

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

### Introduction

• Most collected visual data today are unlabeled/weakly labeled • High-dimensional data often lie approx. in low-dimensional subspaces | • Observation: each data  $\mathbf{x}_i$  can be expressed as a linear combination • Subspace clustering is the problem of clustering data into subspaces • This work addresses sensitivity of subspace clustering to **outliers** 







## Prior Work & Challenges

- Robust PCA methods (e.g. REAPER, Outlier Pursuit) require inliers drawn from a single subspace
- Low-rank methods (e.g., LRR) require inlier subspaces to be independent
- Other methods (TSC, CoP,  $\ell_1$ -thresholding) require dense inliers and/or incoherent outliers

Challenges: multiple subspaces, unknown number of subspaces and their  $\bullet$  Step 2. Define a random walk from  $\{c_i\}$ : dimensions, sparsity of inliers, close-by outliers, etc.

## Contributions

• Outlier detection by using *self-representation* and *random walk* 

Self-representation allows our method to handle multiple subspaces, Step 3. Compute stationary distribution: and the number of subspaces and their dimensions are not required Random walk allows our method to explore contextual information, hence our method can handle sparsity of inliers and close-by outliers Our method is provably correct in identifying outliers

This work was supported by the grant NSF-IIS 1447822.

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

•••/ outliers 

• Task: cluster data  $X = [\mathbf{x}_1, \cdots, \mathbf{x}_N]$  in a union of subspaces of a few others from its own subspace, i.e.  $\mathbf{x}_{i} = X\mathbf{c}_{i}, \mathbf{c}_{ij} = 0$ • Algorithm: find such representation by solving sparse optimization  $\min_{\mathbf{c}_j} \|\mathbf{c}_j\|_1$  s.t.  $\mathbf{x}_j = X\mathbf{c}_j, c_{jj} = 0$ , define affinity between any two points  $\mathbf{x}_i$ ,  $\mathbf{x}_j$  to be  $|[\mathbf{c}_j]_i|$ , then apply spectral clustering • This work: we extend the method to deal with outliers in data X

Input: unlabeled dataset  $X = [\mathbf{x}_1, \cdots, \mathbf{x}_N]$  that contains both inliers and outliers.

• Step 1. Compute data self-representation:

 $\min_{\mathbf{i}} \lambda \|\mathbf{c}_j\|_1 + \frac{1-\lambda}{2} \|\mathbf{c}_j\|_2^2 \quad \text{s.t. } \mathbf{x}_j = X\mathbf{c}_j, c_{jj} = 0$ 

-If  $\mathbf{x}_j$  is inlier:  $[\mathbf{c}_j]_i \neq 0 \rightarrow \mathbf{x}_i$  is inlier -If  $\mathbf{x}_i$  is outlier:  $[\mathbf{c}_i]_i \neq 0 \rightarrow \mathbf{x}_i$  can be either inlier or outlier

 $[P]_{ij} = |[\mathbf{c}_i]_j| / ||\mathbf{c}_i||_1$ 

-There is no transition from inlier to outlier -Any random walker will end up in inliers

$$\bar{\pi}^{(T)} = \frac{1}{T} \sum_{t=1}^{I} \pi^{(0)} P^t$$

where  $\pi^{(0)} = [1/N, \cdots, 1/N]$  is uniform



### IEEE 2017 Conference on Computer Vision and Pattern Recognition

