

# Subspace Clustering

Shihua Zhang

Fall 2019

# Overview

- 1 Background
- 2 Sparse Subspace Clustering
- 3 Other Self-Expressive Models
- 4 Subspace Clustering by Block Diagonal Representation
- 5 More than Linear Model

# Outline

## 1 Background

## 2 Sparse Subspace Clustering

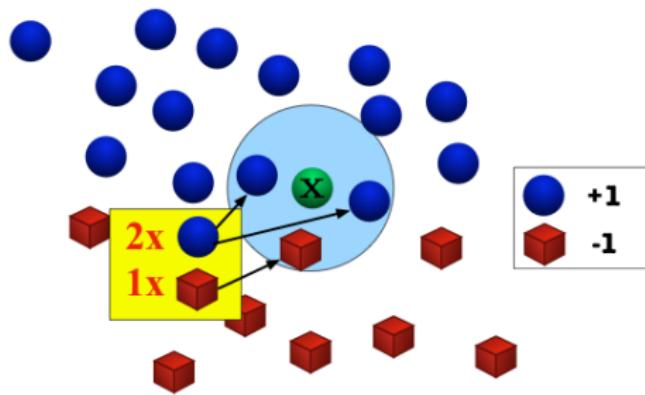
## 3 Other Self-Expressive Models

## 4 Subspace Clustering by Block Diagonal Representation

## 5 More than Linear Model

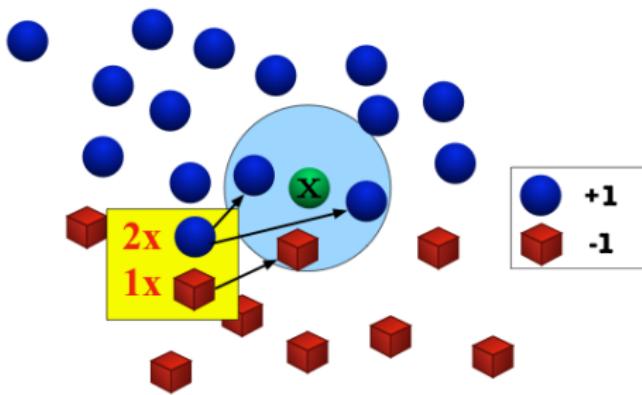
# Let's Begin with $k$ -NN

**$k$ -nearest neighbors algorithm ( $k$ -NN)** [Cover et al., 1967]:



# Let's Begin with $k$ -NN

**$k$ -nearest neighbors algorithm ( $k$ -NN)** [Cover et al., 1967]:



- **Assumption:** Similar inputs have similar outputs
- **Classification rule:** For a test input  $x$ , assign the most common label amongst its  $k$  most similar training inputs

# Bayes Optimal Classifier

## Bayes optimal classifier:

- Assume we know  $P(y|\mathbf{x})$ . What would you predict?
- **Examples:**  $y \in \{-1, 1\}$ ,  $P(+1|\mathbf{x}) = 0.8$ ,  $P(-1|\mathbf{x}) = 0.2$ .
- **Best prediction:**  $y^* = h_{\text{opt}} = \arg \max_y P(y|\mathbf{x})$
- You predict the most likely class.

# Bayes Optimal Classifier

## Bayes optimal classifier:

- Assume we know  $P(y|\mathbf{x})$ . What would you predict?
- **Examples:**  $y \in \{-1, 1\}$ ,  $P(+1|\mathbf{x}) = 0.8$ ,  $P(-1|\mathbf{x}) = 0.2$ .
- **Best prediction:**  $y^* = h_{\text{opt}} = \arg \max_y P(y|\mathbf{x})$
- You predict the most likely class.

## Error of the BayesOpt classifier:

$$\epsilon_{\text{BayesOpt}} = 1 - P(h_{\text{opt}}(\mathbf{x})|\mathbf{x}) = 1 - P(y^*|\mathbf{x})$$

- In our example, that is  $\epsilon_{\text{BayesOpt}} = 0.2$ .
- You can never do better than the Bayes Optimal Classifier.

# 1-NN Convergence Proof

## Theorem 1

*As  $n \rightarrow \infty$ , the 1-NN error is no more than twice the error of the Bayes Optimal classifier. (Similar guarantees hold for  $k > 1$ .)*

# 1-NN Convergence Proof

## Theorem 1

As  $n \rightarrow \infty$ , the 1-NN error is no more than twice the error of the Bayes Optimal classifier. (Similar guarantees hold for  $k > 1$ .)

## Proof.

Let  $\mathbf{x}_{NN}$  be the nearest neighbor of our test point  $\mathbf{x}_t$ . As  $n \rightarrow \infty$ ,  $\text{dist}(\mathbf{x}_{NN}, \mathbf{x}) \rightarrow 0$ , i.e.  $\mathbf{x}_{NN} \rightarrow \mathbf{x}_t$ . Then we return the label of  $\mathbf{x}_{NN}$ . What is the probability that this is not the label of  $\mathbf{x}$ ?

$$\begin{aligned}\epsilon_{NN} &= P(y^*|\mathbf{x}_t)(1 - P(y^*|\mathbf{x}_{NN})) + P(y^*|\mathbf{x}_{NN})(1 - P(y^*|\mathbf{x}_t)) \\ &\leq (1 - P(y^*|\mathbf{x}_{NN})) + (1 - P(y^*|\mathbf{x}_t)) = 2(1 - P(y^*|\mathbf{x}_t)) = 2\epsilon_{\text{BayesOpt}}.\end{aligned}$$



# Test 1-NN on CIFAR-10

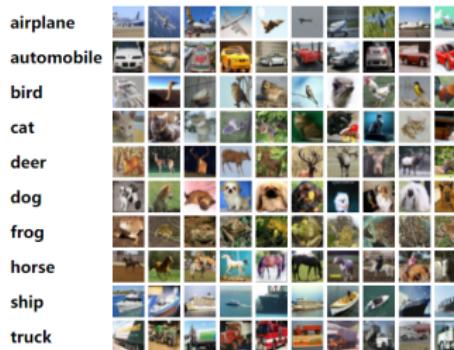
**Good news:** As  $n \rightarrow \infty$ , the 1-NN classifier is only a factor 2 worse than the best possible classifier. It seems that all of classification problems can be solved perfectly via  $k$ -NN!

# Test 1-NN on CIFAR-10

**Good news:** As  $n \rightarrow \infty$ , the 1-NN classifier is only a factor 2 worse than the best possible classifier. It seems that all of classification problems can be solved perfectly via  $k$ -NN!

**Let's have a try on the CIFAR-10 Dataset:**

- 60000  $32 \times 32$  colour images in 10 classes (6000 per class).
- There are 50000 training images and 10000 test images.



# Test 1-NN on CIFAR-10

## Actual performance:

- The accuracy on training set: 100%!
- The accuracy on testing set: 38.6%!

# Test 1-NN on CIFAR-10

## Actual performance:

- The accuracy on training set: 100%!
- The accuracy on testing set: 38.6%!

What's wrong with  $k$ -NN?

# Test 1-NN on CIFAR-10

## Actual performance:

- The accuracy on training set: 100%!
- The accuracy on testing set: 38.6%!

What's wrong with  $k$ -NN?

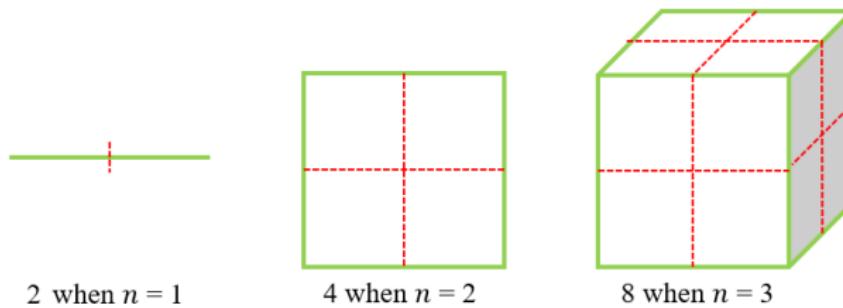
We are cursed by high dimensionality!

- Insufficient Data
- Invalid Distance
- Redundant Information

# Insufficient Data

**60000 images in CIFAR-10, still insufficient?**

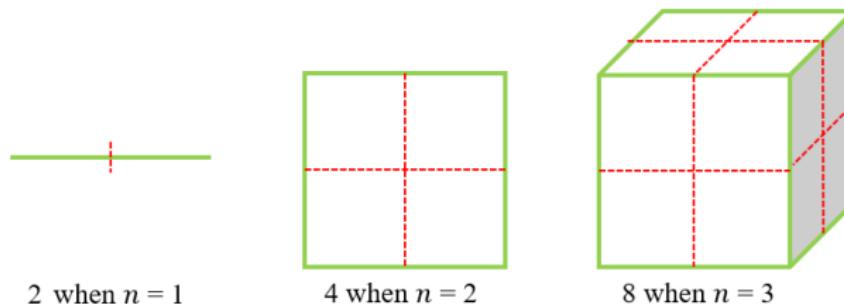
**Problem:** Given a cubic in  $\mathbb{R}^n$ , divide each dimension into 2 parts and put a data point in each sub-cubic. How many data points do we need?



# Insufficient Data

**60000 images in CIFAR-10, still insufficient?**

**Problem:** Given a cubic in  $\mathbb{R}^n$ , divide each dimension into 2 parts and put a data point in each sub-cubic. How many data points do we need?

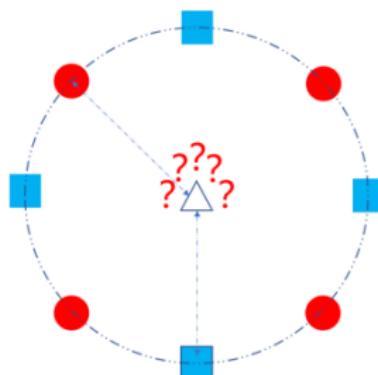
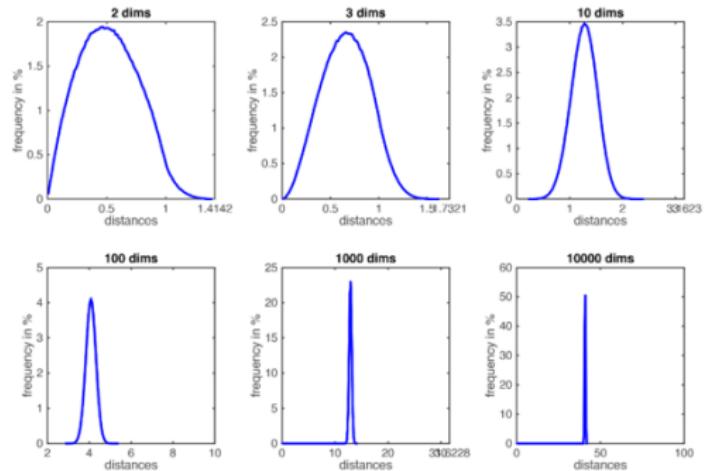


What if  $n = 32 \times 32 \times 3$ ?

**We need  $2^{3072} \simeq 10^{900}$  data points!!!**

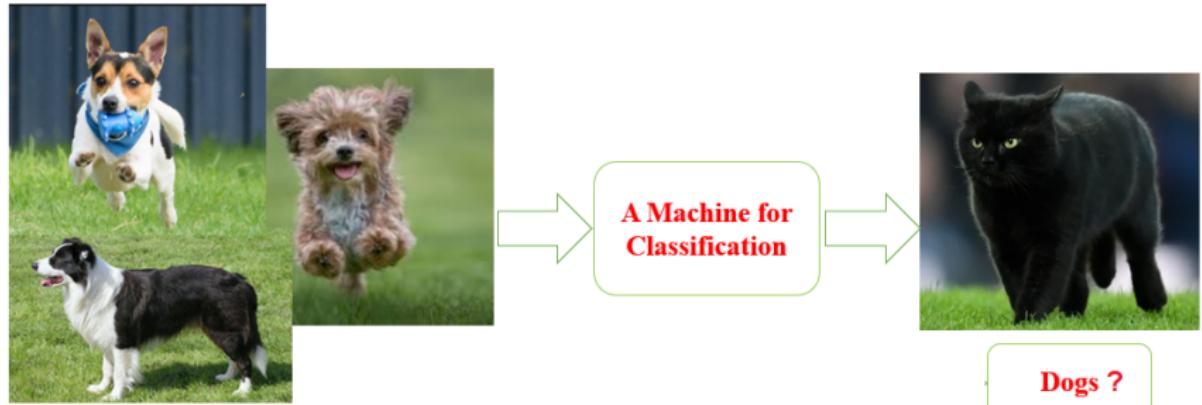
# Invalid Distance

As the number of dimensions  $n$  grows, all distances concentrate within a very small range.



Traditional Euclidean distance will lose its meaning!

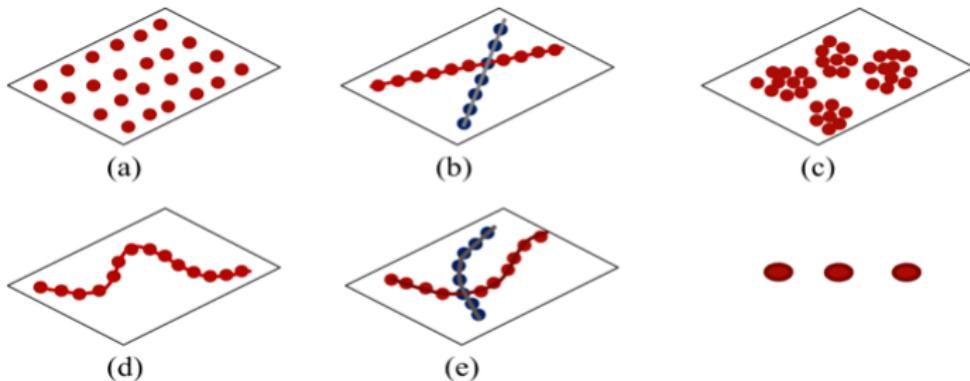
# Redundant Information



Redundant information may 'dilute' the main information!

# Blessings of Dimensionality

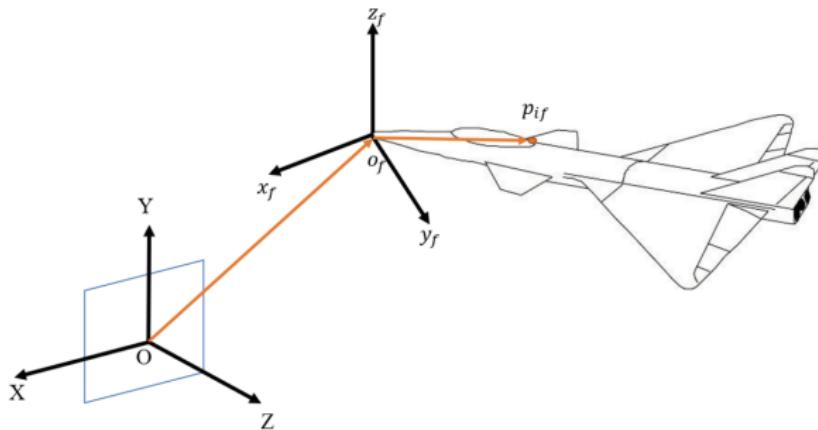
- **Bad news:** We are cursed!!
- **Good news:** Low dimensionality structure is hidden from high dimensionality
  - Principal Component Analysis
  - **Subspace Clustering**
  - Manifold Learning
  - ...



# Motion Segmentation

Let us consider how a static camera observes a single moving object:

- Let  $\mathbf{p}_{if}$  be the 2D coordinate of point  $\mathbf{p}_i$  at  $f$ -th frame.
- $\mathbf{x}_i = [\mathbf{p}_{i1}, \mathbf{p}_{i2}, \dots, \mathbf{p}_{iF}]^T$  represents the trajectory of point  $\mathbf{p}_i$  in  $F$  frames.



# Motion Segmentation

Let  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]_{2F \times N}$  be the data matrix representing trajectories of  $N$  points. Then we have:

$$\text{rank}(X) = 4$$

In other words, despite each point in  $X$  with dimension of  $2F$ , the ambient dimension is only 4 [Costeira and Kanade, 1995]!

# Motion Segmentation

- **Data:** Given feature points on multiple rigidly moving objects tracked in multiple frames of a video.
- **Task:** Separate the feature trajectories according to the moving objects.
- **Assumption:** Point trajectories associated with each moving object across multiple frames lie in a linear subspace of dimension at most 4.



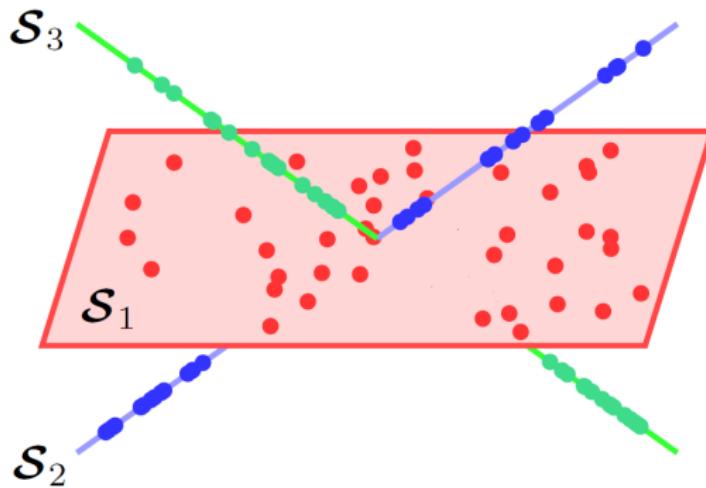
# Face Clustering

- **Data:** Face images of multiple subjects.
- **Task:** Find images that belong to the same subject.
- **Assumption:** The set of all reflectance functions produced by Lambertian objects under distant, isotropic lighting lies close to a 9D linear subspace [Basri and Jacobs, 2003].



# What is Subspace Clustering?

- **Problem:** Given a set of sufficient amount of data points  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$ , drawn from a union of  $L$  linear or affine subspaces  $\{\mathcal{S}\}_{i=1}^L$  of unknown dimensions, segment all data vectors into their respective subspaces.



# Prior Work on Subspace Clustering

Existing algorithms can be divided into four main categories:

- Iterative Methods [Tseng, 2000]
- Algebraic Approaches [Kanatani, 2001]
- Statistical Methods [Tipping and Bishop, 1999]
- **Self-Expressiveness Method** [Elhamifar and Vidal, 2009]

# Prior Work on Subspace Clustering

Existing algorithms can be divided into four main categories:

- Iterative Methods [Tseng, 2000]
- Algebraic Approaches [Kanatani, 2001]
- Statistical Methods [Tipping and Bishop, 1999]
- **Self-Expressiveness Method** [Elhamifar and Vidal, 2009]

**Self-Expressiveness** method can handle noise and outliers in the data well and doesn't need to know the dimension of each subspace in advance.

# Summary

- The Curses of Dimensionality
  - Insufficient Data
  - Invalid Distance
  - Redundant Information

# Summary

- The Curses of Dimensionality
  - Insufficient Data
  - Invalid Distance
  - Redundant Information
- The Blessings of Dimensionality
  - Low dimensionality structure is hidden from high dimensionality

# Summary

- The Curses of Dimensionality
  - Insufficient Data
  - Invalid Distance
  - Redundant Information
- The Blessings of Dimensionality
  - Low dimensionality structure is hidden from high dimensionality
- Subspace Clustering Model
  - Motion Segmentation
  - Face Clustering

# Outline

- 1 Background
- 2 Sparse Subspace Clustering
- 3 Other Self-Expressive Models
- 4 Subspace Clustering by Block Diagonal Representation
- 5 More than Linear Model

# Self-Expressiveness Method

- **Idea:** each data point in a union of subspaces can be efficiently represented as a linear or affine combination of other points.
- **Difficulty:** This is an ill-posed problem.

# Self-Expressiveness Method

- **Idea:** each data point in a union of subspaces can be efficiently represented as a linear or affine combination of other points.
- **Difficulty:** This is an ill-posed problem.

## Example 2

Given a data matrix  $X = \begin{bmatrix} 0 & 0 & 1 & 2 \\ 1 & 2 & 0 & 0 \end{bmatrix}$ , then  $\begin{bmatrix} 0 \\ 2 \end{bmatrix}$  may be represented as  $2 \begin{bmatrix} 0 \\ 1 \end{bmatrix} - 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 \\ 0 \end{bmatrix}$ .

- **Inspiration:** Compressed Sensing

## Inspired by Compressed Sensing

$$y = A x$$

$$\begin{aligned} & \min_{\mathbf{x}} \|\mathbf{x}\|_1 \\ \text{s.t. } & \mathbf{y} = \mathbf{A}\mathbf{x} \end{aligned}$$

- $l_1$  norm not only forces  $\mathbf{x}$  to be sparse, but also guarantee the uniqueness of the solution.

# Sparse Subspace Clustering

**Problem:** Given a set of sufficient amount of data points  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$ , drawn from a union of  $n$  linear or affine subspaces  $\{\mathcal{S}\}_{i=1}^n$  of unknown dimensions, segment all data vectors into their respective subspaces.

# Sparse Subspace Clustering

**Problem:** Given a set of sufficient amount of data points  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$ , drawn from a union of  $n$  linear or affine subspaces  $\{\mathcal{S}\}_{i=1}^n$  of unknown dimensions, segment all data vectors into their respective subspaces.

**Goal:** aim to find a matrix  $Z \in \mathbb{R}^{N \times N}$ , such that:

$$\begin{aligned} & \min_Z \quad \|Z\|_1 \\ \text{s.t. } & X = XZ, \text{ diag}(Z) = 0 \end{aligned}$$

# Sparse Subspace Clustering

**Problem:** Given a set of sufficient amount of data points  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$ , drawn from a union of  $n$  linear or affine subspaces  $\{\mathcal{S}\}_{i=1}^n$  of unknown dimensions, segment all data vectors into their respective subspaces.

**Goal:** aim to find a matrix  $Z \in \mathbb{R}^{N \times N}$ , such that:

$$\begin{aligned} & \min_Z \|Z\|_1 \\ \text{s.t. } & X = XZ, \text{ diag}(Z) = 0 \end{aligned}$$

- $X = XZ$ :  $X$  is a self-expressive dictionary in which each point can be written as a linear combination of other points.
- $\min_Z \|Z\|_1$ : guarantee the sparsity and uniqueness.
- $\text{diag}(Z) = 0$ : avoid trivial solution  $I$ .

# Sparse Subspace Clustering

**Noiseless data in linear subspace:**

$$\begin{aligned} & \min_Z \|Z\|_1 \\ \text{s.t. } & X = XZ, \text{ diag}(Z) = 0 \end{aligned}$$

# Sparse Subspace Clustering

**Noiseless data in linear subspace:**

$$\begin{aligned} & \min_Z \|Z\|_1 \\ \text{s.t. } & X = XZ, \text{ diag}(Z) = 0 \end{aligned}$$

**Noiseless data in affine subspace:**

$$\begin{aligned} & \min_Z \|Z\|_1 \\ \text{s.t. } & X = XZ, Z^T \mathbf{1} = \mathbf{1}, \text{ diag}(Z) = 0 \end{aligned}$$

# Sparse Subspace Clustering

**Noiseless data in linear subspace:**

$$\begin{aligned} & \min_Z \|Z\|_1 \\ \text{s.t. } & X = XZ, \text{ diag}(Z) = 0 \end{aligned}$$

**Noiseless data in affine subspace:**

$$\begin{aligned} & \min_Z \|Z\|_1 \\ \text{s.t. } & X = XZ, Z^T \mathbf{1} = \mathbf{1}, \text{ diag}(Z) = 0 \end{aligned}$$

**Noise and sparse outlying entries:**

$$\begin{aligned} \min & \quad \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|S\|_F^2 \\ \text{s.t. } & X = XZ + E + S, \text{ diag}(Z) = 0 \end{aligned}$$

# Sparse Subspace Clustering

Let's first consider the most general situation:

$$\begin{array}{ll}\min & \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|S\|_F^2 \\ \text{s.t.} & X = XZ + E + S, \text{diag}(Z) = 0\end{array}$$

# Sparse Subspace Clustering

Let's first consider the most general situation:

$$\begin{aligned} \min \quad & \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|S\|_F^2 \\ \text{s.t.} \quad & X = XZ + E + S, \text{diag}(Z) = 0 \end{aligned}$$

- How to optimize?
- Get an optimal  $Z$ , what's next?
- Is there any theoretical guarantee?

# Algorithm for this Sparse Optimization

Adopt an Alternating Direction Method of Multipliers (ADMM) and first consider:

$$\begin{aligned} & \min_{(Z, E, S)} \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|S\|_F^2 \\ \text{s.t. } & X = XZ + E + S, \quad Z^T 1 = 1, \quad \text{diag}(Z) = 0 \end{aligned}$$

Eliminate  $S$  from the optimization problem:

$$\begin{aligned} & \min_{(Z, E)} \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|X - XZ - E\|_F^2 \\ \text{s.t. } & Z^T 1 = 1, \quad \text{diag}(Z) = 0 \end{aligned}$$

Introduce an auxiliary matrix  $A \in \mathbb{R}^{N \times N}$  and consider:

$$\begin{aligned} & \min_{(Z, E, A)} \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|X - XA - E\|_F^2 \\ \text{s.t. } & A^T 1 = 1, \quad A = Z - \text{diag}(Z) \end{aligned}$$

# Algorithm for this Sparse Optimization

Add two penalty terms:

$$\begin{aligned} \min_{(Z, E, A)} & \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|X - XA - E\|_F^2 \\ & + \frac{\rho}{2} \|A^T \mathbf{1} - \mathbf{1}\|_2^2 + \frac{\rho}{2} \|A - (Z - \text{diag}(Z))\|_F^2 \\ \text{s.t. } & A^T \mathbf{1} = \mathbf{1}, A = Z - \text{diag}(Z) \end{aligned}$$

Introduce a vector  $\delta \in \mathbb{R}^N$  and a matrix  $\Delta \in \mathbb{R}^{N \times N}$  of Lagrange multipliers for the two equality constraints:

$$\begin{aligned} \mathcal{L}(Z, A, E, \delta, \Delta) = & \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|X - XA - E\|_F^2 \\ & + \frac{\rho}{2} \|A^T \mathbf{1} - \mathbf{1}\|_2^2 + \frac{\rho}{2} \|A - (Z - \text{diag}(Z))\|_F^2 \\ & + \delta^T (A^T \mathbf{1} - \mathbf{1}) + \text{tr} (\Delta^T (A - Z + \text{diag}(Z))) \end{aligned}$$

All of sub-problems can be solved efficiently!

# ADMM Algorithm

---

**Algorithm : Solving SSC via an ADMM Algorithm**

---

**Initialization:** Set maxIter =  $10^4$ ,  $k = 0$ , and Terminate  $\leftarrow$  False. Initialize  $Z^{(0)}$ ,  $A^{(0)}$ ,  $E^{(0)}$ ,  $\delta^{(0)}$ , and  $\Delta^{(0)}$  to zero.

```
1: while (Terminate == False) do
2:     update  $A^{(k+1)}$  by solving the following system of linear equations
       
$$(\lambda_s Y^\top Y + \rho I + \rho 11^\top) A^{(k+1)} = \lambda_s Y^\top (Y - E^{(k)}) + \rho(11^\top + Z^{(k)}) - 1\delta^{(k)\top} - \Delta^{(k)},$$

3:     update  $Z^{(k+1)}$  as  $Z^{(k+1)} = J - \text{diag}(J)$ , where  $J \triangleq \mathcal{T}_{\frac{1}{\rho}}(A^{(k+1)} + \Delta^{(k)}/\rho)$ ,
4:     update  $E^{(k+1)}$  as  $E^{(k+1)} = \mathcal{T}_{\lambda_e}(Y - YA^{(k+1)})$ ,
5:     update  $\delta^{(k+1)}$  as  $\delta^{(k+1)} = \delta^{(k)} + \rho(A^{(k+1)\top} 1 - 1)$ ,
6:     update  $\Delta^{(k+1)}$  as  $\Delta^{(k+1)} = \Delta^{(k)} + \rho(A^{(k+1)} - Z^{(k+1)})$ ,
7:      $k \leftarrow k + 1$ ,
8:     if ( $\|A^{(k)\top} 1 - 1\|_\infty \leq \epsilon$  and  $\|A^{(k)} - Z^{(k)}\|_\infty \leq \epsilon$  and  $\|A^{(k)} - A^{(k-1)}\|_\infty \leq \epsilon$  and  $\|E^{(k)} - E^{(k-1)}\|_\infty \leq \epsilon$  or ( $k \geq \text{maxIter}$ )) then
9:         Terminate  $\leftarrow$  True
10:    end if
11: end while
```

**Output:** Optimal sparse coefficient matrix  $Z^* = Z^{(k)}$ .

---

# Clustering using Sparse Coefficients

$$\begin{aligned} \min \quad & \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|S\|_F^2 \\ \text{s.t.} \quad & X = XZ + E + S, \quad \text{diag}(Z) = 0 \end{aligned}$$

Obtain a sparse representation for each data point whose nonzero elements ideally correspond to points from the same subspace.

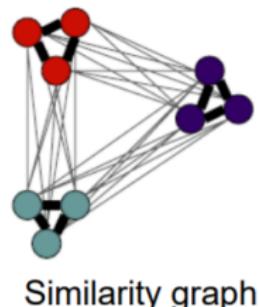
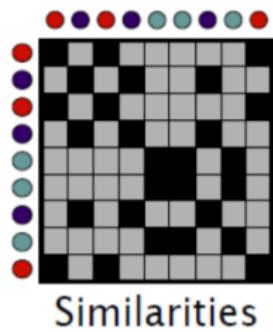
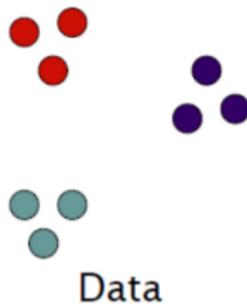
**What's next?**

Spectral clustering!

# What is spectral clustering?

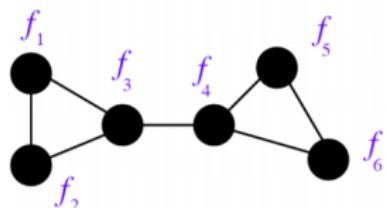
## Key steps:

- Take a similarity graph
- Construct its graph Laplacian matrix
- Do something with its bottom eigenvectors to get clusters



# What is Spectral Clustering?

**Idea: within-similarity high, between similarity low**



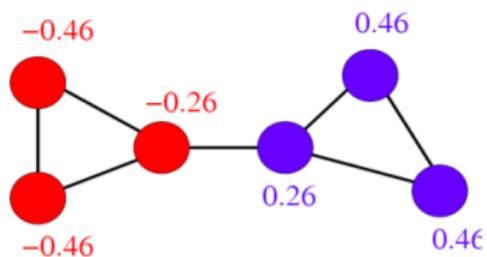
$$\mathbf{L} = \begin{pmatrix} 2 & -1 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 \\ -1 & -1 & 3 & -1 & 0 & 0 \\ 0 & 0 & -1 & 3 & -1 & -1 \\ 0 & 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & -1 & 2 \end{pmatrix} = \mathbf{D} - \mathbf{A}$$

$$\operatorname{argmin}_S \sum_{i \in S, j \in V-S} w_{ij} = \operatorname{argmin}_{f_i \in \{-1, 1\}} \sum_{i \sim j} w_{ij} (f_i - f_j)^2 = \frac{1}{8} \operatorname{argmin}_{f_i \in \{-1, 1\}} \mathbf{f}^T \mathbf{L} \mathbf{f}$$

Relaxation of integrality constraint + optimality give eigenvectors:

$$\mathbf{L} \mathbf{f} = \lambda \mathbf{f}$$

# What is Spectral Clustering?



$$L = \begin{pmatrix} 2 & -1 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 \\ -1 & -1 & 3 & -1 & 0 & 0 \\ 0 & 0 & -1 & 3 & -1 & -1 \\ 0 & 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & -1 & 2 \end{pmatrix}$$

$$L\mathbf{f}_1 = \lambda_1 \mathbf{f}_1 \quad \mathbf{f}_1 = [-0.46, -0.46, -0.26, 0.26, 0.46, 0.46]$$

- This is a traditional spectral clustering [Von Luxburg, 2007].
- It can be easily extended more general situation.

# The Framework of Sparse Subspace Clustering

- 1) Solve the sparse optimization problem:

$$\begin{aligned} \min \quad & \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|S\|_F^2 \\ \text{s.t.} \quad & X = XZ + E + S, \quad \text{diag}(Z) = 0 \end{aligned}$$

- 2) Normalize the columns of  $Z$  as  $\mathbf{z}_i \leftarrow \frac{\mathbf{z}_i}{\|\mathbf{z}_i\|_\infty}$
- 3) Form a similarity graph with  $N$  nodes representing the data points. Set the weights on the edges between the nodes by  $W = |Z| + |Z|^T$ .
- 4) Apply spectral clustering to the similarity graph.

# Subspace Sparse Recovery Theory

**Recovery conditions** for two classes of subspace arrangements:

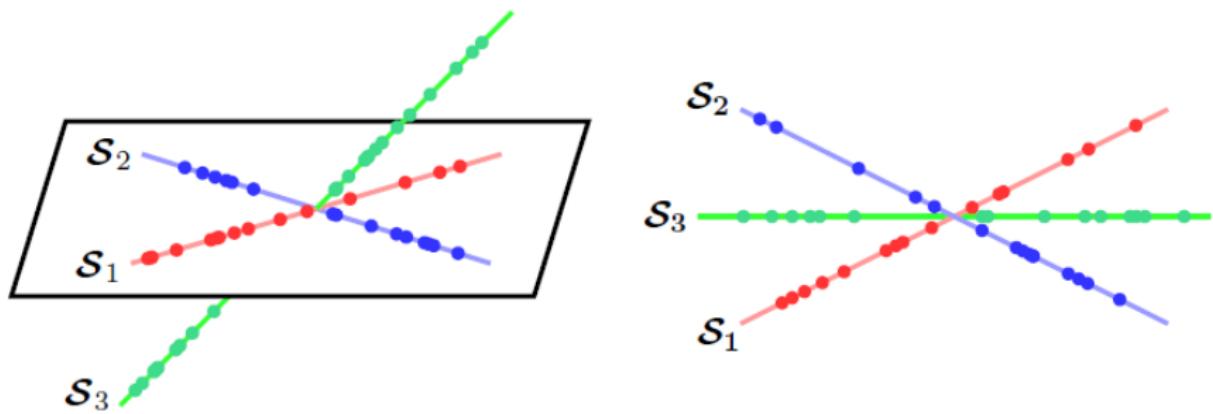
## Definition 3

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

## Definition 4

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  $\dim(\mathcal{S}_i \oplus \mathcal{S}_j) = \dim(\mathcal{S}_i) + \dim(\mathcal{S}_j)$ .

# Subspace Sparse Recovery Theory



**Figure:** Left: the three 1-dimensional subspaces are independent as they span the 3-dimensional space and the sum of their dimensions is also 3. Right: the three 1-dimensional are disjoint as any two subspaces intersect at the origin.

# Subspace Sparse Recovery Theory

**Key:** an important notion that can be used to characterize two disjoint subspaces is the smallest principal angle:

## Definition 5

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

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

# Independent Subspace Model

## Theorem 6

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  $X_i$  denote  $N_i$  data points in  $\mathcal{S}_i$ , where  $\text{rank}(X_i) = d_i$ , and let  $X_{-i}$  denote data points in all subspaces except  $\mathcal{S}_i$ . Then, for every  $\mathcal{S}_i$  and every nonzero  $\mathbf{x}$  in  $\mathcal{S}_i$ , the  $l_q$ -minimization program:

$$\begin{bmatrix} \mathbf{z}^* \\ \mathbf{z}_{-}^* \end{bmatrix} = \operatorname{argmin} \left\| \begin{bmatrix} \mathbf{z} \\ \mathbf{z}_{-} \end{bmatrix} \right\|_q \quad \text{s.t.} \quad \mathbf{x} = [X_i \ X_{-i}] \begin{bmatrix} \mathbf{z} \\ \mathbf{z}_{-} \end{bmatrix}$$

for  $q \geq 1$ , recovers a subspace-sparse representation, i.e.,  $\mathbf{z}^* \neq \mathbf{0}$  and  $\mathbf{z}_{-}^* = \mathbf{0}$ .

# Disjoint Subspace Model

## Theorem 7

Consider a collection of data points drawn from  $n$  disjoint subspaces  $\{\mathcal{S}_i\}_{i=1}^n$  of dimensions  $\{d_i\}_{i=1}^n$ . Let  $X_i$  denote  $N_i$  data points in  $\mathcal{S}_i$ , where  $\text{rank}(X_i) = d_i$ , and let  $X_{-i}$  denote data points in all subspaces except  $\mathcal{S}_i$ . Let  $\mathcal{W}_i$  be the set of all full-rank submatrices  $X_i \in \mathbb{R}^{D \times d_i}$  of  $X_i$ , where  $\text{rank}(X_i) = d_i$ . If the condition:

$$\max_{\tilde{X}_i \in \mathcal{W}_i} \sigma_{d_i}(\tilde{X}_i) > \sqrt{d_i} \|X_{-i}\|_{1,2} \max_{j \neq i} \cos(\theta_{ij})$$

holds, then for every nonzero  $\mathbf{x}$  in  $\mathcal{S}_i$ , the  $l_1$ -minimization recovers a subspace-sparse solution, i.e.,  $\mathbf{z}^* \neq \mathbf{0}$  and  $\mathbf{z}_{-i}^* = \mathbf{0}$ .

# Experiments with Synthetic Data

- 3 disjoint subspaces  $\{\mathcal{S}_i\}_{i=1}^3$
- Principal angles  $\theta$
- Number of data points  $N_g$
- Subspace sparse recovery error (ssr error):

$$\frac{1}{3N_g} \sum_{i=1}^{3N_g} \left( 1 - \frac{\|\mathbf{z}_{ik_i}\|_1}{\|\mathbf{z}_i\|_1} \right)$$

- Subspace clustering error:

$$\frac{\text{\# of misclassified points}}{\text{total \# of points}}$$

# Experiments with Synthetic Data

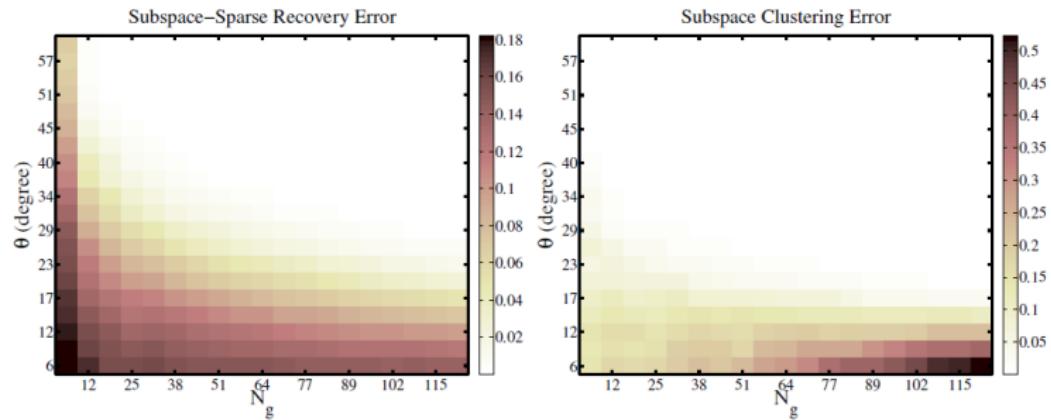


Figure: Subspace-sparse recovery error (left) and subspace clustering error (right) for three disjoint subspaces

- In general, when  $\theta$  or  $N_g$  increases, the errors decrease.
- Note that for small  $\theta$  as we increase  $N_g$ , the subspace sparse recovery error is large and slightly decreases, while the clustering error increases.

# Summary

The framework of sparse subspace clustering:

1) Formulate into the following problem:

$$\begin{aligned} \min \quad & \|Z\|_1 + \lambda_e \|E\|_1 + \frac{\lambda_s}{2} \|S\|_F^2 \\ \text{s.t.} \quad & X = XZ + E + S, \text{ diag}(Z) = 0 \end{aligned}$$

2) Normalize the columns of  $Z$

3) Form a similarity graph

4) Apply spectral clustering

- Self-expressive
- Handel noise and outliers
- Good theoretical guarantee

# Outline

- 1 Background
- 2 Sparse Subspace Clustering
- 3 Other Self-Expressive Models
- 4 Subspace Clustering by Block Diagonal Representation
- 5 More than Linear Model

# Weakness of Sparse Subspace Clustering

**Weakness: maybe too sparse in the same class**



# Weakness of Sparse Subspace Clustering

**Weakness: maybe too sparse in the same class**



**Solution:**

- Robust Subspace Clustering [Soltanolkotabi et al., 2014]
- Low Rank Representation [Liu et al., 2010]
- Least Squares Subspace Clustering [Lu et al., 2012]

# Robust Subspace Clustering

## Definition 8

Fix  $i$  and  $j \in \{1, \dots, n\}$  and let  $Z$  be the similarity matrix. Then we say that  $(i, j)$  obeying  $Z_{ij} \neq 0$  is **a false discovery** if  $\mathbf{x}_i$  and  $\mathbf{x}_j$  do not originate from the same subspace.

## Definition 9

In the same situation,  $(i, j)$  obeying  $Z_{ij} \neq 0$  is **a true discovery** if  $\mathbf{x}_i$  and  $\mathbf{x}_j$  originate from the same subspace.

**Main Idea:** We wish to make few false discoveries (and not link too many pairs belonging to different subspaces). At the same time, we wish to make many true discoveries, whence a natural trade off.

# Robust Subspace Clustering

## LASSO with data-driven regularization:

For each  $\mathbf{x}_i$ , consider a natural sparse regression strategy LASSO:

$$\min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x}_i - \mathbf{X}\mathbf{z}_i\|_2^2 + \lambda \|\mathbf{z}_i\|_1$$

$\lambda$  controls sparsity in  $\mathbf{z}_i$ .

We need to choose  $\lambda$  wisely:

- Take  $\lambda$  as large as possible (as to prevent false discoveries)
- Make sure that the number of true discoveries is also on the order of the dimension  $d$ , typically in the range  $[0.5d, 0.8d]$ , where  $d$  is the ambient dimension of  $\mathbf{x}_i$ .

# Robust Subspace Clustering

## How to Choose $\lambda$ :

Let  $\hat{\mathbf{z}}$  be the minimizer of the LASSO function:

$$K(\mathbf{z}, \lambda) = \frac{1}{2} \|\mathbf{x} - \mathbf{Xz}\|_2^2 + \lambda \|\mathbf{z}\|_1$$

Let  $\hat{\mathbf{z}}_{eq}$  be the solution of the problem when  $\lambda = 0$ .

Then we obtain the following inequation:

$$\frac{1}{2} \|\mathbf{x} - \mathbf{X}\hat{\mathbf{z}}\|_2^2 \leq K(\hat{\mathbf{z}}, \lambda) \leq K(\hat{\mathbf{z}}_{eq}, \lambda) = \lambda \|\hat{\mathbf{z}}_{eq}\|_1$$

# How to Choose $\lambda$

- $\lambda$  has to scale at least like  $1/\sqrt{d}$ 
  - To make sure  $\hat{\mathbf{z}}$  has a number of nonzero components in the range  $[0.5d, 0.8d]$ , then  $\|\mathbf{x} - X\hat{\mathbf{z}}\|_2^2$  has to be greater than a fixed constant.
  - $\|\hat{\mathbf{z}}_{eq}\|_1$  is on the order of  $\sqrt{d}$

# How to Choose $\lambda$

- $\lambda$  has to scale at least like  $1/\sqrt{d}$ 
  - To make sure  $\hat{\mathbf{z}}$  has a number of nonzero components in the range  $[0.5d, 0.8d]$ , then  $\|\mathbf{x} - X\hat{\mathbf{z}}\|_2^2$  has to be greater than a fixed constant.
  - $\|\hat{\mathbf{z}}_{eq}\|_1$  is on the order of  $\sqrt{d}$
- $\lambda$  has to scale at most like  $\sqrt{(\log N)/d}$ 
  - $\hat{\mathbf{z}} = \mathbf{0}$  if  $\lambda \geq \|X^T \mathbf{x}\|_\infty$
  - $\|X^T \mathbf{x}\|_\infty$  scales at most like  $\sqrt{(\log N)/d}$

# Experiments with Synthetic Data

- A single subspace in  $\mathbb{R}^n$  with  $n = 2000$ .
- 5, 4, 3, 4, 4, and 2 subspaces of respective dimensions 200, 150, 100, 50, 20 and 10.
- Noise level  $\sigma = 0.3$ .
- For each data point, we choose different values of  $\lambda$  around the heuristic  $\lambda_0 = 1/\sqrt{d}$ , namely,  $\lambda \in [0.1\lambda_0, 2.5\lambda_0]$ .

# Experiments with Synthetic Data

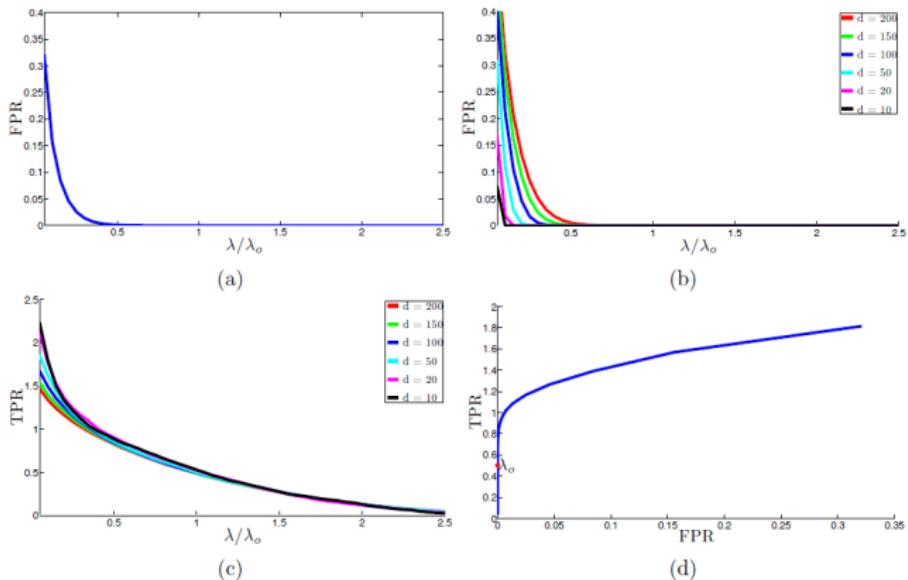


Figure 2: Performance of LASSO for values of  $\lambda$  in an interval including the heuristic  $\lambda_o = 1/\sqrt{d}$ . (a) Average number of false discoveries normalized by  $(n - d)$  (FPR) on all  $m$  sampled data points. (b) FPR for different subspace dimensions. Each curve represents the average FPR over those samples originating from subspaces of the same dimension. (c) Average number of true discoveries per dimension for various dimensions (TPR). (d) TPR vs. FPR (ROC curve). The point corresponding to  $\lambda = \lambda_o$  is marked as a red dot.

# How to Estimate $d$ ?

## Review:

$$\frac{1}{2} \|\mathbf{x} - X\hat{\mathbf{z}}\|_2^2 \leq K(\hat{\mathbf{z}}, \lambda) \leq K(\hat{\mathbf{z}}_{eq}, \lambda) = \lambda \|\hat{\mathbf{z}}_{eq}\|_1$$

- $\|\hat{\mathbf{z}}_{eq}\|_1$  is on the order of  $\sqrt{d}$ .
- Why not use a multiple of  $\|\mathbf{z}\|_1$  as a proxy for  $\sqrt{d}$ ?

# How to Estimate $d$ ?

**A two-step procedure with data-driven regularization:**

- Solve:

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\mathbf{z}\| \quad \text{s.t. } \|\mathbf{x} - \mathbf{Xz}\|_2^2 < \tau$$

- Set  $\lambda = \alpha / \|\mathbf{z}^*\|_1$

- Solve

$$\hat{\mathbf{z}} = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x} - \mathbf{Xz}\|_2^2 + \lambda \|\mathbf{z}\|_1$$

# Experiments with Synthetic Data

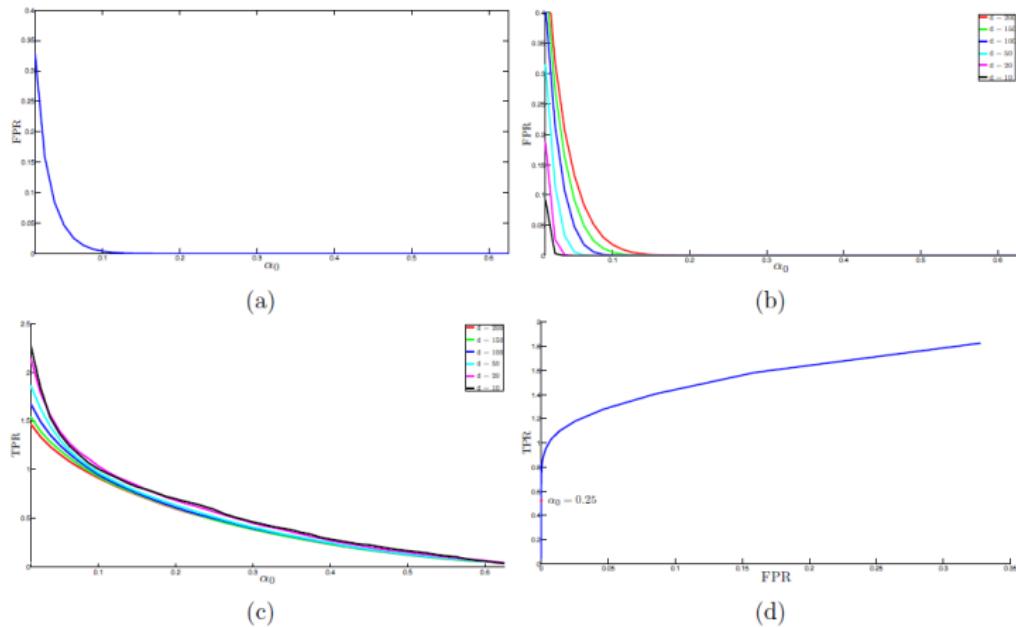


Figure: Performance of the two-step procedure using  $\alpha = 0.25$ . (a) False positive rate (FPR). (b) FPR for various subspace dimensions. (c) True positive rate (TPR). (d) TPR vs. FPR.

# Theoretical Guarantee

## Theorem 10 (No false discoveries)

*Assume that the subspace attached to the  $i$ -th column obeys the affinity and sampling conditions and that the noise level  $\sigma$  is bounded. Then in Robust Subspace Clustering Algorithm, take  $\tau = 2\sigma$  and  $\alpha > 0.707\sigma$ . Then with high probability, there is no false discovery in the  $i$ -th column.*

# Theoretical Guarantee

## Theorem 11 (Many true discoveries)

Consider the same setup as in above Theorem with  $\alpha$  also obeying  $\alpha < \alpha_0$  for some numerical constant  $\alpha_0$ . Then with high probability, there are at least

$$c_0 \frac{d(i)}{\log(N(i)/d(i))}$$

true discoveries in the  $i$ -th column ( $c_0$  is a positive numerical constant).

# Low-rank Representation

**Main Idea:** finding the lowest-rank representation of a collection of vectors jointly to capture the global structure.

# Low-rank Representation

**Main Idea:** finding the lowest-rank representation of a collection of vectors jointly to capture the global structure.

Consider:

$$\begin{aligned} & \min_Z \text{rank}(Z) \\ & \text{s.t. } X = XZ \end{aligned}$$

- NP-hard
- Ill-posed

# Low-rank Representation

Consider the nuclear norm:

$$\begin{aligned} & \min_Z \|Z\|_* \\ & \text{s.t. } X = XZ \end{aligned}$$

- $Z^* = X^+ X$  is the unique minimizer to problem.
- $Z^*$  is also a minimizer to the problem with the constrain of rank.

# Low-rank Representation

Consider the nuclear norm:

$$\begin{aligned} \min_Z & \|Z\|_* \\ \text{s.t. } & X = XZ \end{aligned}$$

- $Z^* = X^+ X$  is the unique minimizer to problem.
- $Z^*$  is also a minimizer to the problem with the constrain of rank.

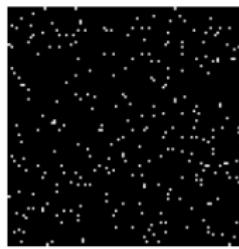
In other words, the power of nuclear norm is far more than convex relaxation!

# Low-rank Representation

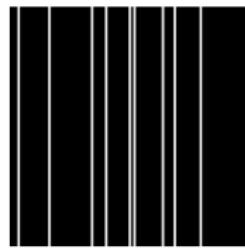
In this topic, we want to handle sample-specific corruptions.



(a) noise



(b) random corruptions



(c) sample-specific corruptions

Consider:

$$\min_{Z, E} \|Z\|_* + \lambda \|E\|_{2,1}$$

$$\text{s.t. } X = XZ + E$$

where  $\|E\|_{2,1} = \sum_{j=1}^n \sqrt{\sum_{i=1}^n (|E|_{ij})^2}$  is called  $l_{2,1}$ -norm.

# Solving the Optimization Problem

$$\min_{Z, E} \|Z\|_* + \lambda \|E\|_{2,1}$$

$$\text{s.t. } X = XZ + E$$

First convert it to the following problem:

$$\min_{Z, E, J} \|J\|_* + \lambda \|E\|_{2,1}$$

$$\text{s.t. } X = XZ + E, \ Z = J$$

# Solving the Optimization Problem

$$\min_{Z, E} \|Z\|_* + \lambda \|E\|_{2,1}$$

$$\text{s.t. } X = XZ + E$$

First convert it to the following problem:

$$\min_{Z, E, J} \|J\|_* + \lambda \|E\|_{2,1}$$

$$\text{s.t. } X = XZ + E, \quad Z = J$$

Solving the Augmented Lagrange Multiplier (ALM) problem:

$$\min_{Z, E, J, Y_1, Y_2} \|J\|_* + \lambda \|E\|_{2,1} + \text{tr} [Y_1^T (X - XZ - E)] + \text{tr} [Y_2^T (Z - J)] + \frac{\mu}{2} (\|X - XZ - E\|_F^2 + \|Z - J\|_F^2)$$

# Solving the Optimization Problem

---

**Algorithm 1** Solving Problem by Inexact ALM

---

Input: data matrix  $X$ , parameter  $\lambda$

Initialize:  $Z = J = 0, E = 0, Y_1 = 0, Y_2 = 0, \mu = 10^{-6}, max_u = 10^{10}, \rho = 1.1, \varepsilon = 10^{-8}$ .

while not converged do

1. fix the others and update  $J$  by

$$J = \arg \min \frac{1}{\mu} \|J\|_* + \frac{1}{2} \|J - (Z + Y_2/\mu)\|_F^2$$

2. fix the others and update  $Z$  by

$$Z = (I + X^t X)^{-1} (X^t X - X^t E + J + (X^t Y_1 - Y_2)/\mu)$$

3. fix the others and update  $E$  by

$$E = \arg \min \frac{\lambda}{\mu} \|E\|_{2,1} + \frac{1}{2} \|E - (X - XZ + Y_1/\mu)\|_F^2$$

4. update the multipliers

$$Y_1 = Y_1 + \mu(X - XZ - E)$$

$$Y_2 = Y_2 + \mu(Z - J)$$

5. update the parameter  $\mu$  by  $\mu = \min(\rho\mu, max_u)$

6. check the convergence conditions

$$\|X - XZ - E\|_\infty < \varepsilon \text{ and } \|Z - J\|_\infty < \varepsilon.$$

end while

# The Framework of Low-rank Representation

- 1) Solve the low-rank optimization program:

$$\begin{aligned} & \min_{Z, E} \|Z\|_* + \lambda \|E\|_{2,1} \\ & \text{s.t. } X = XZ + E \end{aligned}$$

- 2) Normalize the columns of  $Z$  as  $\mathbf{z}_i \leftarrow \frac{\mathbf{z}_i}{\|\mathbf{z}_i\|_\infty}$
- 3) Form a similarity graph with  $N$  nodes representing the data points. Set the weights on the edges between the nodes by  $W = |Z| + |Z|^T$ .
- 4) Apply spectral clustering to the similarity graph.

# Least Squares Subspace Clustering

Review a theorem in sparse subspace clustering:

## Theorem 12

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  $X_i$  denote  $N_i$  data points in  $\mathcal{S}_i$ , where  $\text{rank}(X_i) = d_i$ , and let  $X_{-i}$  denote data points in all subspaces except  $\mathcal{S}_i$ . Then, for every  $\mathcal{S}_i$  and every nonzero  $\mathbf{x}$  in  $\mathcal{S}_i$ , the  $l_q$ -minimization program:

$$\begin{bmatrix} \mathbf{z}^* \\ \mathbf{z}_{-}^* \end{bmatrix} = \operatorname{argmin} \left\| \begin{bmatrix} \mathbf{z} \\ \mathbf{z}_{-} \end{bmatrix} \right\|_q \quad \text{s. t.} \quad \mathbf{x} = [X_i \ X_{-i}] \begin{bmatrix} \mathbf{z} \\ \mathbf{z}_{-} \end{bmatrix}$$

for  $q \geq 1$ , recovers a subspace-sparse representation, i.e.,  $\mathbf{z}^* \neq \mathbf{0}$  and  $\mathbf{z}_{-}^* = \mathbf{0}$ .

# Least Squares Subspace Clustering

- When the subspaces are independent, for  $q \geq 1$ , the  $l_q$  minimization can recover the representation.
- $l_2$ -norm can prevent over sparsity in the same class.
- $l_2$ -norm is more tractable than  $l_1$  norm.

Consider:

$$\min_Z \|X - XZ\|_F^2 + \lambda \|Z\|_F^2 \quad \text{s.t. } \text{diag}(Z) = 0$$

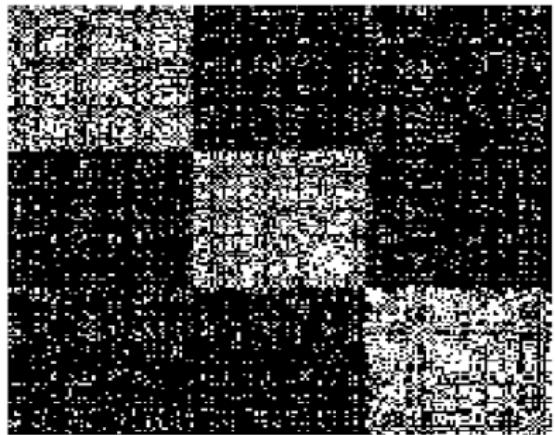
- The above optimization problem can be solved closely!

# Least Squares Subspace Clustering

SSC



LSR



- ✓ Few wrong connections
- ✗ Not well connected

- ✗ Many wrong connections
- ✓ Well-connected

# Summary

- **Robust Subspace Clustering**

Choose  $\lambda$  wisely to guarantee the prevent false discoveries and make sure that the number of true discoveries at the same time.

- **Low-rank Representation**

Lowest-rank representation to capture the global structure.

- **Least Squares Subspace Clustering**

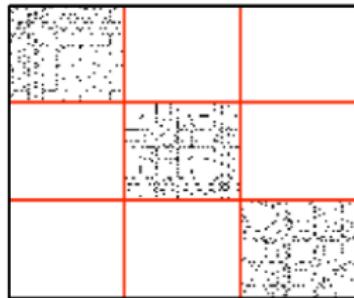
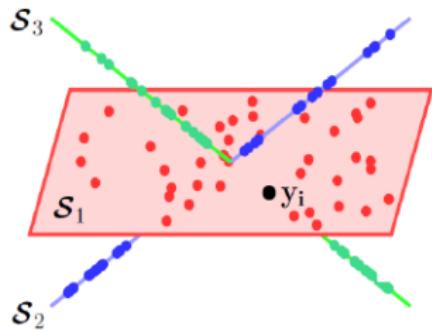
Use  $l_2$ -norm to replace  $l_1$ -norm and prevent over sparsity in the same class.

# Outline

- 1 Background
- 2 Sparse Subspace Clustering
- 3 Other Self-Expressive Models
- 4 Subspace Clustering by Block Diagonal Representation
- 5 More than Linear Model

# Block Diagonality Properties

If  $X$  is arranged, then a good representation matrix  $Z$  should be block diagonal.



# Unifying Previous Methods

Previous methods can be written as:

$$\min f(Z) \text{ s.t. } Z \in \Omega$$

	$f(Z)$	$\Omega$
SSC	$\ Z\ _0$ or $\ Z\ _1$	$\{Z   \textcolor{red}{X} = XZ, \text{diag}(Z) = 0\}$
LRR	$\ Z\ _*$	$\{Z   \textcolor{red}{X} = XZ\}$
MSR	$\ Z\ _1 + \lambda \ Z\ _*$	$\{Z   \textcolor{red}{X} = XZ, \text{diag}(Z) = 0\}$
SSQP	$\ Z^\top Z\ _1$	$\{Z   \textcolor{red}{X} = XZ, Z \geq 0, \text{diag}(Z) = 0\}$
LSR	$\ Z\ _F$	$\{Z   \textcolor{red}{X} = XZ\}$

- All the above methods give **block diagonal** solution for **independent** subspaces.
- However, they are proved case by case.
- More details about MSR and SSQP can be seen in [Luo et al., 2011] and [Wang et al., 2011] respectively.

# Unifying Previous Methods

## What $f$ Gives Block Diagonality?

- Consider the following general SC formulation:

$$\min \textcolor{blue}{f}(Z) \text{ s.t. } Z \in \Omega$$

- When subspaces are independent, What kind of objective functions induce the **block diagonal** solution?
- Enforced Block Diagonal Condition [Lu et al., 2018] will be introduced in the following parts.

# Enforced Block Diagonal Condition

**Enforced Block Diagonal (EBD) Conditions.** Assume the matrix function  $f$  is defined on  $\Omega(\neq \emptyset)$ . For any  $Z = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \in \Omega$ ,  $Z \neq 0$ , where  $A$  and  $D$  are square matrices,  $B$  and  $C$  are of compatible dimension,  $A, D \in \Omega$ . Let  $Z^D = \begin{bmatrix} A & 0 \\ 0 & D \end{bmatrix} \in \Omega$ . The EBD conditions are:

- (1)  $f(Z) = f(P^\top Z P)$ ,  $\forall$  permutation matrix  $P \in \{P | P^\top Z P \in \Omega\}$ .
- (2)  $f(Z) \geq f(Z^D)$ , where the equality holds iff  $B = C = 0$ .
- (3)  $f(Z^D) = f(A) + f(D)$ .

# Enforced Block Diagonal Condition

$f$  satisfies EBD cond.

+

→ block diagonal  $Z$

subspaces are independent

**Theorem** Assume data sampling is sufficient, and subspaces are independent. Consider:

$$\min f(Z) \text{ s.t. } Z \in \Omega = \{Z | X = XZ\}.$$

If  $f$  satisfies the EBD conditions (1)(2), then  $Z^*$  is block diagonal

$$Z^* = \text{blkdiag}(Z_1^*, Z_2^*, \dots, Z_k^*)$$

with  $Z_i^* \in \mathbb{R}^{n_i \times n_i}$  corresponding to  $X_i$ , for each  $i$ . Furthermore, if  $f$  satisfies the EBD conditions (1)(2)(3), for each  $i$ ,  $Z_i^*$  is optimal solution to:

$$\min f(Y) \text{ s.t. } X_i = X_i Y$$

# Enforced Block Diagonal Condition

- The above theorem gives a general guarantee for almost all subspace clustering algorithms when subspaces are independent.
- Usually, the solution is far from being  $k$ -block diagonal since the independent subspaces assumption does not hold due to data noise.
- How to Pursue Block-diagonality directly?
  - How to pursue exactly block-diagonal structure?
  - Efficient optimization?
  - Performance guarantee?

# Block Diagonal Prior

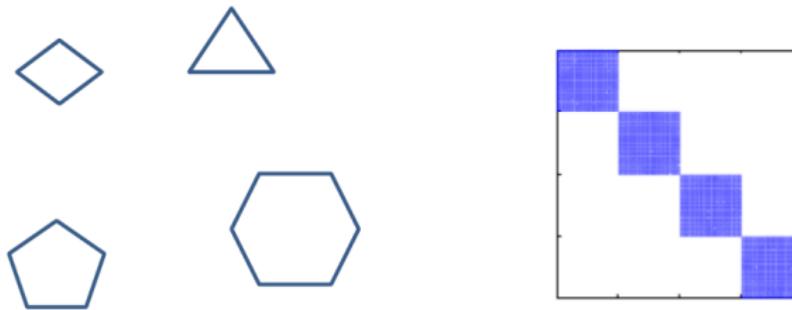
**Aim:** Give a formal definition of the  $k$ -block diagonal matrix.  
Consider the following matrix:

$$B = \begin{bmatrix} B_0 & 0 & 0 \\ 0 & B_0 & 0 \\ 0 & 0 & B_0 \end{bmatrix}, \text{ where } B_0 = \begin{bmatrix} 1 & 0 \\ -1 & 1 \end{bmatrix}$$

- There are 3 blocks intuitively.
- But we can also say that it has 1 or 2 blocks.
- How to define  $k$ -block diagonal property?

# Block Diagonal Prior

Recall a spectral graph theory:



**Theorem** Let  $W$  be an affinity matrix. Then the multiplicity  $k$  of the eigenvalue 0 of the corresponding Laplacian  $L_W$  equals the number of connected components in  $W$ .

Here,  $L_W(j, j') = -W(j, j')$ , if  $j \neq j'$ ;  $\sum_{\ell \neq j} W(j, \ell)$  otherwise.

$W$  is  $k$  block diagonal  $\Leftrightarrow W \in \{W \mid \text{rank}(L_W) = n - k, W \in \mathbb{R}^{n \times n}\}$

# $k$ -Block Diagonal Regularizer

Suppose data lying in  $k$  subspaces:

- **Aim:** design a regularizer to enforce  $B$  to be  $k$ -block diagonal
- $L_B$  has  $k$  zero eigenvalues
- Consider the sum of the smallest  $k$  eigenvalues of  $L_Z$

## Definition 13

For any affinity matrix  $B \in \mathbb{R}^{n \times n}$ , the  $k$ -block diagonal regularizer is defined as the sum of the  $k$  smallest eigenvalues of  $L_B$ , i.e.,

$$\|B\|_{\mathcal{K}} = \sum_{i=n-k+1}^n \lambda_i(L_B)$$

# Block Diagonal Representation

**$k$ -block diagonal regularizer:**

$$\|B\|_{\mathcal{K}} = \sum_{i=n-k+1}^n \lambda_i(L_B)$$

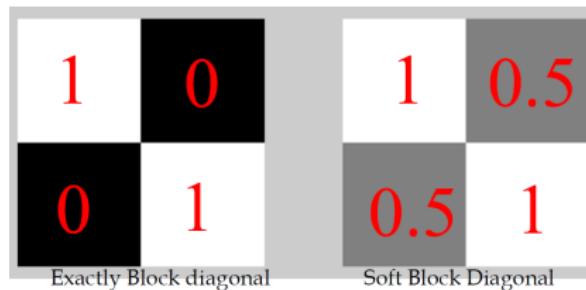
With the proposed  $k$ -block diagonal regularizer, the Block Diagonal Representation (BDR) method for subspace clustering is as follows [Lu et al., 2018]:

$$\begin{aligned} & \min_B \frac{1}{2} \|X - XB\|_F^2 + \gamma \|B\|_{\mathcal{K}} \\ \text{s.t. } & \text{diag}(B) = 0, B \geq 0, B = B^T \end{aligned}$$

# Block Diagonal Representation

$$\begin{aligned} & \min_B \frac{1}{2} \|X - XB\|_F^2 + \gamma \|B\|_{\mathcal{K}} \\ \text{s.t. } & \text{diag}(B) = 0, B \geq 0, B = B^T \end{aligned}$$

- The restrictions on  $B$  will limit its representation capability.
- An exactly block diagonal matrix may not be necessary.



- Both the above two affinity matrices lead to the same clustering results.

# Soft Block Diagonal Representation

$$\min_B \frac{1}{2} \|X - XB\|_F^2 + \gamma \|B\|_{\mathcal{K}}$$

$$\text{s.t. } \text{diag}(B) = 0, B \geq 0, B = B^T$$

Consider the soft version of block diagonal representation:

$$\min_{B,Z} \frac{1}{2} \|X - XZ\|^2 + \frac{\lambda}{2} \|Z - B\|^2 + \gamma \|B\|_{\mathcal{K}}$$

$$\text{s.t. } \text{diag}(B) = 0, B \geq 0, B = B^T$$

# Soft Block Diagonal Representation

$$\begin{aligned} & \min_B \frac{1}{2} \|X - XB\|_F^2 + \gamma \|B\|_{\mathcal{K}} \\ \text{s.t. } & \text{diag}(B) = 0, B \geq 0, B = B^T \end{aligned}$$

Consider the soft version of block diagonal representation:

$$\begin{aligned} & \min_{B,Z} \frac{1}{2} \|X - XZ\|^2 + \frac{\lambda}{2} \|Z - B\|^2 + \gamma \|B\|_{\mathcal{K}} \\ \text{s.t. } & \text{diag}(B) = 0, B \geq 0, B = B^T \end{aligned}$$

**Next problem:** How to optimize?

The key challenge lies in the nonconvex term  $\|B\|_{\mathcal{K}}$ .

# Optimization of BDR

## Theorem 14

Let  $L \in \mathbb{R}^{n \times n}$  and  $L \succeq 0$ . Then

$$\sum_{i=n-k+1}^n \lambda_i(L) = \min_W \langle L, W \rangle, \text{ s.t. } 0 \preceq W \preceq I, \text{tr}(W) = k$$

# Optimization of BDR

## Theorem 14

Let  $L \in \mathbb{R}^{n \times n}$  and  $L \succeq 0$ . Then

$$\sum_{i=n-k+1}^n \lambda_i(L) = \min_W \langle L, W \rangle, \text{ s.t. } 0 \preceq W \preceq I, \text{tr}(W) = k$$

Then the origin problem is equivalent to:

$$\begin{aligned} & \min_{Z, B, W} \frac{1}{2} \|X - XZ\|_F^2 + \frac{\lambda}{2} \|Z - B\|_F^2 + \gamma \langle \text{diag}(B\mathbf{1}) - B, W \rangle \\ & \text{s.t. } \text{diag}(B) = 0, B \geq 0, B = B^T, 0 \preceq W \preceq I, \text{tr}(W) = k \end{aligned}$$

# Optimization of BDR

$$\begin{aligned} & \min_{Z, B, W} \frac{1}{2} \|X - XZ\|_F^2 + \frac{\lambda}{2} \|Z - B\|_F^2 + \gamma \langle \text{diag}(B\mathbf{1}) - B, W \rangle \\ \text{s.t. } & \text{diag}(B) = 0, B \geq 0, B = B^T, 0 \preceq W \preceq I, \text{tr}(W) = k \end{aligned}$$

# Optimization of BDR

$$\begin{aligned} & \min_{Z, B, W} \frac{1}{2} \|X - XZ\|_F^2 + \frac{\lambda}{2} \|Z - B\|_F^2 + \gamma \langle \text{diag}(B\mathbf{1}) - B, W \rangle \\ \text{s.t. } & \text{diag}(B) = 0, B \geq 0, B = B^T, 0 \preceq W \preceq I, \text{tr}(W) = k \end{aligned}$$

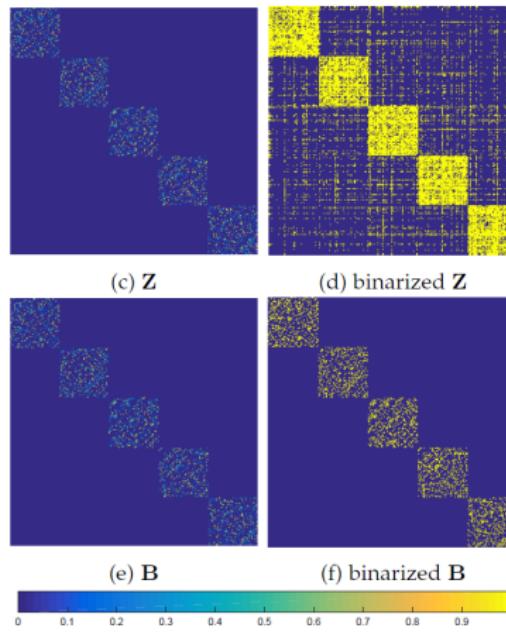
Note that  $W$  is independent from  $Z$ , thus one can group them as a super block  $\{W, Z\}$  and treat  $\{B\}$  as the other block.

Thus, the problem can be solved by alternating updating  $\{W, Z\}$  and  $\{B\}$ .

**All of sub problems can be solved efficiently!**

# Experiments with Synthetic Data

- 5 disjoint subspaces in  $\mathbb{R}^{30}$
- Each subspace of dimension 5
- Pick 50 data points in each subspace



# Convergence Guarantee

$$\begin{aligned} & \min_{Z, B, W} \frac{1}{2} \|X - XZ\|_F^2 + \frac{\lambda}{2} \|Z - B\|_F^2 + \gamma \langle \text{diag}(B\mathbf{1}) - B, W \rangle \\ & \text{s.t. } \text{diag}(B) = 0, B \geq 0, B = B^T, 0 \preceq W \preceq I, \text{tr}(W) = k \end{aligned}$$

## Theorem 15

*The sequence  $\{W^k, Z^k, B^k\}$  generated by above algorithm has at least one limit point and any limit point  $(Z^*, B^*, W^*)$  of  $\{Z^k, B^k, W^k\}$  is a stationary point.*

# Experiment on Motion Segmentation

## Hopkins 155 database:



method	SSC	LRR	LSR	BDR-B	BDR-Z
2-motions					
mean	1.52	3.65	3.24	1.00	<b>0.95</b>
3-motions					
mean	4.40	9.40	5.94	1.95	<b>0.85</b>
all					
mean	2.18	4.95	3.85	1.22	<b>0.93</b>

Table: The mean clustering errors (%) of 155 sequences on Hopkins 155 dataset.

# Experiment on Face Clustering

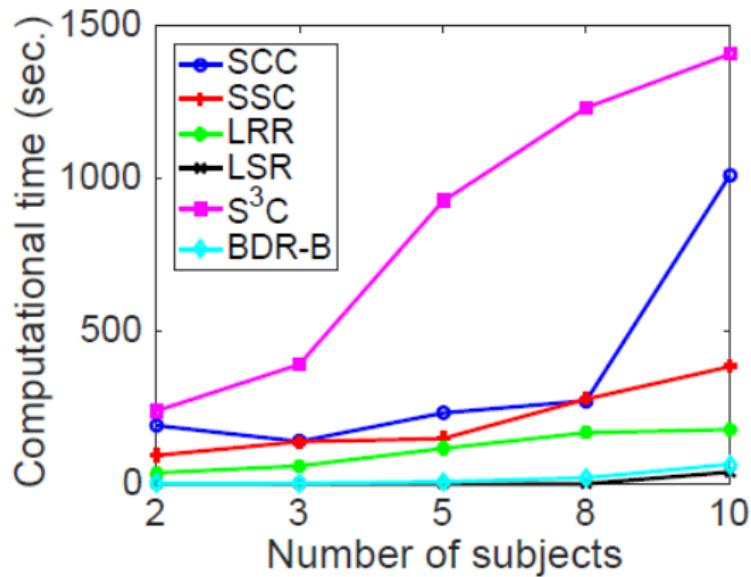
Extended Yale B database:



method	2 subjects			3 subjects			5 subjects			8 subjects			10 subjects		
	mean	median	std	mean	median	std	mean	median	std	mean	median	std	mean	median	std
SCC	24.02	19.92	17.82	42.19	41.93	8.93	61.36	62.34	6.10	71.87	72.27	4.72	72.48	73.28	6.14
SSC	1.64	0.78	2.91	3.26	0.52	7.69	6.30	4.22	5.43	8.94	9.67	6.18	10.09	11.33	4.59
LRR	5.39	0.39	14.50	6.04	1.04	12.34	8.13	2.34	9.61	6.79	3.42	6.50	9.49	12.58	5.38
LSR	3.16	0.78	10.18	3.96	1.56	8.72	7.85	6.72	8.72	28.14	31.05	12.32	33.27	33.12	4.57
S <sup>3</sup> C	1.29	0.00	2.69	2.79	0.52	7.38	4.66	1.88	5.15	6.37	6.35	5.32	6.87	6.17	3.67
BDR-B	3.28	0.78	10.15	3.02	1.30	7.78	4.45	2.19	6.29	3.08	2.93	1.18	2.95	2.81	1.09
BDR-Z	2.97	0.00	10.23	1.15	1.04	0.95	3.00	2.66	2.25	4.46	4.20	2.39	4.04	3.52	1.52

Figure: Clustering error (%) of different algorithms on the Extended Yale B database

# Experiment on Face Clustering



**Figure:** Average computational time (sec.) of the algorithms on the Extended Yale B database as a function of the number of subjects.

# Summary

- Enforced Block Diagonal Condition provided a uniform theoretical guarantee for almost all subspace clustering algorithms when subspaces are independent.
- $K$ -block diagonal regularizer is a powerful tool which encourages a nonnegative symmetric matrix to be  $k$ -block diagonal.
- Block Diagonal Representation (BDR) utilizes the block diagonal regularizer method for subspace clustering.

# Outline

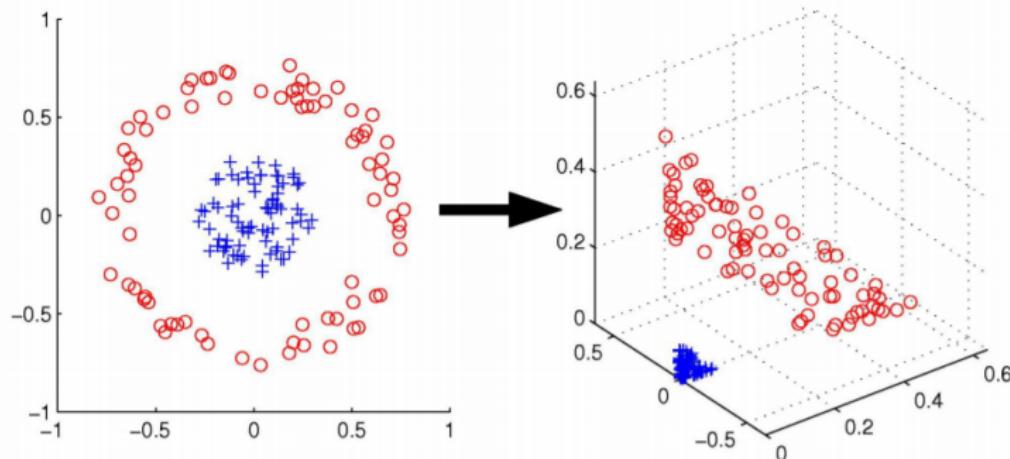
- 1 Background
- 2 Sparse Subspace Clustering
- 3 Other Self-Expressive Models
- 4 Subspace Clustering by Block Diagonal Representation
- 5 More than Linear Model

# More than Linear Model

- While the traditional subspace clustering methods achieve state-of-the-art results on several benchmarks, in practice many datasets are better modeled by non-linear manifolds.
- Several methods extend traditional methods to non-linear manifolds.
  - Kernel Sparse Subspace Clustering [Patel and Vidal, 2014]
  - Deep Subspace Clustering Networks [Ji et al., 2017]

# Kernel Sparse Subspace Clustering

**Main Idea:** Using kernel trick



**Figure:** Non-linear subspaces (or sub-manifolds) are mapped to linear ones in high-dimensional feature space.

# Kernel Sparse Subspace Clustering

Let  $\phi : \mathbb{R}^D \rightarrow \mathcal{H}$  be a mapping from the input space to the reproducing kernel Hilbert space  $\mathcal{H}$ .

Consider the following optimization problem:

$$\begin{aligned} & \min_{C, Z} \|Z\|_1 + \lambda \|\phi(X) - \phi(X)C\|_F^2 \\ & \text{s.t. } C = Z - \text{diag}(Z), \quad C^T \mathbf{1} = \mathbf{1} \end{aligned}$$

This problem can be efficiently solved using the ADMM method.

# Experiment on Alphadigits Dataset



Algorithms (2 digits)	LRSC	LSA	SSC	LRR	KSSC P(3,2)	KSSC P(2,0.2)	KSSC G(10)	KSSC G(7)
Mean	15.81	10.70	5.70	7.76	5.42	5.40	6.13	5.93
Median	8.97	3.85	2.56	3.84	2.56	2.56	2.56	2.56
Algorithms (3 digits)	LRSC	LSA	SSC	LRR	KSSC P(3,2)	KSSC P(2,0.2)	KSSC G(10)	KSSC G(7)
Mean	25.65	22.69	13.58	14.21	12.85	13.30	14.54	13.64
Median	25.64	22.22	8.54	11.11	7.69	8.54	9.40	8.55
Algorithms (5 digits)	LRSC	LSA	SSC	LRR	KSSC P(3,2)	KSSC P(2,0.2)	KSSC G(10)	KSSC G(7)
Mean	37.77	33.81	23.26	23.34	22.64	23.00	24.50	23.84
Median	37.95	33.85	25.12	23.59	23.08	23.85	27.18	26.15

# Deep Subspace Clustering Networks

- Despite kernel trick, the selection of different kernel types is largely empirical, and there is no clear reason to believe that the implicit feature space corresponding to a predefined kernel is truly well-suited to subspace clustering.
- **Main Idea:** Use a novel deep neural network architecture to learn an explicit non-linear mapping of the data that is well-adapted to subspace clustering.

# Deep Subspace Clustering Networks

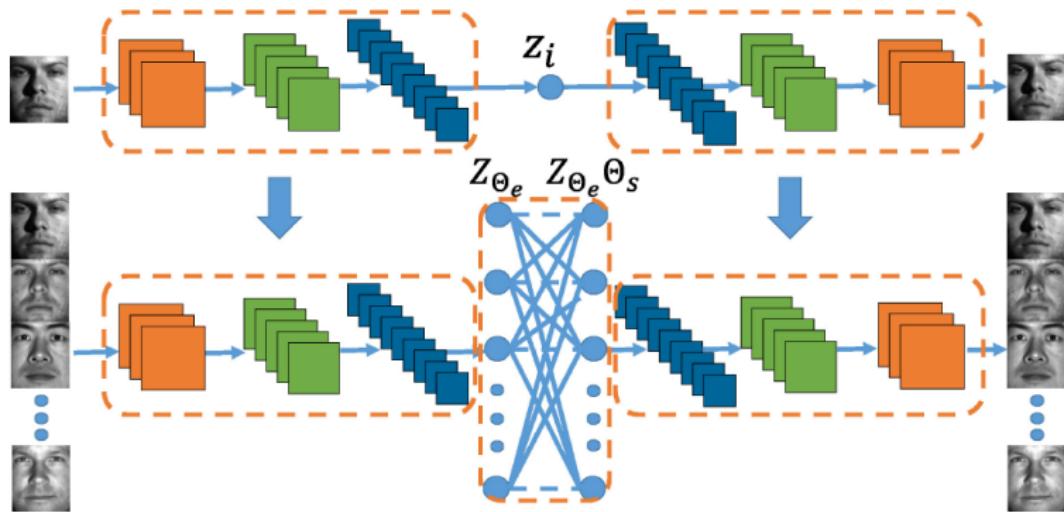


Figure: Deep Subspace Clustering Networks: As an example, we show a deep subspace clustering network with three convolutional encoder layers, one self-expressive layer, and three deconvolutional decoder layers. During training, we first pre-train the deep auto-encoder without the self-expressive layer; we then fine-tune our entire network using this pre-trained model for initialization.

# Deep Subspace Clustering Networks

More details:

- Structure of the networks:
  - Convolution encoder
  - Self-expressive layer
  - Deconvolution decoder

# Deep Subspace Clustering Networks

More details:

- Structure of the networks:

- Convolution encoder
- Self-expressive layer
- Deconvolution decoder

- Loss Function:

$$L(\Theta, C) = \frac{1}{2} \left\| X - \hat{X}_{\Theta} \right\|_F^2 + \lambda_1 \|C\|_p + \frac{\lambda_2}{2} \|Z_{\Theta_e} - Z_{\Theta_e} C\|_F^2$$

s.t.  $\text{diag}(C) = 0$

# Deep Subspace Clustering Networks

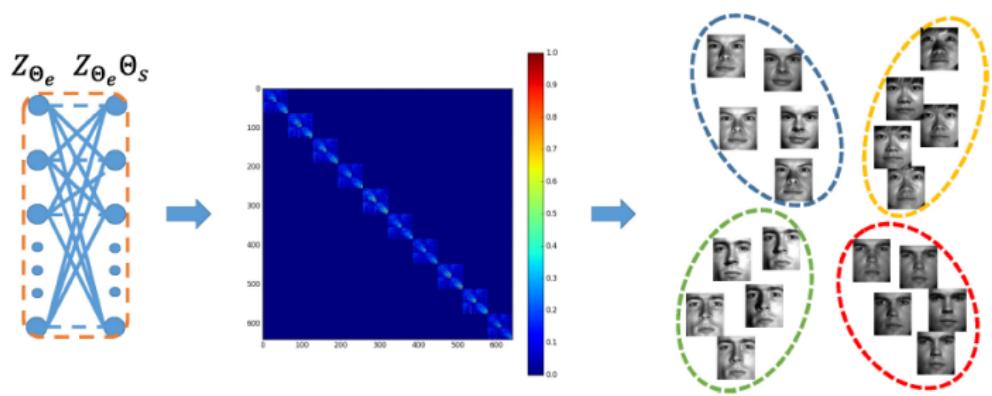


Figure: From the parameters of the self-expressive layer, we construct an affinity matrix, which we use to perform spectral clustering to get the final clusters. Best viewed in color.

# Experiments on COIL20 and COIL100 Datasets

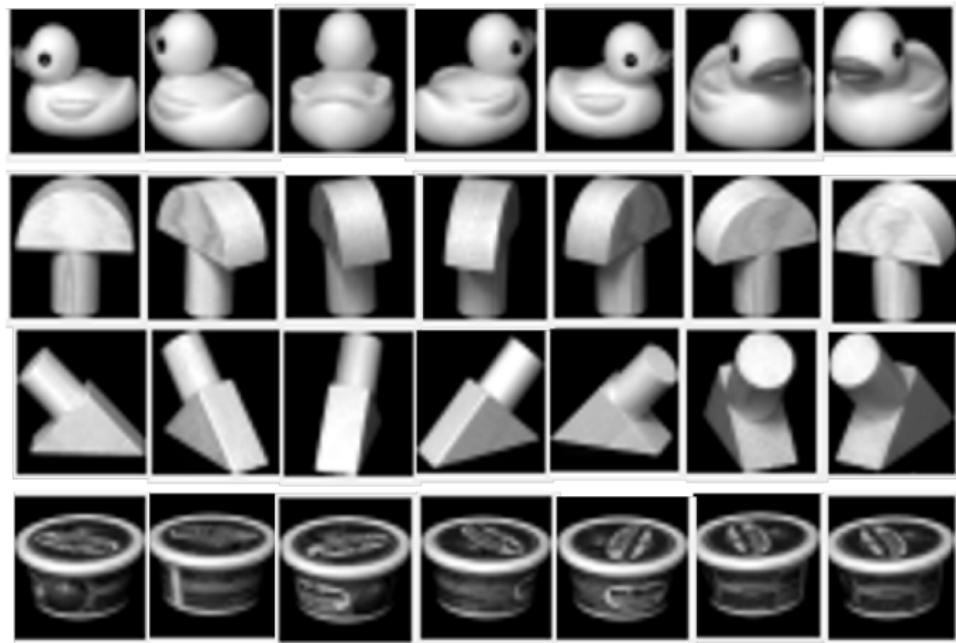


Figure: COIL20 and COIL100

# Experiments on COIL20 and COIL100 Datasets

layers	COIL20			COIL100		
	encoder-1	self-expressive	decoder-1	encoder-1	self-expressive	decoder-1
kernel size	$3 \times 3$	–	$3 \times 3$	$5 \times 5$	–	$5 \times 5$
channels	15	–	15	50	–	50
parameters	150	2073600	136	1300	51840000	1251

Figure: Network settings for COIL20 and COIL100

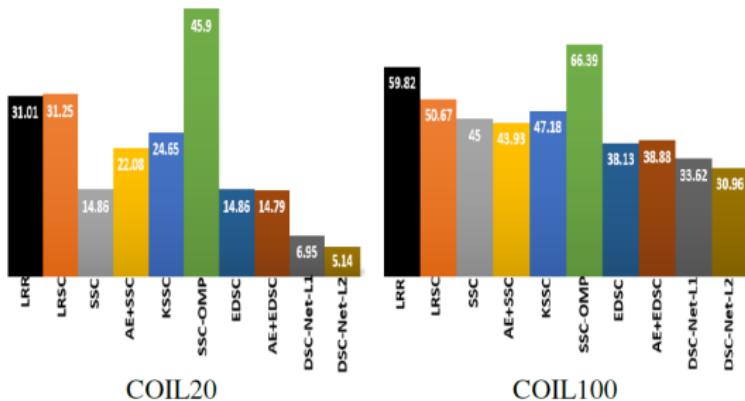


Figure: Subspace clustering error on COIL20 and COIL100 datasets

# Summary

- Traditional subspace clustering methods are restricted when data lying on non-linear manifolds.
- We want to handle manifolds data:
  - Kernel Method
  - Neural Network (Self-Expressive Layer)

# References I

-  Basri, R. and Jacobs, D. W. (2003).  
 Lambertian reflectance and linear subspaces.  
*IEEE Transactions on Pattern Analysis & Machine Intelligence*, (2):218–233.
-  Costeira, J. and Kanade, T. (1995).  
 A multi-body factorization method for motion analysis.  
 In *Proceedings of IEEE International Conference on Computer Vision*, pages 1071–1076. IEEE.
-  Cover, T. M., Hart, P., et al. (1967).  
 Nearest neighbor pattern classification.  
*IEEE transactions on information theory*, 13(1):21–27.
-  Elhamifar, E. and Vidal, R. (2009).  
 Sparse subspace clustering.  
 In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2790–2797. IEEE.
-  Ji, P., Zhang, T., Li, H., Salzmann, M., and Reid, I. (2017).  
 Deep subspace clustering networks.  
 In *Advances in Neural Information Processing Systems*, pages 24–33.
-  Kanatani, K.-i. (2001).  
 Motion segmentation by subspace separation and model selection.  
 In *Proceedings Eighth IEEE International Conference on computer Vision. ICCV 2001*, volume 2, pages 586–591. IEEE.

# References II

-  Liu, G., Lin, Z., and Yu, Y. (2010).  
Robust subspace segmentation by low-rank representation.  
In *ICML*, volume 1, page 8.
-  Lu, C., Feng, J., Lin, Z., Mei, T., and Yan, S. (2018).  
Subspace clustering by block diagonal representation.  
*IEEE transactions on pattern analysis and machine intelligence*, 41(2):487–501.
-  Lu, C.-Y., Min, H., Zhao, Z.-Q., Zhu, L., Huang, D.-S., and Yan, S. (2012).  
Robust and efficient subspace segmentation via least squares regression.  
In *European conference on computer vision*, pages 347–360. Springer.
-  Luo, D., Nie, F., Ding, C., and Huang, H. (2011).  
Multi-subspace representation and discovery.  
In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 405–420. Springer.
-  Patel, V. M. and Vidal, R. (2014).  
Kernel sparse subspace clustering.  
In *2014 IEEE International Conference on Image Processing (ICIP)*, pages 2849–2853. IEEE.
-  Soltanolkotabi, M., Elhamifar, E., Candes, E. J., et al. (2014).  
Robust subspace clustering.  
*The Annals of Statistics*, 42(2):669–699.

# References III

-  Tipping, M. E. and Bishop, C. M. (1999).  
Mixtures of probabilistic principal component analyzers.  
*Neural computation*, 11(2):443–482.
-  Tseng, P. (2000).  
Nearest q-flat to m points.  
*Journal of Optimization Theory and Applications*, 105(1):249–252.
-  Von Luxburg, U. (2007).  
A tutorial on spectral clustering.  
*Statistics and computing*, 17(4):395–416.
-  Wang, S., Yuan, X., Yao, T., Yan, S., and Shen, J. (2011).  
Efficient subspace segmentation via quadratic programming.  
In *Twenty-Fifth AAAI Conference on Artificial Intelligence*.