

# Sparse Subspace Clustering: Algorithm, Theory, and Applications

Elhamifar, E. and Vidal, R.

Presentation for Stat 572

Jerry Wei

Mentor: Vincent Roulet

Department of Statistics  
University of Washington

# Subspace Clustering Problem

High-dimensional data pose challenges to classical clustering methods

**Key observation:** Data in a class or category lie in a low-dimensional subspace

# Subspace Clustering Problem

High-dimensional data pose challenges to classical clustering methods

**Key observation:** Data in a class or category lie in a low-dimensional subspace

**Examples:** rigidly moving object in a video



The video has frames  $f : 1, \dots, F$ , and a set of  $N$  feature points  $\{\mathbf{x}_{fi} \in \mathbb{R}^2\}_{i=1}^N$  is tracked across the frames. For analysis, each feature trajectory  $\mathbf{y}_i$  is taken as a data point, where  $\mathbf{y}_i$  is obtained by stacking the feature points  $x_{fi}$  in the video as

$$\mathbf{y}_i \equiv [\mathbf{x}_{1i}^T, \mathbf{x}_{2i}^T, \dots, \mathbf{x}_{Fi}^T]^T \in \mathbb{R}^{2F}$$

# Subspace Clustering Problem

## Goal:

- find the number of subspaces and their dimensions
- find a basis for each subspace
- group the sample points into subspaces

**Remark:** Extension of traditional clustering: clustering and dimension reduction simultaneously

- Rigid motion: an affine subspace of dimension at most 3
- Images of a subject with fixed pose and varying illumination: lie close to a linear subspace of dimension 9.

# Prior Work on Subspace Clustering

- **Iterative methods** Alternate between assigning points to subspaces and fitting a subspace to each cluster. (K-subspaces, Median K-flats)
- **Algebraic approaches** Based on factorization of the data matrix. (GPCA)
- **Iterative statistical methods** Parametric assumption on data distribution, iterate between clustering and subspace estimation by EM algorithm. (Random Sample Consensus)
- **Information-theoretic statistical approaches** Minimizes the information-theoretic cost to fit a mixture of degenerate Gaussian to the points, up to a given distortion. (Agglomerative Lossy Compression)
- **Spectral clustering-based methods** Construct a similarity graph based on data information and then apply spectral clustering. (Local Subspace Affinity, Low-Rank Subspace Clustering) **Sparse Subspace Clustering (SSC)**

# Notations

Let  $\{\mathcal{S}_l\}_{l=1}^n$  be an arrangement of  $n$  linear subspaces of dimensions  $\{d_l\}_{l=1}^n$  lying in  $\mathbb{R}^D$ .

Let  $\{\mathbf{y}_i\}_{i=1}^N$  be the collection of  $N$  data points that are free of noises and lie in the union of the  $n$  subspaces.

Let  $N_l$  be the number of samples points in subspace  $\mathcal{S}_l$  and the data matrix

$$\mathbf{Y} \equiv [\mathbf{y}_1, \dots, \mathbf{y}_N] = [\mathbf{Y}_1, \dots, \mathbf{Y}_n] \mathbf{\Gamma}$$

where  $\mathbf{Y}_l \in \mathbb{R}^{D \times N_l}$  is a matrix of all the points in  $\mathcal{S}_l$  with  $N_l > d_l$  and  $\mathbf{\Gamma}$  is an unknown permutation matrix.

# Self-expressiveness property

## Definition

*Each data point in a union of subspaces can be efficiently reconstructed by a combination of other points in the dataset.*

**Subspace-Sparse Representation:** represented by a few other points in the same subspace

# Self-expressiveness property

## Definition

*Each data point in a union of subspaces can be efficiently reconstructed by a combination of other points in the dataset.*

**Subspace-Sparse Representation:** represented by **a few** other points **in the same subspace**

Specifically, each data point  $\mathbf{y}_i \in \cup_{l=1}^n \mathcal{S}_l$  can have the representation

$$\mathbf{y}_i = \mathbf{Y}\mathbf{c}_i, c_{ii} = 0$$

where  $\mathbf{c}_i \equiv (c_{i1}, c_{i2}, \dots, c_{iN})^T$  is the vector of constant coefficients

**Goal:**  $\mathbf{c}_i$  has **a few** non-zero entries corresponding to data points **in the same subspace** as  $\mathbf{y}_i$



# Sparse Optimization Problem

**The optimization problem:**

$$\min \|\mathbf{c}_i\|_q \quad s.t. \quad \mathbf{y}_i = \mathbf{Y}\mathbf{c}_i, c_{ii} = 0$$

where  $\mathbf{c}_i$  is a  $N$  dimensional vector of constants on the weights for the other data points and  $\|\cdot\|_q$  is the  $l_q$  norm.

Ideally, want to use  $l_0$  norm, but is NP-hard

# Sparse Optimization Problem

**The optimization problem:**

$$\min \|\mathbf{c}_i\|_q \quad \text{s.t.} \quad \mathbf{y}_i = \mathbf{Y}\mathbf{c}_i, c_{ii} = 0$$

where  $\mathbf{c}_i$  is a  $N$  dimensional vector of constants on the weights for the other data points and  $\|\cdot\|_q$  is the  $l_q$  norm.

Ideally, want to use  $l_0$  norm, but is NP-hard

Use  $l_1$  norm the convex problem is: For all  $i = 1, \dots, N$ ,

$$\min \|\mathbf{c}_i\|_1 \quad \text{s.t.} \quad \mathbf{y}_i = \mathbf{Y}\mathbf{c}_i, c_{ii} = 0$$

or in matrix form,

$$\min \|\mathbf{C}\|_1 \quad \text{s.t.} \quad \mathbf{Y} = \mathbf{Y}\mathbf{C}, \text{diag}(\mathbf{C}) = 0$$

where  $\mathbf{C} \equiv [\mathbf{c}_1, \dots, \mathbf{c}_N] \in \mathbb{R}^{N \times N}$  is the matrix whose  $i$ -th column corresponds to the representation coefficient for  $\mathbf{y}_i$ .

# Clustering using Sparse Coefficients

The optimization gives the sparse coefficient matrix  $\mathbf{C}$  as the output. Use the matrix to construct a graph, and apply spectral clustering gives the final clustering result.

Weighted graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{W})$ , where the  $N$  data points compose the vertices  $\mathcal{V}$ , and  $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$  denotes the set of edges between nodes,  $\mathbf{W}$  denotes the weights for the edges.

$$\mathbf{W} = |\mathbf{C}| + |\mathbf{C}|^T, \mathbf{W}_{ij} = |c_{ij}| + |c_{ji}|.$$

# Sparse Subspace Clustering Algorithm

## Algorithm : Sparse Subspace Clustering (SSC)

**Input:** A set of points  $\{\mathbf{y}_i\}_{i=1}^N$  lying in a union of  $n$  linear subspaces  $\{\mathcal{S}_l\}_{l=1}^n$ .

- 1 Solve the sparse optimization program
- 2 Normalize the columns of the resulting coefficients matrix  $\mathbf{C}$  as  $\mathbf{c}_i \leftarrow \frac{\mathbf{c}_i}{\|\mathbf{c}_i\|_\infty}$
- 3 Form a similarity graph with  $N$  nodes representing the data points. Set the weights on the edges between the nodes by  $\mathbf{W} = |\mathbf{C}| + |\mathbf{C}|^T$
- 4 Apply spectral clustering to the similarity graph.

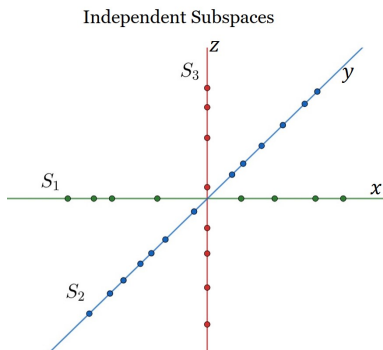
**Output:** Segmentation of the data:  $\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_n$

# Subspace-Sparse Recovery Theory

**Subspace-Sparse Representation:** represented by a few other points in the same subspace

## Definition (Independent Subspaces)

A collection of subspaces  $\{\mathcal{S}_i\}_{i=1}^n$  is said to be independent if  $\dim(\oplus_{i=1}^n \mathcal{S}_i) = \sum_{i=1}^n \dim(\mathcal{S}_i)$  where  $\oplus$  denotes the direct sum operator.



# Subspace-Sparse Recovery Theory

## Theorem (Independent Subspaces Recovery)

Consider a collection of data points drawn from  $n$  independent subspaces  $\{\mathcal{S}_i\}_{i=1}^n$  of dimensions  $\{d_i\}_{i=1}^n$ . Let  $\mathbf{Y}_i$  denote  $N_i$  data points in  $\mathcal{S}_i$ , where  $\text{rank}(\mathbf{Y}_i) = d_i$ , and let  $\mathbf{Y}_{-i}$  denote data points in all subspaces except  $\mathcal{S}_i$ . Then, for every  $\mathcal{S}_i$  and every nonzero  $\mathbf{y}$  in  $\mathcal{S}_i$ , the  $l_q$ -minimization program

$$\begin{bmatrix} \mathbf{c}^* \\ \mathbf{c}_-^* \end{bmatrix} = \underset{\mathbf{c}}{\text{argmin}} \left\| \begin{bmatrix} \mathbf{c} \\ \mathbf{c}_- \end{bmatrix} \right\|_q \quad \text{s.t.} \quad \mathbf{y} = [\mathbf{Y}_i \quad \mathbf{Y}_{-i}] \begin{bmatrix} \mathbf{c} \\ \mathbf{c}_- \end{bmatrix}$$

for  $q < \infty$ , recovers a subspace-sparse representation, i.e.,  $\mathbf{c}^* \neq \mathbf{0}$  and  $\mathbf{c}_-^* = \mathbf{0}$

## Practical Extension: Data Nuisances

Data points contaminated with sparse outlying entries and noise

$$\mathbf{y}_i = \mathbf{y}_i^0 + \mathbf{e}_i^0 + \mathbf{z}_i^0$$

where  $\mathbf{y}_i^0$  is the error-free part of observation and lies perfect on the underlying subspace,  $\mathbf{e}_i^0 \in \mathbb{R}^D$  is sparse outlying entries, and  $\mathbf{z}_i^0 \in \mathbb{R}^D$  is the noise component.

Utilizing self-expressiveness property of  $\mathbf{y}_i^0$

$$\mathbf{y}_i = \sum_{j \neq i} c_{ij} \mathbf{y}_j + \mathbf{e}_i + \mathbf{z}_i$$

$$\mathbf{e}_i = \mathbf{e}_i^0 - \sum_{j \neq i} c_{ij} \mathbf{e}_j^0$$

$$\mathbf{z}_i = \mathbf{z}_i^0 - \sum_{j \neq i} c_{ij} \mathbf{z}_j^0$$

## Practical Extension: Data Nuisances

Let matrices  $\mathbf{E}$  and  $\mathbf{Z}$  have  $\mathbf{e}_i$  and  $\mathbf{z}_i$  as columns, we have the matrix representation:

$$\mathbf{Y} = \mathbf{Y}\mathbf{C} + \mathbf{E} + \mathbf{Z}, \quad \text{diag}(\mathbf{C}) = 0$$

And hence the optimization program for sparse representation becomes

$$\begin{aligned} \min & \|\mathbf{C}\|_1 + \lambda_e \|\mathbf{E}\|_1 + \frac{\lambda_z}{2} \|\mathbf{Z}\|_F^2 \\ \text{s.t. } & \mathbf{Y} = \mathbf{Y}\mathbf{C} + \mathbf{E} + \mathbf{Z}, \quad \text{diag}(\mathbf{C}) = 0 \end{aligned}$$



## Practical Extension: Affine Subspace

Data point  $\mathbf{y}_i$  in an affine subspace  $\mathcal{S}_l$  with dimension  $d_l$  can be written as an affine combination of  $d_l + 1$  other points from  $\mathcal{S}_l$ .

$$\mathbf{y}_i = \mathbf{Y}\mathbf{c}_i, c_{ii} = 0, \mathbf{1}^T \mathbf{c}_i = 1$$

where  $\mathbf{c}_i$  has  $d_l + 1$  nonzero entries corresponding to points also in  $\mathcal{S}_l$ . The more general program:

$$\begin{aligned} \min \quad & \|\mathbf{C}\|_1 + \lambda_e \|\mathbf{E}\|_1 + \frac{\lambda_z}{2} \|\mathbf{Z}\|_F^2 \\ \text{s.t.} \quad & \mathbf{Y} = \mathbf{Y}\mathbf{C} + \mathbf{E} + \mathbf{Z}, \quad \text{diag}(\mathbf{C}) = 0, \quad \mathbf{1}^T \mathbf{C} = \mathbf{1}^T \end{aligned}$$

# Real Data Experiments: Motion Segmentation

**Dataset:** Hopkins 155 dataset. Video sequences along with the features extracted and tracked in all the frames.



The video has frames  $f : 1, \dots, F$ , and a set of  $N$  feature points  $\{\mathbf{x}_{fi} \in \mathbb{R}^2\}_{i=1}^N$  is tracked across the frames. For analysis, each feature trajectory  $\mathbf{y}_i$  is taken as a data point, where  $\mathbf{y}_i$  is obtained by stacking the feature points  $x_{fi}$  in the video as

$$\mathbf{y}_i \equiv [\mathbf{x}_{1i}^T, \mathbf{x}_{2i}^T, \dots, \mathbf{x}_{Fi}^T]^T \in \mathbb{R}^{2F}$$

120 videos of two motions ( $N = 266$  and  $F = 30$ )

35 videos of 3 motions ( $N = 398$  and  $F = 29$ )

The most general version of the SSC is implemented. With/without PCA as preprocessing.

## Clustering error (%) of different algorithms on the Hopkins 155 dataset

Algorithm	SSC	LSA	LRSC	K-Subspace
2 Motions, without PCA				
Mean	<b>2.23</b>	20.88	3.98	28.42
Median	<b>0</b>	14.72	<b>0</b>	33.72
3 Motions, without PCA				
Mean	<b>5.78</b>	21.09	7.96	32.27
Median	<b>0.95</b>	22.81	3.40	29.62
2 Motions, with PCA				
Mean	<b>2.32</b>	3.01	3.89	19.78
Median	<b>0</b>	0.25	<b>0</b>	1.62
3 Motions, with PCA				
Mean	5.78	<b>5.19</b>	7.88	21.34
Median	<b>0.95</b>	1.11	3.39	3.19

# Discussion

## Summary

- SSC algorithm
- Guarantee for obtaining Subspace-Sparse representation
- Practical Extensions
- Real application on motion segmentation

## Advantages:

- Deal with noises, sparse outlying entries, and affine subspaces directly
- Can deal with subspaces with different unknown dimensions

## Further Problem:

- Subspace-sparse recovery guarantee for corrupted data
- Theory for graph connectivity
- Application when number of clusters unknown
- Nonlinear extension
- Application to very large dataset

# Reference

- Elhamifar, E. and Vidal, R., 2013. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE transactions on pattern analysis and machine intelligence*, 35(11), pp.2765-2781.
- G. Liu and S. Yan, "Latent Low-Rank Representation for subspace segmentation and feature extraction," 2011 International Conference on Computer Vision, Barcelona, 2011, pp. 1615-1622, doi: 10.1109/ICCV.2011.6126422.

Question?

# Graph Connectivity

For subspaces of dimensions greater than or equal to 4, under odd distribution of data, it is possible that points in the same subspace form multiple components.

Consider a regularization term

$$\|\mathbf{C}\|_{r,0} \equiv \sum_{i=1}^N I(\|\mathbf{c}^i\|_2 > 0)$$

where  $\mathbf{c}^i$  denotes the  $i$ -th row of the matrix  $\mathbf{C}$ .

The convex relaxation

$$\|\mathbf{C}\|_{r,1} \equiv \sum_{i=1}^N \|\mathbf{c}^i\|_2$$

The optimization program to consider is

$$\min \|\mathbf{C}\|_1 + \lambda_r \|\mathbf{C}\|_{r,1} \quad s.t. \quad \mathbf{Y} = \mathbf{Y}\mathbf{C}, \text{diag}(\mathbf{C}) = 0$$

## Definition (Disjoint Subspaces)

A collection of subspaces  $\{\mathcal{S}_i\}_{i=1}^n$  is said to be disjoint if every pair of subspaces intersect only at the origin. In other words, for every pair of subspaces we have  $\dim(\mathcal{S}_i \oplus \mathcal{S}_j) = \dim(\mathcal{S}_i) + \dim(\mathcal{S}_j)$ , where  $\oplus$  denotes the direct sum operator.

To characterize two disjoint subspaces, we can have:

## Definition (Smallest Principal Angle)

The smallest principal angle between two subspaces  $\mathcal{S}_i$  and  $\mathcal{S}_j$ , denoted by  $\theta_{ij}$ , is defined as

$$\cos(\theta_{ij}) \equiv \max_{\mathbf{v}_i \in \mathcal{S}_i, \mathbf{v}_j \in \mathcal{S}_j} \frac{\mathbf{v}_i^T \mathbf{v}_j}{\|\mathbf{v}_i\|_2 \|\mathbf{v}_j\|_2}$$



# Subspace-Sparse Recovery Theory

For Disjoint subspaces, need to study points in the intersection of subspaces. Let  $\mathbf{x}$  be a non-zero vector in the intersection of  $\mathcal{S}_i$  and  $\bigoplus_{j \neq i} \mathcal{S}_j$ . Let

$$\begin{aligned} \mathbf{a}_i &= \operatorname{argmin} \|\mathbf{a}\|_1 \quad \text{s.t.} \quad \mathbf{x} = \mathbf{Y}_i \mathbf{a} \\ \mathbf{a}_{-i} &= \operatorname{argmin} \|\mathbf{a}\|_1 \quad \text{s.t.} \quad \mathbf{x} = \mathbf{Y}_{-i} \mathbf{a} \end{aligned}$$

## Theorem

Consider a collection of data points drawn from  $n$  independent subspaces  $\{\mathcal{S}_i\}_{i=1}^n$  of dimensions  $\{d_i\}_{i=1}^n$ . Let  $\mathbf{Y}_i$  denote  $N_i$  data points in  $\mathcal{S}_i$ , where  $\operatorname{rank}(\mathbf{Y}_i) = d_i$ , and let  $\mathbf{Y}_{-i}$  denote data points in all subspaces except  $\mathcal{S}_i$ . The  $l_1$  minimization

$$\begin{bmatrix} \mathbf{c}^* \\ \mathbf{c}_-^* \end{bmatrix} = \operatorname{argmin} \left\| \begin{bmatrix} \mathbf{c} \\ \mathbf{c}_- \end{bmatrix} \right\|_1 \quad \text{s.t.} \quad \mathbf{y} = [\mathbf{Y}_i \quad \mathbf{Y}_{-i}] \begin{bmatrix} \mathbf{c} \\ \mathbf{c}_- \end{bmatrix}$$

recovers a subspace-sparse representation, i.e.,  $\mathbf{c}^* \neq \mathbf{0}$  and  $\mathbf{c}_-^* = \mathbf{0}$  if and only if

$$\forall \mathbf{x} \in \mathcal{S}_i \cap \left(\bigoplus_{j \neq i} \mathcal{S}_j\right), \mathbf{x} \neq \mathbf{0} \implies \|\mathbf{a}_i\|_1 < \|\mathbf{a}_{-i}\|_1$$

# Parameter Choice

$$\mu_e \equiv \min_i \max_{j \neq i} \|\mathbf{y}_i\|_1, \quad \mu_z \equiv \min_i \max_{j \neq i} |\mathbf{y}_i^T \mathbf{y}_j|$$

The choose  $\lambda_e \geq \alpha_e / \mu_e$ ,  $\lambda_z \geq \alpha_z / \mu_z$  with  $\alpha_e, \alpha_z \geq 1$ .

## Proposition

*Consider the optimization program with noise and sparse outlying entries. Without the noise term  $\mathbf{Z}$ , if  $\lambda_e \leq 1 / \mu_e$ , then there exists at least one data point  $\mathbf{y}_l$  for which in the optimal solution we have  $(\mathbf{c}_l, \mathbf{e}_l) = (\mathbf{0}, \mathbf{y}_l)$ . That is, the data point is taken entirely as outlying entries and is not represented as combination of other data points. Similarly, without the sparse outlying term  $\mathbf{E}$ , if  $\lambda_z \leq \mu_z$ , then there exists at least one data point  $\mathbf{y}_l$  for which  $(\mathbf{c}_l, \mathbf{z}_l) = (\mathbf{0}, \mathbf{y}_l)$ .*

# Face Clustering Results

Algorithm	SSC	LSA	RANSAC	LRSC	KSubspace
2 Subjects					
Mean	<b>1.87</b>	33.47	37.45	11.41	45.41
Median	<b>0</b>	46.09	39.06	11.25	46.88
Time	57.5	13.7	4240.4	2.0	149.2
3 Subjects					
Mean	<b>3.30</b>	53.03	49.39	13.97	59.18
Median	<b>1.04</b>	51.06	50.52	13.87	59.90
Time	81.1	19.5	6815.0	3.1	536.0
5 Subjects					
Mean	<b>4.32</b>	58.82	62.10	21.58	67.59
Median	<b>2.50</b>	57.19	63.12	21.56	67.19
Time	135.8	29.4	12721.3	11.2	1115.9
8 Subjects					
Mean	<b>5.89</b>	57.41	787	34.73	72.00
Median	<b>4.59</b>	57.81	5415	34.37	71.58
Time	216.0	65.0	15901.7	19.5	2030.1
10 Subjects					
Mean	<b>7.40</b>	56.04	71.4	51.06	72.03
Median	<b>5.63</b>	60.47	72.5	50.78	72.34
Time	326.0	95.3	16393.7	59	3248.0

# ADMM Procedure

## Algorithm 2: ADMM Procedure to solve sparse-optimization program

**Initialization:** Set  $\text{maxIter} = N$  and  $\text{errThres} = \varepsilon$ .  $k = 0$ ,  $\text{Terminate} = \text{False}$ .  
Initialize  $\mathbf{C}^{(0)}$ ,  $\mathbf{E}^{(0)}$ ,  $\mathbf{A}^{(0)}$ ,  $\Delta^{(0)}$ ,  $\delta^{(0)}$  to zero.

**While** ( $\text{Terminate} == \text{False}$ ) **do**

1 Update  $\mathbf{A}^{(k+1)}$  by solving

$$(\lambda_z \mathbf{Y}^T \mathbf{Y} + \rho \mathbf{I} + \rho \mathbf{1}\mathbf{1}^T) \mathbf{A}^{(k+1)} = \lambda_z \mathbf{Y}^T (\mathbf{Y} - \mathbf{E}^{(k+1)}) + \rho (\mathbf{1}\mathbf{1}^T + \mathbf{C}^{(k)}) - \mathbf{1} \delta^{(k)T} - \Delta^{(k)}$$

2 Update  $\mathbf{C}^{(k+1)} = \mathbf{J} - \text{diag}(\mathbf{J})$ , where  $\mathbf{J} = \mathcal{T}_{\frac{1}{\rho}}(\mathbf{C}^{(k+1)} + \Delta^{(k)}/\rho)$ ,

3 Update  $\mathbf{E}^{(k+1)} = \mathcal{T}_{\frac{\lambda_e}{\lambda_z}}(\mathbf{Y}\mathbf{A}^{(k+1)} - \mathbf{Y})$ ,

4 Update  $\Delta^{(k+1)} = \Delta^{(k)} + \rho(\mathbf{A}^{(k+1)} - \mathbf{C}^{(k+1)})$ ,

5 Update  $\delta^{(k+1)} = \delta^{(k)} + \rho(\mathbf{A}^{(k+1)T} \mathbf{1} - \mathbf{1})$ ,

6 **if** ( $\max\{\|\mathbf{A}^{(k+1)} - \mathbf{C}^{(k+1)}\|_{\infty}, \|\mathbf{A}^{(k+1)T} \mathbf{1} - \mathbf{1}\|_{\infty}, \|\mathbf{A}^{(k+1)} - \mathbf{A}^{(k)}\|_{\infty}, \|\mathbf{E}^{(k+1)} - \mathbf{E}^{(k)}\|_{\infty}\} \leq \varepsilon$  or  $k + 1 \geq \text{maxIter}$ ):

$\text{Terminate} = \text{True}$

**end if**

7  $k = k + 1$ ,

**end while Output:** Sparse coefficient matrix  $\mathbf{C} = \mathbf{C}^{(k)}$