

Sparse Subspace Clustering (SSC)

Javier A. Salazar Cavazos

University of Michigan

June 15, 2022

1 Background

Motivation, Subspaces, & Clustering

Subspace Clustering, Challenges, & Literature Review

2 Sparse Subspace Clustering

Affinity Matrix Construction

Spectral Clustering

Results

3 References & Questions

Motivation



Face Clustering [1]

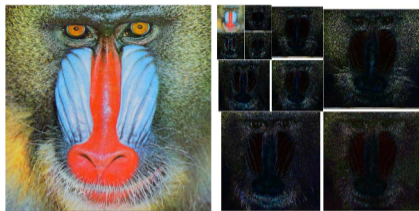


Image Compression [2]

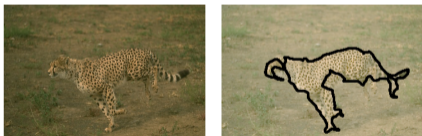
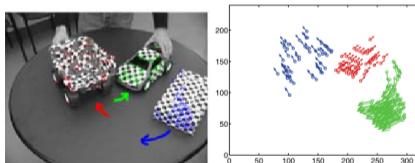


Image Segmentation [3]



Motion Estimation [4]

Image Credit: Respective Papers

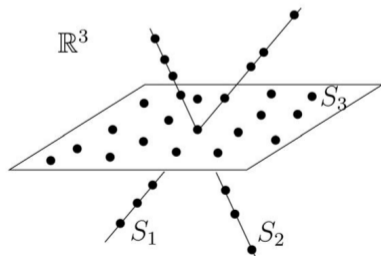
Subspaces

For a vector space \mathcal{V} defined on a field \mathbb{F} , a nonempty set $\mathcal{S} \subseteq \mathcal{V}$ is called a *linear* subspace iff

- \mathcal{S} is closed under vector addition (i.e. for $\mathbf{u}, \mathbf{v} \in \mathcal{S} \implies \mathbf{u} + \mathbf{v} \in \mathcal{S}$)
- \mathcal{S} is closed under scalar multiplication (i.e. $\mathbf{v} \in \mathcal{S} \implies \alpha \mathbf{v} \in \mathcal{S}$)

Key Points:

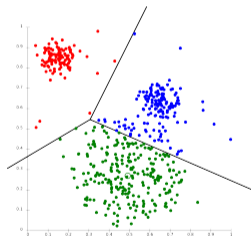
- Every *linear* subspace includes $\mathbf{0}$
- A shifted *linear* subspace is called an affine subspace
- Left singular matrix contains subspace basis



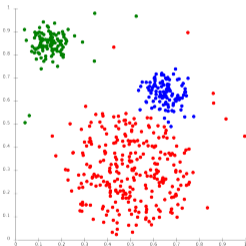
Clustering

An unsupervised machine learning method to identify and group “similar” unlabeled data points in order to find structure or patterns

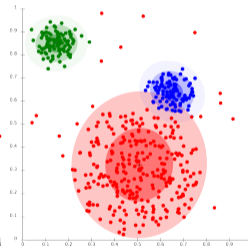
By “similar”, this can mean different things:



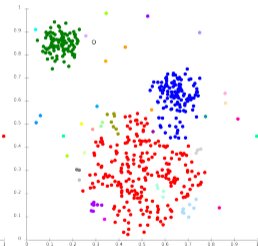
Centroid-based



Density-based

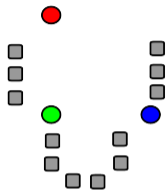


Model-based

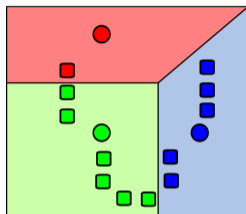


Hierarchical-based

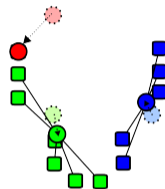
K-Means Clustering



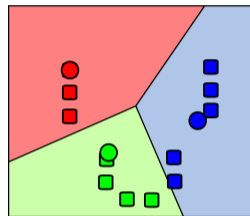
1. Initialize centroids



2. Generate clusters



3. Update centroids



4. Repeat steps 2 & 3

Note!

NP-hard problem \implies initialization is extremely important!

Subspaces + Clustering

$n = \#$ of points, $k = \#$ of subspaces, $D = \dim(\mathbf{y})$, $d_i = \dim(\mathcal{S}_i)$

Points:

$Y = [\mathbf{y}_1, \dots, \mathbf{y}_n]$ for $\{\mathbf{y}_j \in \mathbb{R}^D\}_{j=1}^n$ drawn from $\{\mathcal{S}_i\}_{i=1}^k$

Subspaces:

$\mathcal{S}_i = \{\mathbf{y} \in \mathbb{R}^D : \mathbf{y} = U_i \mathbf{z}\}$ where \mathbf{z} is coordinate in basis U_i

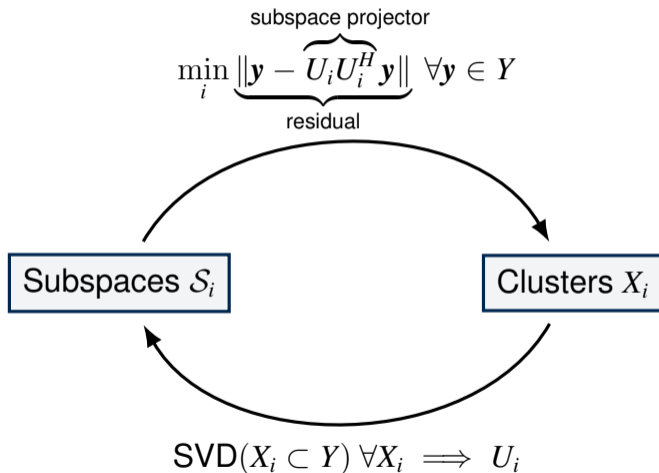
Goal:

Find subspace bases $U_i \in \mathbb{R}^{D \times d_i}$ and segmentation $X_i \subset Y, X_i \in \mathcal{S}_i$

Recall!

The SVD can only return one subspace basis ($k = 1$ case) for Y
How do we solve this when $k > 1$?

The Chicken & Egg Dilemma



Literature Review [5]

Algebraic Methods (e.g. Generalized PCA)

- Geometric and algebraic interpretations

Iterative Methods (e.g. K-Subspaces)

- Alternate between segmentation and subspaces

Statistical Methods (e.g. Mixture of Probabilistic PCA)

- Model assumptions combined with prior beliefs

Affinity-based Methods (e.g. Sparse Subspace Clustering)

- Measures “similarity” between points

Outline

1 Background

Motivation, Subspaces, & Clustering

Subspace Clustering, Challenges, & Literature Review

2 Sparse Subspace Clustering

Affinity Matrix Construction

Spectral Clustering

Results

3 References & Questions

SSC Affinity Measure

Self-Expressiveness Property

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

$$\text{For } y_i \in \bigcup_{l=1}^n \mathcal{S}_l, \text{ self-expressive: } y_i = \underbrace{Yc_i}_{\text{not unique!}} \text{ where } \underbrace{c_{ii}}_{y_i \neq y_i \text{ (trivial)}} = 0$$

self-expressiveness \gg other affinity measures

Recall that $u, v \in \mathcal{S} \implies u + v \in \mathcal{S}$

self-expressiveness = natural result of subspace definition

Question: Why is c_i not unique in general?

Sparse Subspace Clustering [6]

Nonconvex formulation (vector case):

$$\min_{c_i} \|c_i\|_0 \quad \text{s.t.} \quad y_i = Yc_i, c_{ii} = 0$$

Convex formulation (matrix case):

$$\min_C \underbrace{\|C\|_{1,1}}_{\|C\|_{1,1} \triangleq \|\text{vec}(C)\|_1} \quad \text{s.t.} \quad Y = YC, \text{diag}(C) = \mathbf{0}$$

Convex formulation (matrix case + outliers + noise):

$$\min_{C,E,Z} \underbrace{\|C\|_{1,1}}_{\text{sparse affinity}} + \underbrace{\lambda_e \|E\|_{1,1}}_{\text{sparse outliers}} + \frac{\lambda_z}{2} \underbrace{\|Z\|_F^2}_{\text{noise}} \quad \text{s.t.} \quad Y = YC + E + Z, \overbrace{\text{diag}(C) = \mathbf{0}}^{C \neq I} \quad (1)$$

Optimization - ADMM

Introduce auxiliary variable $A = C - \text{mdiag}(C)$ where $\text{mdiag}(C) = \text{diagm}(\text{diag}(C))$ to decouple terms such that:

$$\min_{C,E,A} \|C\|_{1,1} + \lambda_e \|E\|_{1,1} + \frac{\lambda_z}{2} \underbrace{\|Y - YA - E\|_F^2}_Z \quad \text{s.t.} \quad A = C - \text{mdiag}(C)$$

For dual variable Λ , we have augmented Lagrangian function as follows:

$$\mathcal{L}_\mu(C, E, A, \Lambda) = \|C\|_{1,1} + \lambda_e \|E\|_{1,1} + \frac{\lambda_z}{2} \|Y - YA - E\|_F^2 + \langle \Lambda, A - C + \text{mdiag}(C) \rangle + \frac{\mu}{2} \|A - C + \text{mdiag}(C)\|_F^2$$

Optimization - More ADMM

$$E = \min_E \mathcal{L}_\mu(\cdot) = \min_E \|E\|_{1,1} + \frac{1}{2(\lambda_e \lambda_z^{-1})} \|E - (Y - YA)\|_F^2 = \text{prox}_{\lambda_e \lambda_z^{-1} \|\cdot\|_{1,1}}(Y - YX)$$

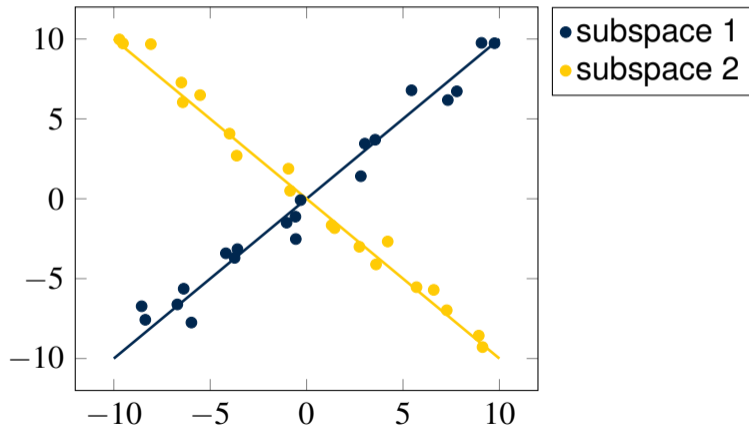
$$\begin{aligned} A &= \min_A \mathcal{L}_\mu(\cdot) = \min_A \frac{\lambda_z}{2} \|Y - YA - E\|_F^2 + \frac{\mu}{2} \|A - C + \text{mdiag}(C)\|_F^2 + \frac{1}{\mu} \|\Lambda\|_F^2 \\ &= (\lambda_z Y^T Y + \mu I)^{-1} [\lambda_z Y^T (Y - E) + \mu (C - \text{mdiag}(C))] - \Lambda \end{aligned}$$

$$\begin{aligned} C &= \min_C \mathcal{L}_\mu(\cdot) = \min_C \|C\|_{1,1} + \frac{\mu}{2} \|A - C + \text{mdiag}(C)\|_F^2 + \frac{1}{\mu} \|\Lambda\|_F^2 \\ &= J - \text{mdiag}(J) \text{ where } J = \text{prox}_{\mu^{-1} \|\cdot\|_{1,1}}(A + \frac{1}{\mu} \Lambda) \end{aligned}$$

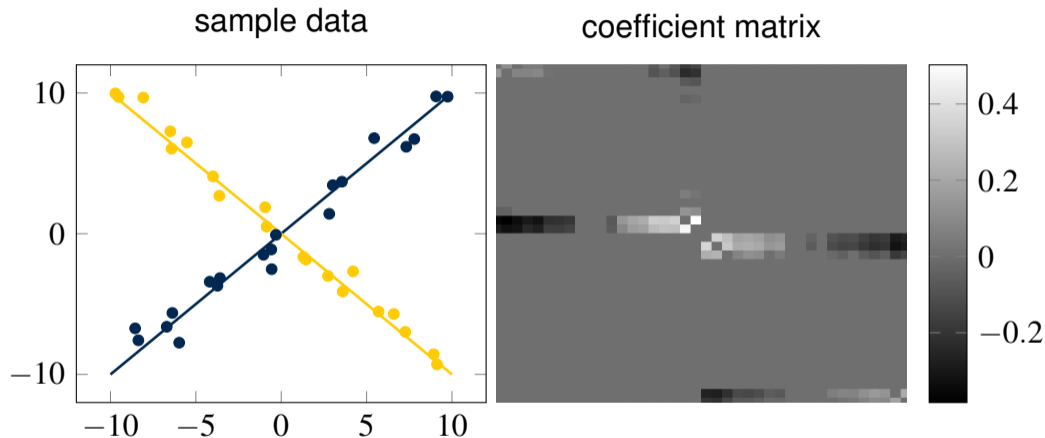
$$\Lambda = \Lambda + \mu(A - C + \text{mdiag}(C))$$

Example Problem

Let $Y = [\underbrace{\mathbf{y}_1^{(1)}, \dots, \mathbf{y}_t^{(1)}}_{\text{subspace 1}}, \underbrace{\mathbf{y}_{t+1}^{(2)}, \dots, \mathbf{y}_n^{(2)}}_{\text{subspace 2}}]$ where $\mathbf{y}_i \in \mathbb{R}^2 \quad \forall i, n = 40$ pts



Example Solution

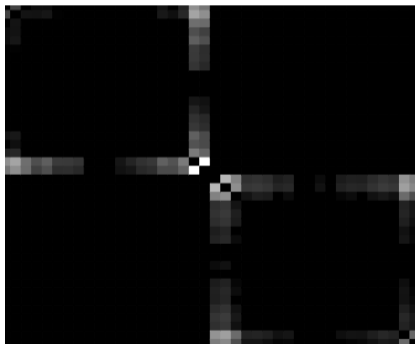


Question: What points did SSC select? What happens near the origin?

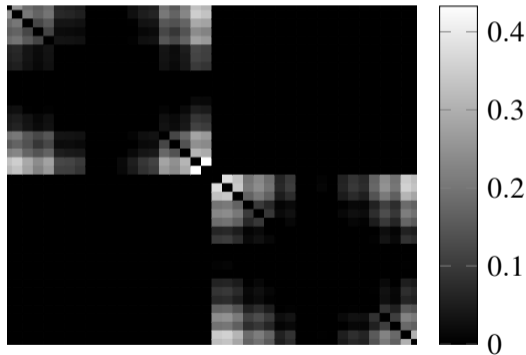
Example Solution

Note: Plots shown are for $W = |C| + |C|^T$ for visualization convenience

high sparsity

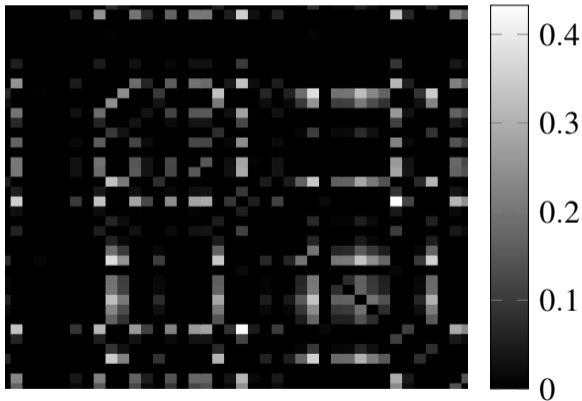


low sparsity



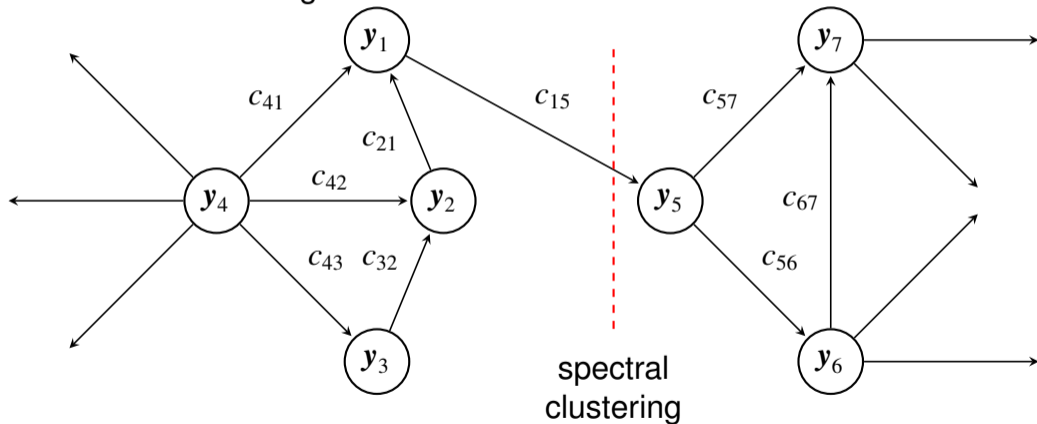
The Problem

$$Y = [\mathbf{y}_1^{(1)}, \mathbf{y}_2^{(2)}, \mathbf{y}_3^{(2)}, \mathbf{y}_4^{(1)}, \mathbf{y}_5^{(1)}, \mathbf{y}_6^{(2)}, \dots] \implies C \text{ below}$$



Spectral Clustering [7]

Data points \implies nodes/vertices
coefficients \implies edges



Graph Notation

$\mathcal{G}(V, C) \implies$ directed graph

Let $W = |C| + |C|^T$ s.t. $\mathcal{G}(V, W) \implies$ undirected graph

W is not self-expressive ($Y \neq YW$) but this is fine

Direction does not matter!

- y_i may use y_j but the other way is not guaranteed for C
- W ensures both nodes are connected together
- connections tell the big picture! (w.r.t. clustering)

Let $D_W = \text{diag}(d_1, \dots, d_n)$ where $d_i = \sum_{j=1}^n w_{ij}$ (degree of nodes)

- If y_i shows up many times \implies high degree node

Graph Laplacians

Let $L = D - W$ be the unnormalized Laplacian matrix

Proposition 1

- 1 For $f \in \mathbb{R}^n$, we have that:

$$f^T L f = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2$$

- 2 L is symmetric and positive semi-definite
- 3 $\lambda_{\min}(L) = 0$ and $\nu_{\min} = \mathbf{1}$
- 4 $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$

Graph Theory Primer

Let A_1, \dots, A_k form a partition on $\mathcal{G}(V, E)$

Let $W(A, B) = \sum_{i \in A, j \in B} w_{ij}$

Then, our problem becomes the following:

$$\min_{A_1, \dots, A_k} \text{cut}(A_1, \dots, A_k) = \frac{1}{2} \sum_{i=1}^k W(A_i, \bar{A}_i)$$

To make things “balanced”, consider this:

$$\min_{A_1, \dots, A_k} \text{RatioCut}(A_1, \dots, A_k) = \frac{1}{2} \sum_{i=1}^k \frac{W(A_i, \bar{A}_i)}{|A_i|}$$

RatioCut

- 1 Define indicator vectors as $h_{i,j} = \{1/\sqrt{|A_j|}$ if $V_i \in A_j$, 0 if $V_i \notin A_j\}$
- 2 Collect vectors into matrix $H \in \mathbb{R}^{n \times k}$ (k clusters)
- 3 Observe that vectors are orthonormal (i.e. $H^T H = I$)
- 4 Without proof, it can be shown that $h_i^T L h_i = (H^T L H)_{ii} = \text{cut}(A_i, \bar{A}_i)/|A_i|$

Thus all together, RatioCut can be written as follows:

$$\text{RatioCut}(A_1, \dots, A_k) = \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{|A_i|} = \sum_{i=1}^k (H^T L H)_{ii} = \text{Tr}(H^T L H)$$

Note: For simplicity, this is assuming unnormalized Laplacian.

More RatioCut

Our optimization problem now becomes this:

$$\min_{A_1, \dots, A_k} \text{RatioCut}(A_1, \dots, A_k) = \text{Tr}(H^T L H) \text{ s.t. } H^T H = I$$

We make this convex by allowing arbitrary real values for H :

$$\min_{H \in \mathbb{R}^{n \times k}} \text{Tr}(H^T L H) \text{ s.t. } H^T H = I$$

The solution is $\hat{H} = [\nu_1(L), \dots, \nu_k(L)]$ assuming $\lambda_1 \leq \dots \leq \lambda_n$

Big Picture!

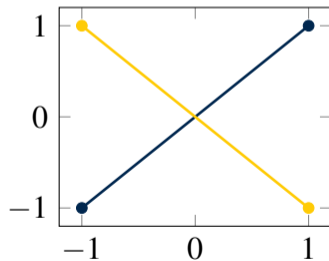
Spectral clustering is a relaxation to various graph cutting problems

Spectral Embedding

For simplicity, assume $n = 4$ and $k = 2$.

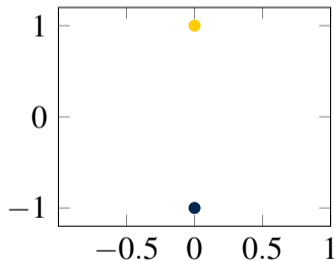
Recall $h_i = \left\{ \begin{array}{l} \sqrt{\frac{|\bar{A}|}{|A|}} \text{ if } V_i \in A, \\ -\sqrt{\frac{|A|}{|\bar{A}|}} \text{ if } V_i \in \bar{A} \end{array} \right\} \implies h^T L h = \text{cut}(A, \bar{A}) / |A|$

sample data



$$h = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix} \implies$$

spectral embedding



Question: Why is it difficult to cluster sample data? Is the right plot easier?

Spectral Clustering

Algorithm SpectralClustering(W, k)

Require: $W \in \mathbb{R}^{n \times n}, k \in \mathbb{N}$

$$D = 0 \in \mathbb{R}^{n \times n}$$

$$D_{i,i} \leftarrow \sum_{j=1}^n w_{ij} \quad \forall i$$

$$L \leftarrow I - D^{-1}W$$

$$E \leftarrow \text{EVD}(L).\text{vectors}$$

$$\nu \leftarrow E_{:,1:k}$$

$$X_1, \dots, X_k \leftarrow \text{kmeans}(\nu^T, k)$$

return $X_1, \dots, X_k \subset Y, X_i \in \mathcal{S}_i$

▷ W = affinity matrix, k = # of subspaces

▷ node degree matrix

▷ normalized laplacian

▷ eigenvalue decomposition

▷ assuming smallest to largest λ

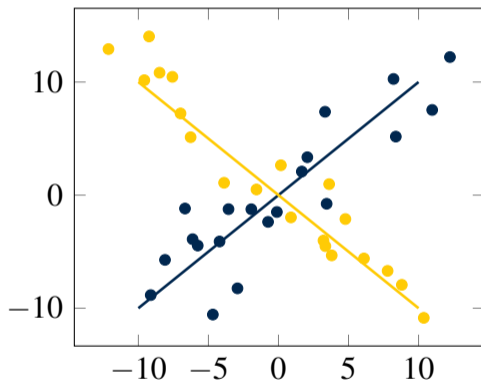
▷ cluster assignments \forall pts

Note!

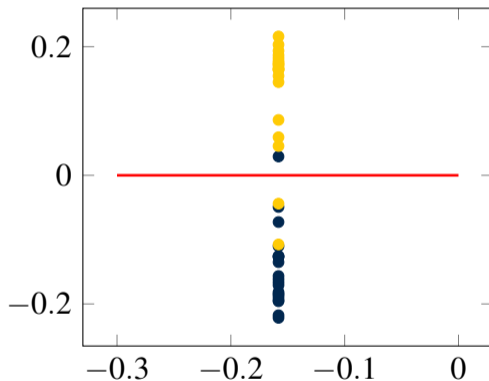
Use partial decomposition methods (e.g. power method) for big data

Example Problem

sample data

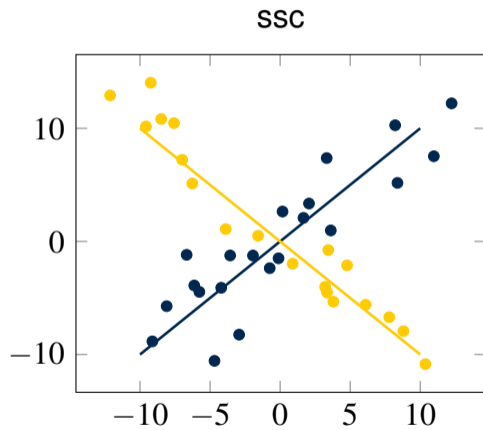
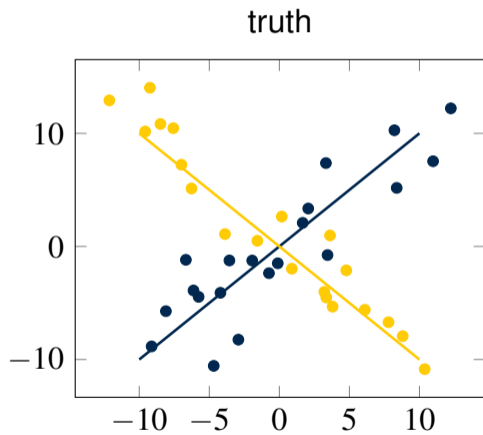


spectral embedding



Question: Conceptually, which points will fail to be classified correctly?

Example Problem



SSC Overview

Algorithm SSC(Y, k)

Require: $Y \in \mathbb{R}^{d \times n}, k \in \mathbb{N}$

$C \leftarrow \text{ADMM}(Y, \lambda, \mu)$

$C_{:,i} \leftarrow C_{:,i} / \|C_{:,i}\|_{\infty} \quad \forall i$

$W \leftarrow |C| + |C|^T$

$\{X_1, \dots, X_k\} \leftarrow \text{SpectralClustering}(W, k)$

for $i = 1, \dots, k$ **do**

$U_i \leftarrow \text{PCA}(X_i \subset Y)$

end for

return $\{X_1, \dots, X_k\}, \{U_1, \dots, U_k\}$

▷ Y = data matrix, k = # of subspaces

▷ solve optimization EQ1

▷ normalize edge weights

▷ undirected graph matrix

▷ see pseudocode 1

▷ for all clusters

▷ use SVD, weighted PCA, etc...

▷ return clustered points and subspaces

Results

SSC: $\|E\|_{1,1} + \|C\|_{1,1} + \|Z\|_F^2$ s.t. $Y = YC + E + Z$

LRR: $\|E\|_{1,1} + \|C\|_* + \|Z\|_F^2$ s.t. $Y = YC + E + Z$

LRSC: $\|E\|_{1,1} + \|C\|_* + \|X - Y\|_F^2 + \|Z\|_F^2$ s.t. $X = XC + E + Z$

LSA: Compute k-NN for y_j , apply PCA to get \hat{S}_j , get affinity by doing $A_{jk} = \exp \left[- \sum_{m=1}^{\min(d_j, d_k)} (\sin \theta_{jk}^m)^2 \right]$ where θ is principal angle

SCC: Randomly take points, form volume of the simplex formed by the points, and get affinity by taking advantage of polar curvature

Dataset:

Extended Yale B dataset \rightarrow face clustering

Results

Clustering Error (%)					
Mean Value	Algorithms				
	LSA	SCC	LRR	LRSC	SSC
2 subjects	32.80	16.62	9.52	5.32	1.86
3 subjects	52.29	38.16	19.52	8.47	3.10
5 subjects	58.02	58.90	34.16	6.90	4.31
8 subjects	59.19	66.11	41.19	23.72	5.85
10 Subjects	60.42	73.02	38.85	30.36	10.94

Note: No pre-processing of the data for these results

Thanks for listening!

Questions?

References

- [1] J. Ho, M.-H. Yang, J. Lim, K.-C. Lee, and D. Kriegman, "Clustering appearances of objects under varying illumination conditions," in *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.*, IEEE, vol. 1, 2003, pp. I–I.
- [2] W. Hong, J. Wright, K. Huang, and Y. Ma, "Multiscale hybrid linear models for lossy image representation," *IEEE Transactions on Image Processing*, vol. 15, no. 12, pp. 3655–3671, 2006.
- [3] A. Y. Yang, J. Wright, Y. Ma, and S. S. Sastry, "Unsupervised segmentation of natural images via lossy data compression," *Computer Vision and Image Understanding*, vol. 110, no. 2, pp. 212–225, 2008.
- [4] R. Vidal, R. Tron, and R. Hartley, "Multiframe motion segmentation with missing data using powerfactorization and gpca," *International Journal of Computer Vision*, vol. 79, no. 1, pp. 85–105, 2008.
- [5] R. Vidal, "Subspace clustering," *IEEE Signal Processing Magazine*, vol. 28, no. 2, pp. 52–68, 2011.
- [6] E. Elhamifar and R. Vidal, "Sparse subspace clustering: Algorithm, theory, and applications," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 11, pp. 2765–2781, 2013.
- [7] U. Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.

Proposition 2

Definition - Connected Components

For $\mathcal{G}(V, E)$, nonempty set $A \subset V$ is a connected component iff

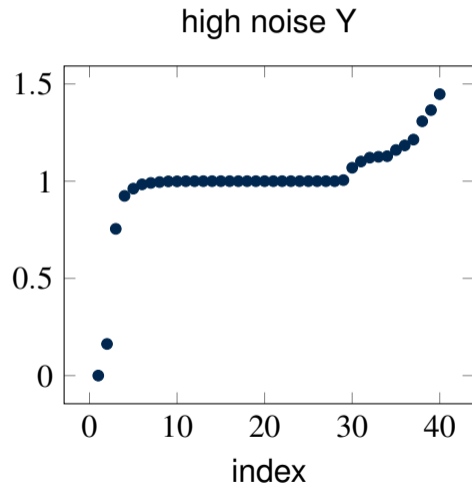
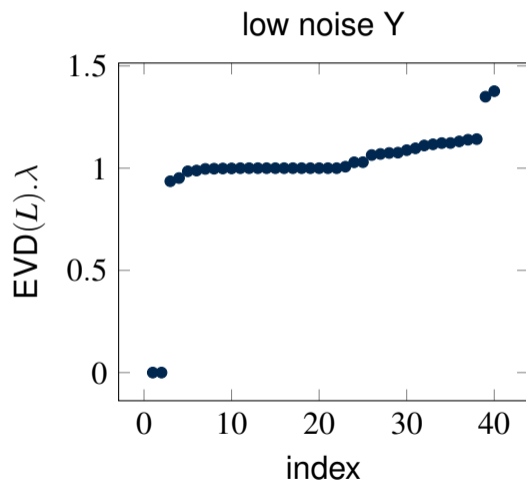
- 1 Any two vertices in A can be joined s.t. all intermediate points also lie in A
- 2 No connections between vertices in A and \bar{A}

Proposition 2 [7]

Let $\mathcal{G}(\cdot)$ be an undirected graph with non-negative weights. For L ,

- 1 Multiplicity k of $\lambda = 0$ is # of connected components
- 2 $\nu_{\lambda=0}(L) = \{\mathbb{1}_{A_1}, \dots, \mathbb{1}_{A_k}\}$

Spectral Gap Heuristic



Elbow Method Heuristic

Consider K-means perspective instead of graph theory

μ_j : sample mean for cluster j

c_j : cluster j

k : number of subspaces

$V_{i,:}^T$: spectral embedding point i

Define within cluster sum of squares as the following:

$$WCSS(k) = \sum_{j=1}^k \sum_{V_{i,:}^T \in c_j} \|V_{i,:}^T - \mu_j\|_2^2$$

Goal: Plot $WCSS(k)$ and find the “elbow” or “knee” of the plot

Optimization - Proximal Approach

For simplicity, assume no outliers exist. Then,

$$\min_C \|C\|_{1,1} + \frac{\lambda_z}{2} \|Y - YC\|_F^2 \quad \text{s.t.} \quad \text{diag}(C) = \mathbf{0}$$

Introduce binary mask $\Omega = \mathbf{1}'\mathbf{1} - I$ and let $\mathcal{P}(X) = (\Omega \odot X)$ s.t.

$$\min_C \|\mathcal{P}(C)\|_{1,1} + \frac{\lambda_z}{2} \|Y - Y\mathcal{P}(C)\|_F^2$$

Assuming $C_0 = \{C : \text{diag}(C) = \mathbf{0}\}$, then the solution is as follows:

$$\left. \begin{array}{l} \nabla(C) = \mathcal{P}(-Y^T(Y - YC)) \\ \text{prox} = \text{soft}(C, \lambda_z^{-1}) \end{array} \right\} \quad \text{Note: } \nabla_C \left(\frac{1}{2} \|C \odot \Omega\|_F^2 \right) = \Omega \odot C \odot \Omega$$

Graph Laplacians Part 2

From proposition 1, we see that:

$$\mathbf{f}^T \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2$$

Assume our similarity measure $w_{ij} = 1/d_{ij}^2$
where $d_{ij} = \text{dist}(\mathbf{y}_i, \mathbf{y}_j) = \|\mathbf{y}_i - \mathbf{y}_j\|_2$

Then, we have the following:

$$\mathbf{f}^T \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2 \approx \frac{1}{2} \sum_{i,j=1}^n \left(\frac{f_i - f_j}{d_{ij}} \right)^2 \approx \int |\nabla g|^2 \, dx$$

Laplacian & RatioCut Relationship

Assume $k = 2$ and $f_i = \left\{ \begin{array}{l} \sqrt{\frac{|\bar{A}|}{|A|}} \text{ if } v_i \in A, \\ -\sqrt{\frac{|A|}{|\bar{A}|}} \text{ if } v_i \in \bar{A} \end{array} \right\}$

$$\begin{aligned} f^T L f &= \frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2 \\ &= \frac{1}{2} \sum_{i \in A, j \in \bar{A}} w_{ij} \left(\sqrt{\frac{|\bar{A}|}{|A|}} + \sqrt{\frac{|A|}{|\bar{A}|}} \right)^2 + \frac{1}{2} \sum_{i \in \bar{A}, j \in A} w_{ij} \left(\sqrt{\frac{|\bar{A}|}{|A|}} - \sqrt{\frac{|A|}{|\bar{A}|}} \right)^2 \\ &= \text{cut}(A, \bar{A}) \left(\frac{|\bar{A}|}{|A|} + \frac{|A|}{|\bar{A}|} + 2 \right) = \text{cut}(A, \bar{A}) \left(\frac{|A| + |\bar{A}|}{|A|} + \frac{|A| + |\bar{A}|}{|\bar{A}|} \right) \\ &= (|A| + |\bar{A}|) \left(\frac{\text{cut}(A, \bar{A})}{|A|} + \frac{\text{cut}(A, \bar{A})}{|\bar{A}|} \right) = |V| \cdot \text{RatioCut}(A, \bar{A}) \end{aligned}$$

Note: For simplicity, this is assuming unnormalized Laplacian.

Proposition 1 Proof

Given that $L = D - W$, we have that:

$$\begin{aligned} f^T L f &= f^T D f - f^T W f = \sum_{i=1}^n d_i f_i^2 - \sum_{i,j=1}^n f_i f_j w_{ij} \\ &= \frac{1}{2} \left(\sum_{i=1}^n d_i f_i^2 - 2 \sum_{i,j=1}^n f_i f_j w_{ij} + \sum_{j=1}^n d_j f_j^2 \right) \\ &= \frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2 \end{aligned}$$

Affine Subspace Extension

① $k = 1$

Get sample mean value, de-mean data, & take SVD

② $k > 1$

Let $\mathbf{x}_1, \dots, \mathbf{x}_n \in A = \{\mathbf{d} + \mathbf{z} : \mathbf{z} \in \mathcal{S}\}$ where \mathbf{d} is displacement vector

Let $\lambda_1, \dots, \lambda_n \in \mathbb{R}$ be scalar values

$$\sum_i^n \lambda_i = 1 \implies \lambda_1 \mathbf{x}_1 + \dots + \lambda_n \mathbf{x}_n \in A$$

For example, if $\lambda_i = \frac{1}{n}$, then we get the centroid location

Question: Given this, how would one generalize for SSC?