bar{\pi}^{(100)}$. The use of thresholding on this probability distribution will correctly distinguish outliers from inliers.

contain inliers that are drawn from an unknown number of unknown subspaces $\{\mathcal{S}_\ell\}_{\ell=1}^n$, and outliers that are outside of $\cup_{\ell=1}^n \mathcal{S}_\ell$, the goal is to identify the set of outliers.

Recall that motivation for our method is that ideally there will be no edge going from an inlier to an outlier in the representation graph. This motivates us to assume that a random walk starting at any inlier will eventually return to itself, i.e. inliers are *essential states* of the Markov chain, while outliers are those that have a chance of never coming back to itself, i.e. outliers are *inessential states*. Formally, we work with a (time homogeneous) Markov chain with state space $\Omega = \{1, \dots, N\}$, in which each state j corresponds to data x_j , and the transition probability P is given by (5). Given $\{i, j\} \subset \Omega$, we say that j is accessible from i , denoted as $i \rightarrow j$, if there exists some $t > 0$ such that the (i, j) -th entry of P^t is positive. Intuitively, $i \rightarrow j$ if a random walk can move from i to j in finitely many steps.

Definition 4.1 (Essential and inessential state [23]). *A state $i \in \Omega$ is essential if for all j such that $i \rightarrow j$ it is also true that $j \rightarrow i$. A state is inessential if it is not essential.*

Our aim in this section is to establish that if inliers connect to themselves, i.e. they are *subspace-preserving* (Section 4.1), and the representation R satisfies certain connectivity conditions (Section 4.2), then inliers are essential states of the Markov chain and outliers are inessential states. Subsequently, in Section 4.3 we show that the Cesàro mean (6) identifies essential and inessential states, thus establishing the correctness of Algorithm 1 for outlier detection.

4.1. Subspace-preserving representation

We first establish that inliers express themselves with only other inliers when they lie in a union of low dimensional subspaces. This property is well-studied in the subspace clustering literature. We will borrow terminologies and results from prior work and modify them for our current task of outlier detection.

Definition 4.2 (Subspace-preserving representation [47]). *If $\mathbf{x}_j \in \mathcal{S}_\ell$ is an inlier, then the representation $\mathbf{r}_j \in \mathbb{R}^N$ is called subspace-preserving if the nonzero entries of \mathbf{r}_j correspond to points in \mathcal{S}_ℓ , i.e. $r_{ij} \neq 0$ only if $\mathbf{x}_i \in \mathcal{S}_\ell$. The representation matrix $R = [\mathbf{r}_1, \dots, \mathbf{r}_N] \in \mathbb{R}^{N \times N}$ is called subspace-preserving if \mathbf{r}_j is subspace-preserving for every inlier \mathbf{x}_j .*

A representation matrix R is subspace-preserving if each inlier uses points in its own subspace for representation. Given X , a subspace-preserving representation R can be obtained by solving (4) when certain geometric conditions hold. The following result is modified from [56]. It assumes that columns of X are normalized to have unit ℓ_2 -norm.

Theorem 4.1. *Let $\mathbf{x}_j \in \mathcal{S}_\ell$ be an inlier. Define the oracle point of \mathbf{x}_j to be $\bar{\delta}_j := \gamma \cdot (\mathbf{x}_j - X_{-j}^\ell \cdot \mathbf{r}_j^\ell)$, where X_{-j}^ℓ is the matrix containing all points in \mathcal{S}_ℓ except \mathbf{x}_j and*

$$\mathbf{r}_j^\ell := \arg \min_{\mathbf{r}} \lambda \|\mathbf{r}\|_1 + \frac{1-\lambda}{2} \|\mathbf{r}\|_2^2 + \frac{\gamma}{2} \|\mathbf{x}_j - X_{-j}^\ell \mathbf{r}\|_2^2.$$

The solution \mathbf{r}_j to (4) is subspace-preserving if

$$\max_{k \neq j, \mathbf{x}_k \in \mathcal{S}_\ell} |\langle \mathbf{x}_k, \bar{\delta}_j \rangle| - \max_{k: \mathbf{x}_k \notin \mathcal{S}_\ell} |\langle \mathbf{x}_k, \bar{\delta}_j \rangle| > \frac{1-\lambda}{\lambda}, \quad (7)$$

where $\bar{\delta}_j := \delta_j / \|\delta_j\|_2$.

An outline of the proof is given in the appendix. Note that the oracle point $\bar{\delta}_j$ lies in \mathcal{S}_ℓ and that its definition only depends on points in \mathcal{S}_ℓ . The first term in condition (7) captures the distribution of points in \mathcal{S}_ℓ near $\bar{\delta}_j$, and is expected to be large if the neighborhood of $\bar{\delta}_j$ is well-covered by points from \mathcal{S}_ℓ . The second term characterizes the similarity between the oracle point $\bar{\delta}_j$ and all other data points,

which includes the outliers and the inliers from other subspaces. The condition requires the former to be larger than the latter by a margin of $\frac{1-\lambda}{\lambda}$, which is close to zero if λ is close to 1. Overall, condition (7) requires that points in \mathcal{S}_ℓ are dense around $\bar{\delta}_j$, which is itself in \mathcal{S}_ℓ , and that outliers and inliers from other subspaces do not lie close to $\bar{\delta}_j$.

Even if (7) holds for all j so that the representation R is subspace-preserving, we cannot automatically establish an equivalence between inliers/outliers and essential/inessential states because of potential complications related to the graph's *connectivity*. This is addressed next.

4.2. Connectivity considerations

In the context of sparse subspace clustering, the well-known connectivity issue [36, 53, 30, 56, 51] refers to the problem that points in the same subspace may not be well-connected in the representation graph, which may cause oversegmentation of the true clusters. Thus, one has to make the assumption that each true cluster is connected to guarantee correct clustering. For the outlier detection problem, it may happen that an inlier is inessential and thus classified as an outlier when the inliers are not well-connected; similarly, an outlier may be essential and thus classified as an inlier if it is not connected to at least one inlier. In fact, the situation is even more involved since the representation graph is directed and inliers and outliers behave differently.

Suppose, as a first example, that there exists an inlier that is never used to express any other inliers. This is equivalent to saying that there is no edge going into this point from any other inliers. Note that the subspace-preserving property can still hold if this inlier expresses itself using other inliers. Yet, since a random walk leaving this point would never return it can not be identified as an inlier. To avoid such cases, we need the following assumption.

Assumption 4.1. *For any inlier subspace \mathcal{S}_ℓ , the vertices $\{\mathbf{x}_j \in \mathcal{S}_\ell\}$ of the representation graph are strongly connected, i.e. there is a path in each direction between each pair of vertices.*

Assumption 4.1 requires good connectivity between points from the same inlier subspace. We also need good connectivity between outliers and inliers. Consider the example when there is a subset of outliers for which all of their outgoing edges lead only to points within that same subset. In this case, the subset of points can not be detected as outliers since their representation pattern is the same as for the inliers. The next assumption rules out this case.

Assumption 4.2. *For each subset of outliers there exists an edge in the representation graph that goes from a point in this subset to an inlier or to an outlier outside this subset.*

4.3. Main theorem: guaranteed outlier detection

We can now establish guaranteed outlier detection by our representation graph based method stated as Algorithm 1.

Theorem 4.2. *If the representation R is subspace-preserving and satisfies Assumptions 4.1 and 4.2, then Algorithm 1 with $T = \infty$ and $\epsilon = 0$ correctly identifies outliers.*

Theorem 4.2 is a direct consequence of the following two facts whose proofs are provided in the appendix.

Lemma 4.1. *If the representation R is subspace-preserving and Assumptions 4.1 and 4.2 hold, then inliers and outliers correspond to essential and inessential states, respectively.*

Lemma 4.2. *For any probability transition matrix P , the averaged probability distribution in (6) satisfies $\lim_{T \rightarrow \infty} \bar{\pi}^{(T)} = \pi$, where π is such that $\pi_j = 0$ if and only if state j is inessential.*

Theorem 4.2 shows that Problem 4.1 is solved by Algorithm 1 if the data X satisfies the geometric conditions in (7) and the representation graph satisfies the required connectivity assumptions.

We note that the random walk by the Cesàro mean adopted here is different from the popular random walk with restart as adopted by PageRank, for example. The benefit of PageRank is that the random walk converges to the unique stationary distribution. However, it is not clear whether this stationary distribution identifies the outliers. In fact, all states in the random walk of PageRank are essential, so that outliers do not converge to zero probabilities. In contrast, the random walk in our method does not necessarily have a unique stationary distribution, but the Cesàro mean does converge to one of the stationary distributions, which we have shown can be used to identify outliers. A detailed discussion is in the Appendix.

5. Experiments

We use several image databases (see Figure 3) to evaluate our outlier detection method (Algorithm 1). For computing the representation r_j in (4), we use the solver in [18] with $\lambda = 0.95$ and $\gamma = \alpha \cdot \frac{\lambda}{\max_{i: i \neq j} |x_j^\top x_i|}$, where α is a parameter tuned to each dataset. In particular, the solution to (4) is nonzero if and only if $\alpha > 1$. The number of iterations T is set to be 1,000.

5.1. Experimental setup

Databases. We construct outlier detection tasks from three publicly available databases. The Extended Yale B [15] dataset contains frontal face images of 38 individuals each under 64 different illumination conditions. The face images are of size 192×168 , for which we downsample to 48×42 . The Caltech-256 [16] is a database that contains

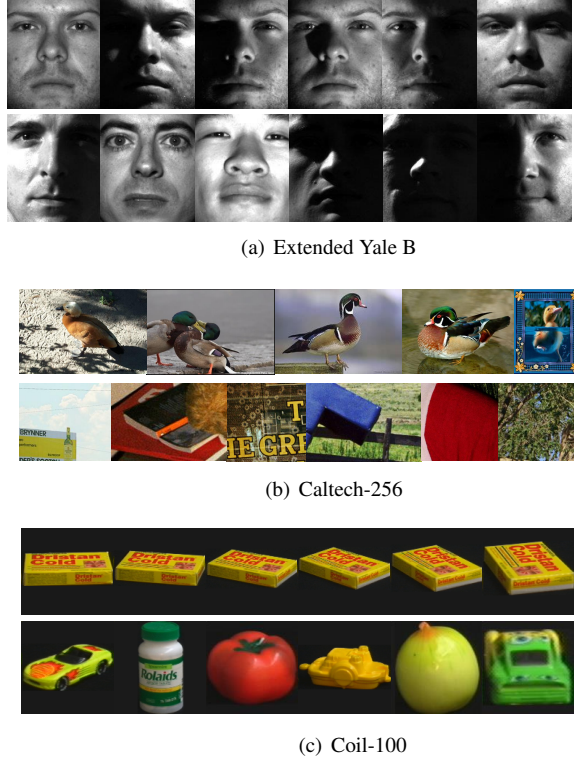


Figure 3. Examples of data used for outlier detection. For each database, the top row shows examples of the inlier set and the bottom row shows examples from the outlier set.

images from 256 categories that have more than 80 images each. There is also an additional “clutter” category in this database that contains 827 images of different varieties, which are used as outliers. The Coil-100 dataset [37] contains 7,200 images of 100 different objects. Each object has 72 images taken at pose intervals of 5 degrees, with the images being of size 32×32 . For the Extended Yale B and Coil-100 datasets we use raw pixel intensity as the feature representation. Images in Caltech-256 are represented by a 4,096-dimensional feature vector extracted from the last fully connected layer of the 16-layer VGG network [42].

Baselines. We compare with 6 other representative methods that are designed for detecting outliers in one or multiple subspaces: CoP [40], OutlierPursuit [55], REAPER [21], DPCP [46], LRR [28] and ℓ_1 -thresholding [43]. We also compare with a graph based method: OutRank [34, 35]. We implement the inexact ALM [26] for solving the optimization in OutlierPursuit. For LRR, we use the code available online at <https://sites.google.com/site/guangcanliu/>. For DPCP, we use the code provided by the authors. All other methods are implemented according to the description in their respective papers.

Evaluation metric. Each outlier detection method generates a numerical value for each data point that indicates its

Table 1. Results on the Extended Yale B database. Inliers are taken to be the images of either one or three randomly chosen subjects, and outliers are randomly chosen from the other subjects (at most one from each subject). For R-graph we set $\alpha = 5$ in the definition of γ .

	OutRank	CoP	REAPER	OutlierPursuit	LRR	DPCP	ℓ_1 -thresholding	R-graph (ours)
<i>Inliers: all images from one subject Outliers: 35%, taken from other subjects</i>								
AUC	0.536	0.556	0.964	0.972	0.857	0.952	0.844	0.986
F1	0.552	0.563	0.911	0.918	0.797	0.885	0.763	0.951
<i>Inliers: all images from three subjects Outliers: 15%, taken from other subjects</i>								
AUC	0.519	0.529	0.932	0.968	0.807	0.888	0.848	0.985
F1	0.288	0.292	0.758	0.856	0.509	0.653	0.545	0.878

“outlierness”, and a threshold value is required for determining inliers and outliers. A Receiver Operating Characteristic (ROC) curve plots the true positive rate and false positive rate for all threshold values. We use the area under the curve (AUC) as a metric of performance in terms of the ROC. The AUC is always between 0 and 1, with a perfect model having an AUC of 1 and a model that guesses randomly having an AUC of approximately 0.5.

As a second metric, we provide the F1-score, which is the harmonic mean of precision and recall. The F1-score is dependent upon the threshold, and we report the largest F1-score across all thresholds. An F1-score of 1 means there exists a threshold that gives both precision and recall equal to 1, i.e. a perfect separation of inliers and outliers.

The reported numbers for all experiments discussed in this section are the averages over 50 trials.

5.2. Outliers in face images

Suppose we are given a set of images of one or more individuals but that the data set is also corrupted by face images of a variety of other individuals. The task is to detect and remove those outlying face images. It is known that images of a face under different lighting conditions lie approximately in a low dimensional subspace. Thus, this task can be modeled as the problem of outlier detection in one subspace or in a union of subspaces.

We use the extended Yale B database. In the first experiment, we randomly choose a single individual from the 38 subjects and use all 64 images of this subject as the inliers. We then choose images from the remaining 37 subjects as outliers with at most one image from each subject. The overall data set has 25% outliers. The average AUC and F1 measures over 50 trials are reported in Table 1. For a fair comparison, we fine-tuned the parameters for all methods.

Comparing to state of the art. We see that our representation graph based method R-graph outperforms the other methods. Besides our method, the REAPER, Outlier Pursuit and DPCP algorithms all perform well. These three methods learn a single subspace and treat those that do not fit the subspace as outliers, thus making them well suited for this data (the images of one individual can be well-approximated by a single low dimensional subspace).

The LRR and ℓ_1 -thresholding methods use data self-representation, which is also the case for our method. However, LRR does not give good outlier detection results, probably because its algorithm for solving the LRR model is not guaranteed to converge to a global optimum. The ℓ_1 -thresholding also does not give good results, showing that the magnitude of the representation vector is not a robust measure for classifying outliers. By considering the connection patterns in the representation graph, our method achieves significantly better results.

The performance of OutRank and CoP is significantly worse than that of the other methods. This poor performance can be explained by the use of a coherence-based distance, which fails to capture similarity between data points when the data lie in subspaces. For example, it can be argued that the coherence between two faces with the same illumination condition can be higher than two images of the same face under different illumination conditions.

Dealing with multiple inlier groups. In order to test the ability of the methods to deal with multiple inlier groups, we designed a second experiment in which inliers are taken to be images of 3 randomly chosen subjects, and outliers are randomly drawn from other subjects as before. For all methods, we use the same parameters as in the previous experiment to test the robustness to parameter tuning. The results of this experiment are reported in Table 1.

We can see that Outlier Pursuit and our R-graph are the two best methods. Although Outlier Pursuit only models a single low dimensional subspace, it can still deal with this data since the union of the three subspaces corresponding to the three subjects in the inlier set is still low dimensional and can be treated as a single low dimensional subspace. However, we postulate that Outlier Pursuit will eventually fail as we increase the number of inlier groups, since the union of low dimensional subspaces will no longer be low rank. Our method does not have this limitation.

Similar to Outlier Pursuit, both REAPER and DPCP can, in principle, handle multiple inlier groups by fitting a single subspace to their union. However, REAPER and DPCP require as input the dimension of the union of the inlier subspaces, which can be hard to estimate in practice. Indeed, in Table 1, we observe that the performances of REAPER and

Table 2. Results on the Caltech-256 database. Inliers are taken to be images of one, three, or five randomly chosen categories, and outliers are randomly chosen from category 257-clutter. For R-graph we set $\alpha = 20$ in the definition of γ .

	OutRank	CoP	REAPER	OutlierPursuit	LRR	DPCP	ℓ_1 -thresholding	R-graph (ours)
<i>Inliers: one category of images Outliers: 50%</i>								
AUC	0.897	0.905	0.816	0.837	0.907	0.783	0.772	0.948
F1	0.866	0.880	0.808	0.823	0.893	0.785	0.772	0.914
<i>Inliers: three categories of images Outliers: 50%</i>								
AUC	0.574	0.676	0.796	0.788	0.479	0.798	0.810	0.929
F1	0.682	0.718	0.784	0.779	0.671	0.777	0.782	0.880
<i>Inliers: five categories of images Outliers: 50%</i>								
AUC	0.407	0.487	0.657	0.629	0.337	0.676	0.774	0.913
F1	0.667	0.672	0.716	0.711	0.667	0.715	0.762	0.858

Table 3. Results on the Coil-100 database. Inliers are taken to be images of one, four, or seven randomly chosen categories, and outliers are randomly chosen from other categories (at most one from each category). For R-graph we set $\alpha = 10$ in the definition of γ .

	OutRank	CoP	REAPER	OutlierPursuit	LRR	DPCP	ℓ_1 -thresholding	R-graph (ours)
<i>Inliers: all images from one category Outliers: 50%</i>								
AUC	0.836	0.843	0.900	0.908	0.847	0.900	0.991	0.997
F1	0.862	0.866	0.892	0.902	0.872	0.882	0.978	0.990
<i>Inliers: all images from four categories Outliers: 25%</i>								
AUC	0.613	0.628	0.877	0.837	0.687	0.859	0.992	0.996
F1	0.491	0.500	0.703	0.686	0.541	0.684	0.941	0.970
<i>Inliers: all images from seven categories Outliers: 15%</i>								
AUC	0.570	0.580	0.824	0.822	0.628	0.804	0.991	0.996
F1	0.342	0.346	0.541	0.528	0.366	0.511	0.897	0.955

DPCP are less competitive in comparison to Outlier Pursuit and our R-graph for the three subspace case.

5.3. Outliers in images of objects

We test the ability of the methods to identify one or several object categories that frequently appear in a set of images amidst outliers that consist of objects that rarely occur. For Caltech-256, images in $n \in \{1, 3, 5\}$ randomly chosen categories are used as inliers in three different experiments. From each category, we use the first 150 images if the category has more than 150 images. We then randomly pick a certain number of images from the “clutter” category as outliers such that there are 50% outliers in each experiment. For Coil-100, we randomly pick $n \in \{1, 4, 7\}$ categories as inliers and pick at most one image from each of the remaining categories as outliers.

The results are reported in Table 2 and Table 3. We see that our R-graph method achieves the best performance. The two geometric distance based methods, OutRank and CoP, achieve good results when there is one inlier category, but deteriorate when the number of inlier categories increases. The performance of REAPER, Outlier Pursuit and DPCP are similar to each other and worse than our method. This may be because they all try to fit a linear subspace to the data, while the data in these two databases may be better modeled by a nonlinear manifold. The ℓ_1 -thresholding

and the representation graph method are all based on data self-expression, and seem to be more powerful for this data.

6. Conclusion

We presented an outlier detection method that combined data self-representation and random walks on a representation graph. Unlike many prior methods for robust PCA, our method is able to deal with multiple subspaces and does not require the number of subspaces or their dimensions to be known. Our analysis showed that the method is guaranteed to identify outliers when certain geometric conditions are satisfied and two connectivity assumptions hold. In our experiments on face image and object image databases, our method achieves the state-of-the-art performance.

Acknowledgment

This work was supported by NSF BIGDATA grant 1447822. The authors also thank Manolis Tsakiris, Conner Lane and Chun-Guang Li for helpful comments.

Appendices

The appendix is organized as follows. In Section A we discuss subspace-preserving representations and give an outline of the proof for Theorem 4.1. Section B contains relevant background on Markov chain theory, which is then used in Section C for proving Lemma 4.1 and Lemma 4.2, as well as providing an in-depth discussion of the Cesàro mean used for outlier detection. In Section D we provide some additional results for experiments on the Extended Yale B database that provide additional insight into the behavior of the methods.

A. Subspace-preserving representation and proof of Theorem 4.1

The idea of a subspace-preserving representation has been extensively studied in the literature of subspace clustering to guarantee the correctness of clustering [12, 43, 44, 31, 27, 10, 38, 59, 17, 58, 56, 50, 52, 24]. Concretely, the data in a subspace clustering task are assumed to lie in a union of low dimensional subspaces, without any outliers that lie outside of the subspaces. A data self-representation matrix is called subspace-preserving if each point uses only points that are from its own subspace in its representation.

Theoretical results in subspace clustering can be adapted to study subspace-preserving representations in the presence of outliers. Here, we use the analysis and result from [56], which studied the elastic net representation (4) for subspace clustering, to prove a subspace-preserving representation result in the presence of outliers, i.e. Theorem 4.1. We also present a corollary of Theorem 4.1 which allows us to compare our result with other subspace clustering results.

A.1. Proof of Theorem 4.1

The proof of Theorem 4.1 follows mostly from the work [56]. We provide an outline of the proof for completeness.

Consider the vector \mathbf{r}_j^ℓ , which is the solution of the problem in the statement of Theorem 4.1. Notice that the entries of \mathbf{r}_j^ℓ correspond to columns of the data matrix X_{-j}^ℓ . One can subsequently construct a representation vector by padding additional zeros to \mathbf{r}_j^ℓ at entries corresponding to points in X that are not in X_{-j}^ℓ . Note that this vector is trivially subspace-preserving by construction. The idea of the proof is to show that this constructed vector, which is subspace-preserving by construction, is a solution to the optimization problem (4) (and no other vector is). A sufficient condition for this to hold is that δ_j , which is computed from \mathbf{r}_j^ℓ , needs to have low correlation with all points $\mathbf{x}_k \notin \mathcal{S}_\ell$. More precisely, we have the following lemma.

Lemma A.1 ([56, Lemma 3.1]). *The vector \mathbf{r}_j is subspace-preserving if $|\langle \mathbf{x}_k, \delta_j \rangle| < \lambda$ for all $\mathbf{x}_k \notin \mathcal{S}_\ell$.*

Lemma A.1 can be proved by using the optimality condition of the optimization problem in (4). Equivalently, it suggests that \mathbf{r}_j is subspace-preserving if

$$\max_{k: \mathbf{x}_k \notin \mathcal{S}_\ell} |\langle \mathbf{x}_k, \bar{\delta}_j \rangle| < \frac{\lambda}{\|\delta_j\|_2}. \quad (\text{A.1})$$

To get more meaningful results, we need an upper bound on $\|\delta_j\|_2$. This is provided by the following lemma.

Lemma A.2 ([57, Lemma C.2]). *If κ_j be the maximum coherence between the oracle point δ_j and columns of X_{-j}^ℓ , i.e. $\kappa_j = \max_{k \neq j: \mathbf{x}_k \in \mathcal{S}_\ell} |\langle \mathbf{x}_k, \bar{\delta}_j \rangle|$, then*

$$\|\delta_j\|_2 \leq \frac{\lambda \kappa_j + 1 - \lambda}{\kappa_j^2}. \quad (\text{A.2})$$

Combining (A.1) and (A.2), \mathbf{r}_j is subspace-preserving if

$$\max_{k: \mathbf{x}_k \notin \mathcal{S}_\ell} |\langle \mathbf{x}_k, \bar{\delta}_j \rangle| < \frac{\kappa_j^2}{\kappa_j + \frac{1-\lambda}{\lambda}}. \quad (\text{A.3})$$

To simplify the result, note that

$$\begin{aligned} \frac{\kappa_j^2}{\kappa_j + \frac{1-\lambda}{\lambda}} &= \kappa_j \cdot \left(\frac{1}{1 + \frac{1-\lambda}{\lambda} \frac{1}{\kappa_j}} \right) \\ &\geq \kappa_j \cdot \left(1 - \frac{1-\lambda}{\lambda} \frac{1}{\kappa_j} \right) = \kappa_j - \frac{1-\lambda}{\lambda}. \end{aligned}$$

Therefore, a sufficient condition for \mathbf{r}_j to be subspace-preserving is that

$$\max_{k: \mathbf{x}_k \notin \mathcal{S}_\ell} |\langle \mathbf{x}_k, \bar{\delta}_j \rangle| < \kappa_j - \frac{1-\lambda}{\lambda}. \quad (\text{A.4})$$

Since (A.4) is the same as (7), the proof has been completed.

A.2. Discussions

Another commonly used geometric quantity for characterizing when representations will be subspace-preserving is the inradius of sets of points [43, 44, 59, 58, 53, 52, 50]. In order to understand the relationship to the results found in these works, we present a corollary of Theorem 4.1.

Definition A.1 (inradius). *The (relative) inradius of a convex body \mathcal{P} , denoted as $\rho(\mathcal{P})$, is the radius of the largest ℓ_2 ball in the span of \mathcal{P} that can be inscribed in \mathcal{P} .*

Corollary A.1. *If $\mathbf{x}_j \in \mathcal{S}_\ell$ is an inlier, then \mathbf{r}_j computed from (4) is subspace-preserving if*

$$\rho_j - \max_{k: \mathbf{x}_k \notin \mathcal{S}_\ell} |\langle \mathbf{x}_k, \bar{\delta}_j \rangle| > \frac{1-\lambda}{\lambda}, \quad (\text{A.5})$$

where $\bar{\delta}_j$ is defined in Theorem 4.1, and ρ_j is the inradius of the convex hull of the symmetrized points in X_j^ℓ , i.e.

$$\rho_j := \rho(\text{conv}\{\pm \mathbf{x}_k : \mathbf{x}_k \in \mathcal{S}_\ell, k \neq j\}). \quad (\text{A.6})$$

The inradius captures the distribution of the columns of X_{-j}^ℓ , i.e. it is large if points are well spread out in S_ℓ . Thus, the condition in (A.5) is easier to be satisfied if the set of points in S_ℓ is dense and well covers the entire subspace. Note that this requirement is stronger than that in Theorem 4.1, which only requires points in S_ℓ to be dense around the oracle point δ_j (i.e. it requires $\max_{k \neq j: \mathbf{x}_k \in S_\ell} |\langle \mathbf{x}_k, \bar{\delta}_j \rangle|$ to be large). In fact, it is established in [56] that $\max_{k \neq j: \mathbf{x}_k \in S_\ell} |\langle \mathbf{x}_k, \bar{\delta}_j \rangle| \geq \rho_j$, so that the condition in (A.5) is a stronger requirement than that of (7) in Theorem 4.1.

B. Background on Markov chain theory

We present background material on Markov chain theory that will help us understand the Cesàro mean (6) used for outlier detection in our method. The following material is organized from textbooks [41, 14, 45, 23] and the website <http://www.math.uah.edu/stat>.

We consider a Markov chain (X_0, X_1, \dots) on a finite state space Ω with transition probabilities p_{ij} for $i, j \in \Omega$. The t -step transition probabilities are defined to be $p_{ij}^{(t)} := P\{X_t = j | X_0 = i\}$.

B.1. Decomposition of the state space

A Markov chain can be decomposed into more basic and manageable parts.

Definition B.1. State j is accessible from state i , denoted as $i \rightarrow j$, if $p_{ij}^{(t)} > 0$ for some $t > 0$. We say that the states i and j communicate with each other, denoted by $i \leftrightarrow j$, if $i \rightarrow j$ and $j \rightarrow i$.

Since it can be shown that \leftrightarrow is an equivalence relation, it induces a partition of the state space Ω into disjoint equivalence classes known as *communicating classes*. We are interested in each of the *closed* communicating classes.

Definition B.2. A non-empty set $C \subseteq \Omega$ is called a *closed* set if $p_{ij} = 0$ for $i \in C$ and $j \notin C$.

Note that states in a closed communicating class are essential while states in other communicating classes are inessential [23].

Theorem B.1 ([41]). The state space Ω has the unique decomposition $\Omega = \mathcal{I} \cup \mathcal{E}_1 \cup \dots \cup \mathcal{E}_n$, where \mathcal{I} is the set of inessential states, and $\mathcal{E}_1, \dots, \mathcal{E}_n$ are closed communicating classes containing essential states.

By Theorem B.1, the state space of any Markov chain is composed of the essential states and inessential states, and the essential states can be further decomposed into a union of communicating classes. Therefore, the probability

transition matrix P can be written in the following form (up to permutation of the states):

$$P = \begin{bmatrix} P_{\mathcal{E}_1 \rightarrow \mathcal{E}_1} & \mathbf{0} & \mathbf{0} \\ & \ddots & \vdots \\ \mathbf{0} & P_{\mathcal{E}_n \rightarrow \mathcal{E}_n} & \mathbf{0} \\ P_{\mathcal{I} \rightarrow \mathcal{E}_1} & \dots & P_{\mathcal{I} \rightarrow \mathcal{E}_n} & P_{\mathcal{I} \rightarrow \mathcal{I}} \end{bmatrix} \quad (\text{B.1})$$

B.2. Stationary distribution

A nonnegative row vector π is called a *stationary distribution* for the Markov chain if it satisfies $\pi = \pi P$.

Theorem B.2 ([23, Proposition 1.14, Corollary 1.17]). A Markov chain consisting of one closed communicating class has a unique stationary distribution. Moreover, each entry of the stationary distribution is positive.

By Theorem B.2, each component \mathcal{E}_ℓ for $\ell = 1, \dots, n$ in the decomposition of the Markov chain in Theorem B.1 has a unique positive stationary distribution $\pi_{\mathcal{E}_\ell}$, i.e.

$$\pi_{\mathcal{E}_\ell} = \pi_{\mathcal{E}_\ell} \cdot P_{\mathcal{E}_\ell \rightarrow \mathcal{E}_\ell} \quad \text{with } \pi_{\mathcal{E}_\ell} > 0 \text{ and } \sum_j (\pi_{\mathcal{E}_\ell})_j = 1. \quad (\text{B.2})$$

We may then define a stationary distribution for P as

$$[\alpha_1 \pi_{\mathcal{E}_1}, \dots, \alpha_n \pi_{\mathcal{E}_n}, \mathbf{0}] \quad \text{for any } \alpha_\ell \geq 0, \sum_{\ell=1}^n \alpha_\ell = 1. \quad (\text{B.3})$$

Note that there is not a unique stationary distribution for P when $n \geq 2$.

B.3. Convergence of the Cesàro mean $\frac{1}{T} \sum_{t=1}^T P^t$

Let $f_{ij}^{(t)} := P\{X_t = j, X_{t'} \neq j \text{ for } 1 \leq t' < t | X_0 = i\}$ be the probability that the chain starting at i enters j for the first time at the t -th step. The *hitting probability* $f_{ij} = P\{X_t = j \text{ for some } t > 0 | X_0 = i\}$ is the probability that the random walk ever makes a transition to state j when started at i , i.e.

$$f_{ij} = \sum_{t=1}^{\infty} f_{ij}^{(t)}. \quad (\text{B.4})$$

The *mean return time* $\mu_j := \sum_{t=1}^{\infty} t f_{jj}^{(t)}$ is the expected time for a random walk starting from state j will return to state j . A general convergence result is stated as follows.

Theorem B.3 ([45, Theorem 3.3.1]). For any $i, j \in \Omega$,

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T p_{ij}^{(t)} = \frac{f_{ij}}{\mu_j}. \quad (\text{B.5})$$

This result can be simplified by using the decomposition in Theorem B.1, which leads to the following lemma.

Lemma B.1. *If $i, j \in \Omega$ are in the same closed communicating class, then $f_{ij} = f_{ji} = 1$. Also, if $i \in \Omega$ is an inessential state and $\mathcal{E}_\ell \subseteq \Omega$ is a closed communicating class, then $f_{ij} = f_{i \rightarrow \mathcal{E}_\ell}$ for all $j \in \mathcal{E}_\ell$, where $f_{i \rightarrow \mathcal{E}_\ell}$ is the hitting probability from state i to class \mathcal{E}_ℓ .*

The following result relates the mean return time with the stationary distribution.

Lemma B.2. *For every closed communicating class $\mathcal{E}_\ell \subseteq \Omega$, it holds that $\mu_{\mathcal{E}_\ell} = 1/\pi_{\mathcal{E}_\ell}$ (entry-wise division), where $\mu_{\mathcal{E}_\ell}$ is the vector of mean return times of states in \mathcal{E}_ℓ . If $i \in \Omega$ is an inessential state, then $\mu_i = \infty$.*

By combining Theorem B.3 with Lemma B.1 and Lemma B.2, the Cesàro limit of a probability transition matrix of the form in (B.1) can be written as

$$\lim_{\frac{1}{T}} \sum_{t=1}^T P^t = \begin{bmatrix} \mathbf{1} \cdot \pi_{\mathcal{E}_1} & & \mathbf{0} & \mathbf{0} \\ & \ddots & & \vdots \\ \mathbf{0} & & \mathbf{1} \cdot \pi_{\mathcal{E}_n} & \mathbf{0} \\ \mathbf{f}_{\mathcal{I} \rightarrow \mathcal{E}_1} \cdot \pi_{\mathcal{E}_1} & \cdots & \mathbf{f}_{\mathcal{I} \rightarrow \mathcal{E}_n} \cdot \pi_{\mathcal{E}_n} & \mathbf{0} \end{bmatrix}, \quad (\text{B.6})$$

in which $\mathbf{f}_{\mathcal{I} \rightarrow \mathcal{E}_\ell}$ is a column vector of hitting probability from each state in \mathcal{I} to class \mathcal{E}_ℓ .

We note that while the Cesàro mean converges, the t -step transition probability P^t does not necessarily converge. Consider, for example, the probability transition matrix $P = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$. In this case, $p_{12}^{(t)} = 1$ when t is odd and $p_{12}^{(t)} = 0$ when t is even, i.e. $p_{12}^{(t)}$ is oscillating and never converges. In general, P^t converges if and only if each of the closed communicating classes \mathcal{E}_ℓ for $\ell = 1, \dots, n$ is *aperiodic*.

C. Guaranteed outlier detection

Our outlier detection method by representation graph is guaranteed to correctly identify outliers in a union of subspaces when the representation is subspace-preserving and that the connectivity assumptions are satisfied. In this section, we first prove that the inliers and outliers in the data correspond to essential and inessential states, respectively, of the Markov chain associated with the representation graph (Lemma 4.1). Then, we show that the average of the first T t -step probability distributions $\frac{1}{T} \sum_{t=1}^T \pi_0 P^t$ identifies essential and inessential states (Lemma 4.2), thus establishing the correctness of our method.

C.1. Proof of Lemma 4.1

Recall that we work with a Markov chain with state space $\Omega = \{1, \dots, N\}$, in which each state i corresponds to the point \mathbf{x}_i in the data matrix X .

First, we show that any inlier point \mathbf{x}_i corresponds to an essential state of the Markov chain. Let \mathbf{x}_j be any point such that $i \rightarrow j$. Since the representation matrix is subspace-preserving, we know that \mathbf{x}_i and \mathbf{x}_j lie in the same subspace. Furthermore, by Assumption 4.1, all points in the same subspace are strongly connected, which implies that $j \rightarrow i$. Thus, i is an essential state.

Second, we show that any outlier point \mathbf{x}_i corresponds to an inessential state of the Markov chain. Consider the set $\Omega_i = \{k : i \rightarrow k\}$, i.e. the set of points that are accessible from \mathbf{x}_i . By Assumption 4.2, the set Ω_i cannot contain only outliers. Thus, there exists \mathbf{x}_j such that $i \rightarrow j$ and \mathbf{x}_j is an inlier. However, since the representation is subspace-preserving, we know that $j \not\rightarrow i$. Therefore, i is not an essential state, i.e. it is an inessential state.

C.2. Proof of Lemma 4.2

According to Theorem B.1, the state space of the Markov chain can be decomposed into $\mathcal{I} \cup \mathcal{E}_1 \cup \dots \cup \mathcal{E}_n$, in which \mathcal{I} contains the inessential states and each \mathcal{E}_ℓ is a closed communicating class containing essential states. Assume, without loss of generality, that the transition probability matrix has the form of (B.1). By using (B.6), the Cesàro mean in (6) has the following limiting behavior:

$$\begin{aligned} \pi &:= \lim_{T \rightarrow \infty} \bar{\pi}^{(T)} = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \pi_0 P^t \\ &= \left[\frac{N_1 + \sum \mathbf{f}_{\mathcal{I} \rightarrow \mathcal{E}_1} \cdot \pi_{\mathcal{E}_1}}{N}, \dots, \frac{N_n + \sum \mathbf{f}_{\mathcal{I} \rightarrow \mathcal{E}_n} \cdot \pi_{\mathcal{E}_n}}{N}, \mathbf{0} \right], \end{aligned} \quad (\text{C.1})$$

where N_ℓ for $\ell = 1, \dots, n$ is the number of states in class \mathcal{E}_ℓ , each $\mathbf{f}_{\mathcal{I} \rightarrow \mathcal{E}_\ell}$ is a vector of hitting probabilities for each state in \mathcal{I} to class \mathcal{E}_ℓ , and $\mu_{\mathcal{E}_\ell}$ is a positive vector of the stationary distributions of states in \mathcal{E}_ℓ . Therefore, π_j is zero if and only if j is an inessential state. This finishes the proof.

C.3. Discussion

In this section, we provide additional comments on using the Cesàro mean $\bar{\pi}^{(T)}$ in (6) for outlier detection.

Stationary distributions. By (C.1), the vector that $\bar{\pi}^{(T)}$ converges to is a stationary distribution of the Markov chain (see (B.3)). In fact, any convex combination of the stationary distribution of each closed communicating class is a stationary distribution of the Markov chain, and the particular stationary distribution that $\bar{\pi}^{(T)}$ converges to depends on the choice of the initial state distribution π_0 .

A T -step probability distribution and PageRank. Traditionally, PageRank and many other spectral ranking algorithms use the limit of the T -step probability distribution $\pi^{(T)}$ rather than $\bar{\pi}^{(T)}$ as adopted in our method. However,

Table C.1. Running time of experiments on Extended Yale B data with three inlier groups and 15% outliers

	OutRank	CoP	REAPER	OutlierPursuit	LRR	DPCP	ℓ_1 -thresholding	R-graph (ours)
Time (sec.)	0.019	0.003	0.079	1.186	3.502	0.182	0.312	0.272

the sequence $\pi^{(T)}$ converges if and only if each closed communicating class of the Markov chain is aperiodic, which is not necessarily satisfied in many cases. To address this, PageRank adopts a random walk with restart algorithm. It can be interpreted as a random walk on a transformed Markov chain that adds a small probability of transition from each state to the other states on the transition probability of the original Markov chain. By doing so, the transformed Markov chain contains a single communicating class that is aperiodic. Therefore, the stationary distribution necessarily becomes unique, and the sequence $\pi^{(T)}$ for the transformed Markov chain converges to the unique stationary distribution regardless of the initial state distribution.

Despite the advantages of the random walk used by PageRank, all states of the Markov Chain are essential, so that outliers do not converge to zero probabilities. Therefore, it is less clear whether the stationary distribution that the algorithm converges to can effectively identify outliers.

D. Additional experimental results

D.1. Computational time comparison

Table C.2 reports the average running time of the experiment on the Extended Yale B database with three inlier groups and 15% outliers (226 images in total). From the table we observe that the running times of OutRank and CoP are much smaller than the other methods. This comes from the fact that OutRank and CoP are based on computing data pairwise inner products, which is efficient for small scale data. In contrast, the other methods solve optimization problems. In particular, REAPER, OutlierPursuit and LRR require computing an eigendecomposition of a matrix of size $D \times D$ (D is the ambient dimension) during each iteration, which is time consuming when D is large. In our experiments we observe that REAPER converges much faster than OutlierPursuit and LRR, thus the running time of REAPER is typically much smaller. The ℓ_1 -thresholding method and R-graph method (our algorithm) both compute the representation matrix by solving an ℓ_1 optimization problem for each of the data points with all other data points as the dictionary. Subsequently, ℓ_1 -thresholding rejects outliers simply by computing the ℓ_1 norms of the representations, while R-graph requires a random walk on the graph defined from the representation. We note that the random walk for R-graph is computationally efficient because of the sparsity of the representation matrix. In each step of the random walk, the computational complexity is on the order of sN where N is the number of data points and $s \ll N$ is the average number of nonzero entries in the

representation vectors $\{r_j\}$.

D.2. Influence of the algorithm parameters

The first step of our method is to compute the data self-representation matrix using the optimization problem (4). In this section, we illustrate the effect that the parameter γ in (4) has on the performance of our method. Recall that for our numerical experiments we set $\gamma = \alpha \cdot \frac{\lambda}{\max_{i:i \neq j} |x_j^\top x_i|}$ and that the solution to (4) is nonzero if and only if $\alpha > 1$. We run experiments on Extended Yale B database with 3 inlier groups and 15% outliers while varying α in the range $[1, 50]$; the results are shown in Figure 1(a). We can see that the R-graph performs well over a wide range of the parameter α . For comparison, Figure 1(a) also plots the performance of the other methods on the same dataset.

D.3. Influence of the percentage of outliers

In this experiment, we fix the number of inlier groups to be 3 and vary the percentage of outliers from 1% to 15%. The performances of the different methods are reported in Figure 1(b). Note that the parameters for all methods are fixed across the different percentages of outliers. We see that the performance of our method is stable with respect to the percentage of outliers. Moreover, our method also achieves the best performance among all methods.

D.4. Visualization of the outliers

To supplement the AUC and F1 measures previously provided, and also to better understand the outliers returned by our outlier detection method, we conducted additional experiments that display the top outliers detected in each experiment. The set of inliers is taken to be the 64 images of the first subject of the Extended Yale B database, and the outlier set is chosen as 10 images randomly chosen from the remaining 37 subjects (see Figure D.2). The top 10 outliers returned by different methods are reported in Figure D.3. Images with red boxes are outliers (i.e. true positives) and images with green boxes are inliers (i.e. false positives).

False positives for all methods are mostly images taken under extreme illumination conditions. Such images have large shadows, which has the effect of removing them from the underlying subspace associated with the individual thus making them more likely to be detected as outliers. The results show that REAPER, Outlier Pursuit, DPCP and R-graph are relatively robust. In particular, R-graph is significantly better than ℓ_1 -thresholding even though both are sparse representation based methods. This shows that while the magnitude of the representation vector adopted by ℓ_1 -

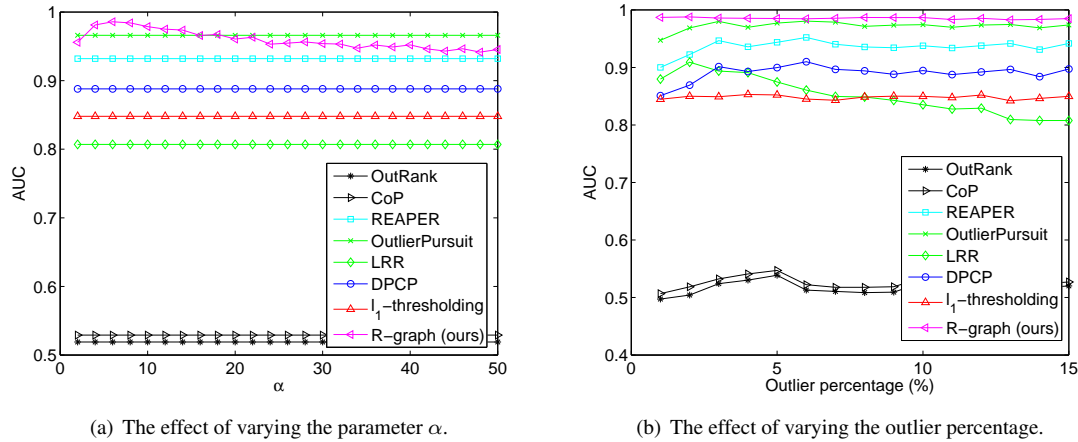


Figure D.1. Additional results for experiments on Extended Yale B with three inlier groups and 15% outliers.

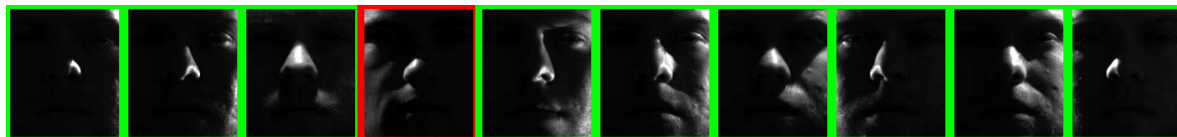


Figure D.2. An outlier detection dataset for visualizing the top 10 outliers returned by different methods.

thresholding can be sensitive to corruptions, the connectivity behavior explored by R-graph is more robust.

References

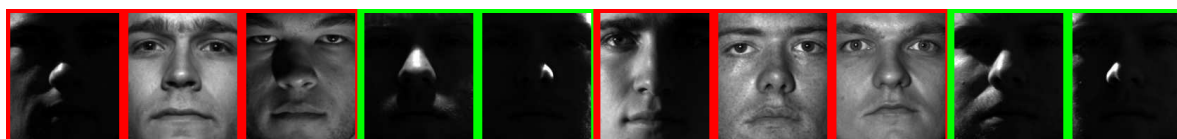
- [1] E. Arias-Castro, G. Chen, and G. Lerman. Spectral clustering based on local linear approximations. *Electron. J. Statist.*, 5:1537–1587, 2011. 1
- [2] R. Basri and D. Jacobs. Lambertian reflection and linear subspaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(2):218–233, 2003. 1
- [3] S. Brin and L. Page. The anatomy of a large-scale hyper-textual web search engine. *Computer Networks and ISDN Systems*, 30:107–117, 1998. 3
- [4] E. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis. *Journal of the ACM*, 58, 2011. 1
- [5] G. Chen and G. Lerman. Spectral curvature clustering (SCC). *International Journal of Computer Vision*, 81(3):317–330, 2009. 1
- [6] Y. Cherapanamjeri, P. Jain, and P. Netrapalli. Thresholding based efficient outlier robust pca. *arXiv preprint arXiv:1702.05571*, 2017. 1
- [7] T.-J. Chin, Y. Heng Kee, A. Eriksson, and F. Neumann. Guaranteed outlier removal with mixed integer linear programs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5858–5866, 2016. 3
- [8] Y. Cong, J. Yuan, and J. Liu. Sparse reconstruction cost for abnormal event detection. In *The 24th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2011, Colorado Springs, CO, USA, 20-25 June 2011*, pages 3449–3456, 2011. 2
- [9] C. Ding, D. Zhou, X. He, and H. Zha. R_1 -pca: rotational invariant l_1 -norm principal component analysis for robust subspace factorization. In *Proceedings of the 23rd international conference on Machine learning*, pages 281–288. ACM, 2006. 1
- [10] E. L. Dyer, A. C. Sankaranarayanan, and R. G. Baraniuk. Greedy feature selection for subspace clustering. *Journal of Machine Learning Research*, 14(1):2487–2517, 2013. 9
- [11] E. Elhamifar and R. Vidal. Sparse subspace clustering. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2790–2797, 2009. 2
- [12] E. Elhamifar and R. Vidal. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(11):2765–2781, 2013. 3, 9
- [13] M. A. Fischler and R. C. Bolles. RANSAC random sample consensus: A paradigm for model fitting with applications to



(a) Top 10 outliers by OutRank



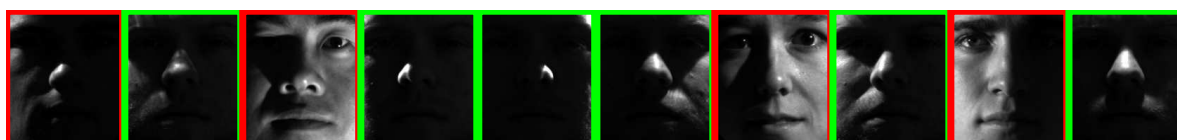
(b) Top 10 outliers by CoP



(c) Top 10 outliers by REAPER



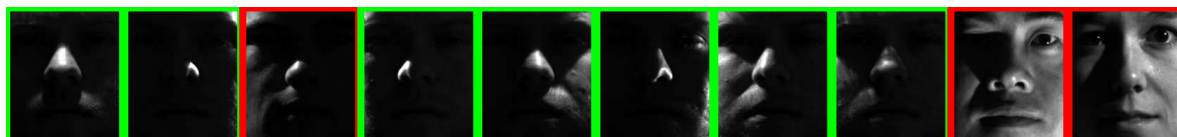
(d) Top 10 outliers by OutlierPursuit



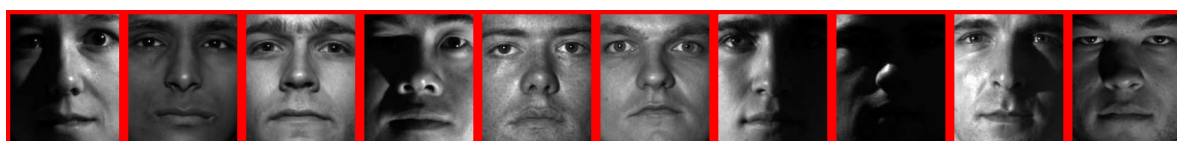
(e) Top 10 outliers by LRR



(f) Top 10 outliers by DPCP



(g) Top 10 outliers by ℓ_1 -thresholding



(h) Top 10 outliers by R-graph (ours)

Figure D.3. Visualizing the top 10 outliers from different methods. Image in red box: true outlier. Image in green box: true inlier.

- image analysis and automated cartography. *Communications of the ACM*, 26:381–395, 1981. 1
- [14] R. G. Gallager. *Stochastic processes: theory for applications*. Cambridge University Press, 2013. 10
- [15] A. Georgiades, P. Belhumeur, and D. Kriegman. From few to many: Illumination cone models for face recognition under variable lighting and pose. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(6):643–660, 2001. 6
- [16] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. 2007. 6
- [17] R. Heckel and H. Bölcskei. Robust subspace clustering via thresholding. *IEEE Transactions on Information Theory*, 61(11):6320–6342, 2015. 1, 9
- [18] B. Jin, D. Lorenz, and S. Schiffler. Elastic-net regularization: error estimates and active set methods. *Inverse Problems*, 25(11), 2009. 6
- [19] K. Kanatani. Motion segmentation by subspace separation and model selection. In *IEEE International Conference on Computer Vision*, volume 2, pages 586–591, 2001. 1
- [20] G. Lerman and T. Maunu. Fast, robust and non-convex subspace recovery. *arXiv preprint arXiv:1406.6145*, 2014. 1
- [21] G. Lerman, M. B. McCoy, J. A. Tropp, and T. Zhang. Robust computation of linear models by convex relaxation. *Foundations of Computational Mathematics*, 15(2):363–410, 2015. 1, 6
- [22] G. Lerman and T. Zhang. Robust recovery of multiple subspaces by geometric ℓ_p minimization. *Annals of Statistics*, 39(5):2686–2715, 2011. 1
- [23] D. A. Levin, Y. Peres, and E. L. Wilmer. *Markov chains and mixing times*. American Mathematical Soc., 2009. 5, 10
- [24] J. Li, Y. Kong, and Y. Fu. Sparse subspace clustering by learning approximation ℓ_0 codes. *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, 2189–2195, 2017. 9
- [25] X. Li and J. Haupt. Identifying outliers in large matrices via randomized adaptive compressive sampling. *IEEE Transactions on Signal Processing*, 63(7):1792–1807, 2015. 1
- [26] Z. Lin, M. Chen, L. Wu, and Y. Ma. The augmented Lagrange multiplier method for exact recovery of corrupted low-rank matrices. *arXiv:1009.5055v2*, 2011. 6
- [27] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma. Robust recovery of subspace structures by low-rank representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2012. 9
- [28] G. Liu, Z. Lin, and Y. Yu. Robust subspace segmentation by low-rank representation. In *International Conference on Machine Learning*, pages 663–670, 2010. 2, 6
- [29] G. Liu, H. Xu, and S. Yan. Exact subspace segmentation and outlier detection by low-rank representation. In *AISTATS*, pages 703–711, 2012. 3
- [30] C. Lu, Z. Lin, and S. Yan. Correlation adaptive subspace segmentation by trace lasso. In *IEEE International Conference on Computer Vision*, pages 1345–1352, 2013. 5
- [31] C.-Y. Lu, H. Min, Z.-Q. Zhao, L. Zhu, D.-S. Huang, and S. Yan. Robust and efficient subspace segmentation via least squares regression. In *European Conference on Computer Vision*, pages 347–360, 2012. 9
- [32] M. McCoy, J. A. Tropp, et al. Two proposals for robust pca using semidefinite programming. *Electronic Journal of Statistics*, 5:1123–1160, 2011. 1
- [33] K. Mitra, A. Veeraraghavan, and R. Chellappa. Analysis of sparse regularization based robust regression approaches. *IEEE Transactions on Signal Processing*, 61(5):1249–1257, 2013. 3
- [34] H. Moonesinghe and P.-N. Tan. Outlier detection using random walks. In *2006 18th IEEE International Conference on Tools with Artificial Intelligence (ICTAI'06)*, pages 532–539. IEEE, 2006. 3, 6
- [35] H. Moonesinghe and P.-N. Tan. Outrank: a graph-based outlier detection framework using random walk. *International Journal on Artificial Intelligence Tools*, 17(01):19–36, 2008. 3, 6
- [36] B. Nasihatkon and R. Hartley. Graph connectivity in sparse subspace clustering. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2011. 5
- [37] S. A. Nene, S. K. Nayar, and H. Murase. Columbia object image library (COIL-100). *Technical Report CUCS-006-96*, 1996. 6
- [38] D. Park, C. Caramanis, and S. Sanghavi. Greedy subspace clustering. In *Neural Information Processing Systems*, 2014. 9
- [39] Q. Qu, J. Sun, and J. Wright. Finding a sparse vector in a subspace: Linear sparsity using alternating directions. In *Advances in Neural Information Processing Systems*, pages 3401–3409, 2014. 3
- [40] M. Rahmani and G. Atia. Coherence pursuit: Fast, simple, and robust principal component analysis. *arXiv preprint arXiv:1609.04789*, 2016. 1, 6
- [41] R. Serfozo. *Basics of applied stochastic processes*. Springer Science & Business Media, 2009. 10
- [42] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 6
- [43] M. Soltanolkotabi and E. J. Candès. A geometric analysis of subspace clustering with outliers. *Annals of Statistics*, 40(4):2195–2238, 2012. 2, 3, 6, 9
- [44] M. Soltanolkotabi, E. Elhamifar, and E. J. Candès. Robust subspace clustering. *Annals of Statistics*, 42(2):669–699, 2014. 9
- [45] H. C. Tijms. *A first course in stochastic models*. John Wiley and Sons, 2003. 10
- [46] M. Tsakiris and R. Vidal. Dual principal component pursuit. In *ICCV Workshop on Robust Subspace Learning and Computer Vision*, pages 10–18, 2015. 3, 6
- [47] R. Vidal, Y. Ma, and S. Sastry. *Generalized Principal Component Analysis*. Springer Verlag, 2016. 5
- [48] U. von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, 2007. 2
- [49] Y. Wang, C. Dicle, M. Sznajder, and O. Camps. Self scaled regularized robust regression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3261–3269, 2015. 3
- [50] Y. Wang, Y. Wang, and A. Singh. A deterministic analysis of noisy sparse subspace clustering for dimensionality-reduced

- data. In *International Conference on Machine Learning*, pages 1422–1431, 2015. 3, 9
- [51] Y. Wang, Y.-X. Wang, and A. Singh. Graph connectivity in noisy sparse subspace clustering. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, pages 538–546, 2016. 5
- [52] Y.-X. Wang and H. Xu. Noisy sparse subspace clustering. *Journal of Machine Learning Research*, 17(12):1–41, 2016. 3, 9
- [53] Y.-X. Wang, H. Xu, and C. Leng. Provable subspace clustering: When LRR meets SSC. In *Neural Information Processing Systems*, 2013. 5, 9
- [54] J. Wright, A. Ganesh, S. Rao, Y. Peng, and Y. Ma. Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization. In *NIPS*, 2009. 1
- [55] H. Xu, C. Caramanis, and S. Sanghavi. Robust pca via outlier pursuit. In *Advances in Neural Information Processing Systems*, pages 2496–2504, 2010. 1, 6
- [56] C. You, C.-G. Li, D. Robinson, and R. Vidal. Oracle based active set algorithm for scalable elastic net subspace clustering. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3928–3937, 2016. 3, 5, 9, 10
- [57] C. You, C.-G. Li, D. Robinson, and R. Vidal. Oracle based active set algorithm for scalable elastic net subspace clustering. *Arxiv*, 2016. 9
- [58] C. You, D. Robinson, and R. Vidal. Scalable sparse subspace clustering by orthogonal matching pursuit. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3918–3927, 2016. 9
- [59] C. You and R. Vidal. Geometric conditions for subspace-sparse recovery. In *International Conference on Machine learning*, pages 1585–1593, 2015. 3, 9
- [60] T. Zhang and G. Lerman. A novel m-estimator for robust pca. *The Journal of Machine Learning Research*, 15(1):749–808, 2014. 1
- [61] T. Zhang, A. Szlam, and G. Lerman. Median k -flats for hybrid linear modeling with many outliers. In *Workshop on Subspace Methods*, pages 234–241, 2009. 1
- [62] T. Zhang, A. Szlam, Y. Wang, and G. Lerman. Hybrid linear modeling via local best-fit flats. *International Journal of Computer Vision*, 100(3):217–240, 2012. 1
- [63] Y. Zheng, S. Sugimoto, and M. Okutomi. Deterministically maximizing feasible subsystem for robust model fitting with unit norm constraint. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1825–1832. IEEE, 2011. 3
- [64] H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society, Series B*, 67:301–320, 2005. 3