

Automatic Model Selection in Subspace Clustering via Triplet Relationships

Jufeng Yang¹, Jie Liang¹, Kai Wang¹, Yong-Liang Yang², Ming-Ming Cheng¹

¹College of Computer and Control Engineering
Nankai University
No.38 Tongyan Road, Tianjin, China

²Department of Computer Science
University of Bath
Claverton Down, Bath, United Kingdom

Abstract

This paper addresses both the model selection (*i.e.* estimating the number of clusters K) and subspace clustering problems in a unified model. The real data always distribute on a union of low-dimensional sub-manifolds which are embedded in a high-dimensional ambient space. In this regard, the state-of-the-art subspace clustering approaches firstly learn the affinity among samples, followed by a spectral clustering to generate the segmentation. However, arguably, the intrinsic geometrical structures among samples are rarely considered in the optimization process. In this paper, we propose to simultaneously estimate K and segment the samples according to the local similarity relationships derived from the affinity matrix. Given the correlations among samples, we define a novel data structure termed the Triplet, each of which reflects a high relevance and locality among three samples which are aimed to be segmented into the same subspace. While the traditional pairwise distance can be close between inter-cluster samples lying on the intersection of two subspaces, the wrong assignments can be avoided by the hyper-correlation derived from the proposed triplets due to the complementarity of multiple constraints. Sequentially, we propose to greedily optimize a new model selection reward to estimate K according to the correlations between inter-cluster triplets. We simultaneously optimize a fusion reward based on the similarities between triplets and clusters to generate the final segmentation. Extensive experiments on the benchmark datasets demonstrate the effectiveness and robustness of the proposed approach.

Introduction

State-of-the-art subspace clustering methods model high-dimensional data samples $\mathbf{X} = \{\mathbf{x}_i \in \mathbb{R}^D\}_{i=1}^N$ into a union of low-dimensional subspaces $\{S_j\}_{j=1}^K$ (Elhamifar and Vidal 2009; Vidal 2011), where D and N denote the dimensionality and scale of the given dataset, respectively, $K \geq 1$. The primary step of the spectral based subspace clustering methods is to calculate the coefficient matrix \mathbf{C} by solving an optimization problem as follows:

$$\min_{\mathbf{C}} L(\mathbf{X}\mathbf{C}, \mathbf{X}) + \lambda \|\mathbf{C}\|_{\xi}, \quad (1)$$

where $L(\cdot, \cdot) : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^+$ denotes the reconstruction loss, λ is the trade-off parameter and $\|\cdot\|_{\xi}$ denotes the regularization term where different ξ 's lead to ℓ_0 , ℓ_1 , ℓ_2 , ℓ_{∞} or the

nuclear norm (Vidal 2011; Yang et al. 2016). The algorithm then employs the spectral clustering (Shi and Malik 2000) on the affinity matrix derived from \mathbf{C} for final assignments. Note that in practice, both the number of subspaces K and their dimensions $\{d_j\}_{j=1}^K$ are always unknown (Vidal 2011; Wang and Zhu 2015). Hence, the goals of subspace clustering include finding the appropriate K and assigning data points into K clusters (Elhamifar and Vidal 2013; Li et al. 2016).

However, modeling the cluster number in the optimization framework is difficult (Li et al. 2016), since the definition of clusters is unquantifiable on the complex data space. Sequentially, most of the spectral based subspace clustering algorithms set the parameter K manually, which achieve state-of-the-art performance on the confined applications where the number of clusters is fixed and given.

Clustering aims to group the similar patterns into the same cluster by maximizing the inter-cluster dissimilarity and the intra-cluster similarity. An effective way for estimating K is to map the original samples into intrinsic correlation space, followed by an iterative optimization according to the local similarity relationships among samples. Elhamifar et al. (Elhamifar and Vidal 2013) propose that the FC achieves a block-diagonal structure with K blocks corresponding to K clusters, where \mathbf{C} is derived from (1) and Γ denotes a proper permutation matrix. Moreover, Peng et al. (Peng et al. 2016) verify the intra-subspace projection dominance (IPD) of \mathbf{C} derived from (1) with various ξ 's, *i.e.*, for all $\mathbf{x}_p, \mathbf{x}_q \in S$ and $\mathbf{x}_k \notin S$, we have $c_{pq} \geq c_{pk}$. Hence, considering a graph $\mathcal{G} = (\mathbf{X}, \mathbf{C})$ (Nasihatkon and Hartley 2011), *i.e.*, $\mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}$ being the vertices and c_{ij} being the weight between \mathbf{x}_i and \mathbf{x}_j , the automatic segmentation can be greedily generated via the following two steps inspired by the density based algorithms (Rodriguez and Laio 2014): 1) finding \hat{K} initialized clusters $\{\mathcal{C}_i\}_{i=1}^{\hat{K}}$ by minimizing the weighted sum of all inter-cluster connections. 2) assigning the remaining samples \mathbf{x} to a proper \mathcal{C} by maximizing the weighted connections between \mathbf{x} and \mathcal{C} .

Yet, since the low-dimensional sub-manifolds can be very dense (Peng, Zhang, and Yi 2013), the points \mathbf{x}_i and \mathbf{x}_j which are close regarding the pairwise distance may not belong to the same subspace, especially near the intersection of two subspaces. A hypergraph $\hat{\mathcal{G}}$ in which one edge can link up more than two vertices (Gao, Tsang, and Chia 2013;

Kim et al. 2014) is then proposed to replace the pairwise \mathcal{G} . Besides, the local geometry of each \mathbf{x}_i can be linearly reconstruct from its correlated points (*i.e.* \mathbf{x}_j for which the c_{ij} is large) (Yin, Gao, and Lin 2016). In this paper, we further introduce a novel data structure termed the *Triplet*, *i.e.*, τ , to explore the local geometry with hyper-correlations rather than pairwise similarity. Each τ contains three points, *i.e.*, $\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k$, together with their correlations, *i.e.*, c_{ij}, c_{jk} and c_{ki} , for which we considered as a meta-element for clustering. All the correlations are large enough to ensure that the three points are strongly connected according to the IPD property (Peng et al. 2016). The proposed triplet relationship achieves favorable theoretical guarantee against the pairwise one in the following two folds. On one hand, it is more robust when partitioning the inter-cluster samples near the intersection of two sub-manifolds due to the complementarity of multiple constraints. On the other hand, it evokes mutual restrictions of neighbored samples thus depicts a local geometrical structure, by which we can calculate the segmentation greedily.

Sequentially, in this paper, we propose a framework termed the autoSC to simultaneously estimate the number of clusters and segment the samples by exploring the local geometrical structure derived from the given coefficient matrix \mathbf{C} . Specifically, we first generate \mathbf{C} by solving the subspace representation problem in (1), followed by extracting the triplet relationship τ . Then, we greedily initialize \hat{K} clusters, *i.e.*, $\{\mathcal{C}_i\}_{i=1}^{\hat{K}}$, to maximize the inter-cluster dissimilarity among triplets, which is obtained via a new model selection reward. Finally, we assign each of the remaining samples into \mathcal{C} to maximize the intra-cluster similarity by optimizing a new fusion reward.

We have mainly two contributions. First, we define the triplet relationship τ which induces a high relevance and locality among three samples to explore the local similarity, and verify its favorable performance against the traditional pairwise relation. Second, we design a greedy framework for joint model selection and clustering utilizing the intrinsic geometrical structures depicted by the proposed triplet relationships. Extensive experiments on benchmarks indicate that our proposed autoSC outperforms the state-of-the-art methods.

Related Work

Subspace Representations

Automatically finding the clusters of samples is a crucial issue in computer vision (Elhamifar and Vidal 2009; Zhao et al. 2017; Schroff, Kalenichenko, and Philbin 2015). Traditional subspace clustering first constructs a sparse linear representation of each data sample using the remaining as a dictionary (Li and Vidal 2015; Cheng et al. 2016). It generates a coefficient matrix $\mathbf{C} \in \mathbb{R}^{N \times N}$ by solving the optimization problem in (1). By modifying the regularization term, *i.e.*, the value of ξ in $\|\mathbf{C}\|_{\xi}$, researchers generate coefficient matrix \mathbf{C} considering different intrinsic properties of the data space \mathbf{X} .

The ℓ_0 and ℓ_1 norm based algorithms (Yang et al. 2016) intend to eliminate most non-zero values of \mathbf{C} , since c_{ij} denotes the similarity between samples \mathbf{x}_i and \mathbf{x}_j and the con-

nections only exist when \mathbf{x}_i and \mathbf{x}_j belong to same subspace. Meanwhile, there are ℓ_2 and nuclear norm based methods (Hu et al. 2014; Liu, Latecki, and Yan 2010) intending to preserve as many the non-zero values in \mathbf{C} . These frameworks interpret the coefficient matrix by that each column of \mathbf{C} , *i.e.*, \mathbf{c}_i , is the self-representation of \mathbf{x}_i , thus they preserve all values to ensure the mapping invariance (also termed as the grouping effect (Lu et al. 2012)). Recently, the mixed norms such as trace Lasso (Lu et al. 2015) and elastic net (You et al. 2016) are applied to the optimization model for the tradeoff between the sparsity and the grouping effect. There also arise frameworks which incorporate various constraints into the model. For example, (Guo, Gao, and Li 2013) considers the similarity between neighbors in sequential data, of which the new penalty $\|\mathbf{C}\mathbf{R}\|_1$ forces consecutive columns of \mathbf{C} to be similar where \mathbf{R} is a lower triangular matrix with -1 on the diagonal and 1 on the second diagonal.

In this paper, we generate the triplet relationship by exploring the intrinsic geometrical structures depicted in the coefficient matrix \mathbf{C} from off-the-shelf subspace representation modules. We verify that the proposed method is robust to the combination of \mathbf{C} with various properties.

Estimating the Number of Clusters

However, in most real applications, K is unknown to the clustering algorithms (Elhamifar and Vidal 2009). Considering the block-diagonal structure of the coefficient matrix (Feng et al. 2014), Liu et al. (Liu et al. 2013) propose to utilize the heuristic estimator: $\hat{K} = N - \text{round}(\sum f_{\epsilon}(\sigma_i))$, where ϵ is a cut-off threshold, σ_i denotes a singular value of normalized Laplacian matrix and f_{ϵ} is a summation function which counts different values regarding that $\sigma_i < \epsilon$. The singular based methods (Favaro, Vidal, and Ravichandran 2011; Elhamifar and Vidal 2009) rely on the large gap between singulars, which is effective only when the sub-manifolds are sparse in the ambient space.

The density based methods always iteratively find the optimal number of clusters and the optimal assignment by analyzing the density among samples. In (Rodriguez and Laio 2014), Rodriguez et al. assume that the cluster centers are characterized by a higher density than its neighbors and different centers should be far enough. For each sample \mathbf{x}_i , its local density $\rho_i = \sum_j \mathcal{X}(d_{ij} - d_c)$ and the distance $\delta_i = \min_{\mathbf{x}_j: \rho_j > \rho_i} (d_{ij})$ from points of higher density are iteratively calculated for comparison. The algorithm finds a tradeoff between ρ and δ to decide the cluster centers and the assignment of the remained samples. Wang et al. (Wang and Zhu 2015) propose the Dirichlet process based Bayesian non-parametric method (DP-space) to exploit the tradeoff between data fitness and model complexity, which is more tolerate to noisy and outlier values than the alternative algebraic and geometry solutions. Li et al. propose the SCAMS (Li, Cheong, and Zhou 2014) which penalizes the clustering cost by minimizing the Frobenius inner product $-\langle \mathbf{C}, \mathbf{B} \rangle$ and estimates K by minimizing $\text{rank}(\mathbf{B})$, where \mathbf{B} is a binary relationship matrix encoding the pairwise relationships among samples. Correlation clustering (CC) (Beier, Hamprecht, and Kappes 2015) minimizes the sum of the weights of the cut edges on

an undirected graph with positive and negative edge weights, such that $y^* = \arg \min_y \sum w_{ij} y_{ij}$ where the label $y_{ij} = 1$ as cut and $y_{ij} = 0$ as uncut. However, these algorithms only take the pairwise correlation into consideration, and they can not be robust when subspaces are very dense and samples of different subspaces are not completely distinguished.

The hyper-graph relation (Lu et al. 2016; Purkait et al. 2014) avoids such drawbacks and the literature follows two different directions. Some transform the hyper relationship into another pairwise graph (Gao, Tsang, and Chia 2013; Schölkopf, Platt, and Hofmann 2006), followed by the conventional graph clustering method (Shi and Malik 2000) to generate the segmentation. There are also generalization methods (Liu, Latecki, and Yan 2010; Li et al. 2016) extending the pairwise graph to the hyper-graph or tensor analysis. For example, the tensor affinity variant of SCAMS, i.e., SCAMSTA (Li et al. 2016), exploits the higher order mathematical structures by providing multi groups of nodes in \mathcal{Z} , i.e., $\mathcal{Z} = \sum_k z_r \circ z_r \circ \dots \circ z_r \in \{0, 1\}^{N \times N \times \dots \times N}$, where $z_r \circ z_r$ denotes the outer product and $z_r \in \{0, 1\}^N$ is the indicator vector.

In this paper, we estimate K by initializing clusters with optimal inter-cluster dissimilarities. We obtain the estimation by exploring the local correlations induced by the proposed triplet relationships, each of which depicts a hyper-similarity among three samples. Both theoretical analysis and experimental results demonstrate the effectiveness of the proposed method.

Methodology

Notations and Problem Formulation

Given the data matrix $\mathbf{X} = \{\mathbf{x}_i \in \mathbb{R}^D\}_{i=1}^N$, subspace clustering (SC) algorithms solve the optimization problem in (1) to generate the coefficient matrix \mathbf{C} , of which each entry c_{ij} reflects the similarity of \mathbf{x}_i and \mathbf{x}_j . For each data sample $\mathbf{x}_j \in \mathcal{S}$, the SC algorithms take advantage of the self-expressive property, i.e., each data sample can be reconstructed by a linear combination of other points in the dataset (Elhamifar and Vidal 2013; Belkin and Niyogi 2001). As a result, \mathbf{x}_j can be written as

$$\mathbf{x}_j = \mathbf{X} \mathbf{c}_j, \text{ s.t. } c_{jj} = 0, \quad (2)$$

where the data matrix \mathbf{X} is considered as a self-expression dictionary and $\mathbf{c}_j = [c_{1j}, c_{2j}, \dots, c_{Nj}]$ contains the coefficients of the combination. Considering a regularization with various well-designed norms on \mathbf{c}_j , the calculation system in (2) has the ability to preserve only the combinations among samples in \mathcal{S} . Inspired by (Peng et al. 2016), we first collect the nearest neighbors for \mathbf{x}_j with top m coefficients in \mathbf{c}_j .

Definition 1. (m Nearest Neighbors) *The m nearest neighbors for each data point \mathbf{x}_j , i.e. $N_m(\mathbf{x}_j) \in \mathbb{R}^{1 \times m}$, are defined as follows:*

$$N_m(\mathbf{x}_j) = \arg \max_{\{\mathbf{x}_{i_l}\}} \sum_{l=1}^m |c_{i_l, j}|, \quad (3)$$

where i_l is the set of ordinals for the nearest neighbors, $c_{i_l, j}$ denotes the coefficient of \mathbf{x}_{i_l} and \mathbf{x}_j .

Based on the generated nearest neighbors, we define the triplet relationship to explore the local hyper-correlation among samples.

Definition 2. (Triplet Relationship) *The three samples \mathbf{x}_i , \mathbf{x}_j and \mathbf{x}_k form a triplet relationship if and only if they satisfy:*

$$\mathbf{1}_{\mathbf{x}_i \in N_m(\mathbf{x}_j)} \times \mathbf{1}_{\mathbf{x}_j \in N_m(\mathbf{x}_k)} \times \mathbf{1}_{\mathbf{x}_k \in N_m(\mathbf{x}_i)} = 1, \quad (4)$$

where $\mathbf{1}_{\mathbf{x} \in N_m}$ is the indicator function which equals to 1 if $\mathbf{x} \in N_m$ and 0 otherwise.

For easy illustration, we introduce the triplet matrix $\mathbf{T} \in \mathbb{R}^{n \times 3}$, where n denotes the number of triplets. Each row of \mathbf{T} , i.e., $\tau = \{\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k\}$, denotes a triplet where \mathbf{x}_i , \mathbf{x}_j and \mathbf{x}_k satisfy the requirement in (4).

Proposition 1. *Given arbitrary three samples in one triplet, i.e., $\tau = \{\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k\} \in \mathbf{T}$, we have $\{\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k\} \subset S$, where \mathbf{T} denotes the set of triplets and S denotes a specific subspace.*

Proof. Let $\tilde{\mathbf{c}}$ be the optimal solution of the optimization function:

$$\min_{\mathbf{c}} \|\mathbf{x} - \mathbf{X} \mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_\xi, \quad (5)$$

where $\|\mathbf{c}\|_\xi$ denotes the ℓ_ξ norm of \mathbf{c} . We divide the vector $\tilde{\mathbf{c}}$ into the following two parts:

$$\tilde{\mathbf{c}} = [\tilde{\mathbf{c}}_{D_{\mathbf{x}}}, \tilde{\mathbf{c}}_{D_{-\mathbf{x}}}]^\top, \quad (6)$$

where $D_{\mathbf{x}}$ and $D_{-\mathbf{x}}$ denote the collections of intra-cluster samples and inter-cluster samples of the sample \mathbf{x} , respectively. According to the Intra-subspace Projection Dominance (IPD) which is proved in (Peng et al. 2016), we have

$$[\tilde{\mathbf{c}}_{D_{\mathbf{x}}}]_{r_{\mathbf{x}}, 1} > [\tilde{\mathbf{c}}_{D_{-\mathbf{x}}}]_{1, 1}, \quad (7)$$

for ℓ_ξ being equal to ℓ_1 , ℓ_2 , ℓ_∞ and nuclear norm, where $[\tilde{\mathbf{c}}_{D_{\mathbf{x}}}]_{r_{\mathbf{x}}, 1}$ denotes the $r_{\mathbf{x}}$ -th largest absolute value of the entries of $[\tilde{\mathbf{c}}_{D_{\mathbf{x}}}]$, $r_{\mathbf{x}}$ denotes the dimensionality of the subspace of the intra-cluster samples of \mathbf{x} .

Assume a binary coefficient matrix \mathbf{C}^* satisfies the following function:

$$c_{ij}^* = \begin{cases} 1, & c_{ij} \in \mathbf{c}_{r_{\mathbf{x}}, m}, \\ 0, & \text{otherwise,} \end{cases} \quad (8)$$

where $\mathbf{c}_{r_{\mathbf{x}}, m}$ denotes the m largest values of \mathbf{c} . According to (4), the triplet $\tau = \{\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k\}$ satisfies $c_{ij}^* \times c_{jk}^* \times c_{ki}^* = 1$. Besides, we have $m \ll \frac{N}{K}$, hence the preserved m values are of the intra-cluster samples of \mathbf{x} . Since $c_{ij}^* \times c_{jk}^* \times c_{ki}^* = 1$ and c^* can only take the value of 0 and 1, we have $c_{ij}^* = c_{jk}^* = c_{ki}^* = 1$. Therefore, the three samples belong to same subspace, i.e., $\{\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k\} \subset S$. \square

In contrast to the traditional pairwise relationship, the triplet is more robust for partitioning the inter-cluster samples near the intersection of two sub-manifolds due to the complementarity derived from (4). Meanwhile, the triplet relationship depicts a local geometrical structure which enables us to jointly estimate the K and clustering using a greedy framework according to the local density among triplets.

In the clustering process, we iteratively fuse samples into the clusters, *i.e.*, $\{\mathcal{C}_i\}_{i=1}^{\widehat{K}}$, where \widehat{K} denotes the estimated number of clusters. Therefore, In the I -th iteration, we define the set of “in-clusters” triplets \mathbf{T}_{in}^I which are already assigned into clusters, and the set of “out-of-clusters” triplets \mathbf{T}_{out}^I which should be assigned in the following iterations. For easy illustration, we reshape the matrices $\mathbf{T}_{in}^I \in \mathbb{R}^{p \times 3}$ and $\mathbf{T}_{out}^I \in \mathbb{R}^{q \times 3}$ to the vectors $\mathbf{X}_{in}^I \in \mathbb{R}^{3p}$ and $\mathbf{X}_{out}^I \in \mathbb{R}^{3q}$, both of which preserve the frequency of each sample. We then propose to optimize two new rewards, *i.e.*, the model selection and the fusion reward.

Definition 3. (Model Selection Reward) *The model selection reward $R_m(\mathcal{C})$ for the clusters $\{\mathcal{C}_i\}_{i=1}^{\widehat{K}}$ is defined as:*

$$R_m(\mathcal{C}) = \sum_i f(\mathcal{C}_i | \mathbf{X}_{out}^I) - \lambda_m \sum_i f(\mathcal{C}_i | \mathbf{X}_{in}^I), \quad (9)$$

where $f(\mathcal{C} | \mathbf{X})$ is a counting function on the frequency that $\mathbf{x} \in \mathbf{X}_{out}^I$ or $\mathbf{x} \in \mathbf{X}_{in}^I$ for all $\mathbf{x} \in \mathcal{C}_i$, λ_m denotes the trade-off.

By maximizing the model selection reward $R_m(\mathcal{C})$, we generate the initialized cluster $\{\mathcal{C}_i\}_{i=1}^{\widehat{K}}$ with the following two advantages where \widehat{K} is the estimated number of clusters. 1) Samples in \mathcal{C} are of high density, *i.e.*, have large amount of correlated samples in \mathbf{X}_{out} , which enables to merge as many in the next iteration; 2) Each \mathcal{C} is of little correlation with samples in \mathbf{X}_{in} , which eliminates the overlap of the inter-cluster. Thus, we can simultaneously estimated \widehat{K} and initialize the clusters by optimizing the model selection reward R_m .

Definition 4. (Fusion Reward) *The fusion reward optimizes the probability that $\mathbf{x}_j \in \mathbf{X}_{out}$ being assigned into the cluster \mathcal{C}_i , which is defined as:*

$$R_f^i(\mathcal{C}_i | \mathbf{x}_j \in \mathbf{X}_{out}) = f(\mathbf{x}_j | \mathcal{C}_i) + \lambda_f f(\mathbf{N}_m(\mathbf{x}_j) | \mathbf{N}_m(\mathcal{C}_i)), \quad (10)$$

where $\mathbf{N}_m(\mathbf{x}_j)$ denotes the m nearest neighbors of \mathbf{x}_j and $\mathbf{N}_m(\mathcal{C}_i)$ denotes the set of m nearest neighbors of samples in \mathcal{C}_i , λ_f denotes the trade-off.

We calculate \widehat{K} fusion rewards $\{R_f^i\}_{i=1}^{\widehat{K}}$ for each \mathbf{x}_j , which represent the probabilities that \mathbf{x}_j being assigned into clusters $\{\mathcal{C}_i\}_{i=1}^{\widehat{K}}$, respectively. We then merge \mathbf{x}_j into the cluster with the largest fusion reward, and move the \mathbf{x}_j from \mathbf{X}_{out} to \mathbf{X}_{in} .

To determine the first triplet for construct a new cluster, we propose to maximize the local density defined as follows.

Definition 5. (Local Density) *The local density of the triplet τ regarding to the \mathbf{X}_{out} is defined as follows:*

$$\rho(\tau, \mathbf{X}_{out}) = \sum_{j=1}^{|\mathbf{n}|} f(\mathbf{x}_{n_j} | \mathbf{X}_{out}), \quad (11)$$

where \mathbf{x}_{n_j} denotes the sample in current triplet τ and \mathbf{n} is the set of their ordinals, $|\mathbf{n}|$ denotes the scale of \mathbf{n} .

Also, to determine the optimal triplet to merge into the initialized clusters, we define the connection score as follows.

Algorithm 1 : Automatic Subspace Clustering (autoSC)

Input: $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$.

- 1: Calculate the correlation matrix \mathbf{C} by (1);
- 2: **for** $i = 1 : N$ **do**
- 3: Calculate the m Nearest Neighbors $\mathbf{N}_m(\mathbf{x}_i)$ by (3);
- 4: **end for**
- 5: Generate the triplet matrix $\mathbf{T} \in \mathbb{R}^{n \times 3}$ by (14);
- 6: Reshape \mathbf{T} to $\mathbf{X}_{out} \in \mathbb{R}^{3n}$, $\mathbf{X}_{in} = \emptyset$;
- 7: $\widehat{K} = 1$;
- 8: Calculate $\tau_{ini}^{\widehat{K}}$ by (15);
- 9: **while** $\rho(\tau_{ini}^{\widehat{K}}, \mathbf{X}_{out}) > \rho(\tau_{ini}^{\widehat{K}}, \mathbf{X}_{in})$ **do**
- 10: $\mathcal{C}_{\widehat{K}} = \tau_{ini}^{\widehat{K}}$;
- 11: $\mathcal{C}_{\widehat{K}} = \mathcal{C}_{\widehat{K}} \cup \{\tau^*\}$ where τ^* is calculated by (16);
- 12: $\mathbf{X}_{in} = \mathbf{X}_{in} \cup \tau_{ini}^{\widehat{K}} \cup \{\tau^*\}$;
- 13: $\mathbf{X}_{out} = \mathbf{X}_{out} / (\tau_{ini}^{\widehat{K}} \cup \{\tau^*\})$;
- 14: $\widehat{K} = \widehat{K} + 1$;
- 15: Calculate $\tau_{ini}^{\widehat{K}}$ by (15);
- 16: **end while**
- 17: Merge \mathcal{C}_i and \mathcal{C}_j if we have (18); Get \widehat{K} clusters;
- 18: **for** $j = 1 : |\mathbf{X}_{out}|$ **do**
- 19: Calculate \mathcal{C}^* for \mathbf{x}_j by (19);
- 20: **end for**

Output: The cluster assignment $\{\mathcal{C}_i\}_{i=1}^{\widehat{K}}$.

Definition 6. (Connection Score) *The connection score of the sample \mathbf{x}_i towards the sample \mathbf{x}_j is defined as:*

$$s(\mathbf{x}_i, \mathbf{x}_j) = f \left(\mathbf{x}_i \mid \bigcap_{k=1}^{n'} (\mathbf{1}_{\mathbf{x}_j \in \tau_k} \times \tau_k) \right), \quad (12)$$

where $\mathbf{1}_{\mathbf{x}_j \in \tau_k}$ is equal to 0 when $\mathbf{x}_j \in \tau_k$ and otherwise equal to 1, n' is the number of all triplets in \mathbf{T}_{out} .

Note both definitions can be easily extended with different components, *e.g.*, τ 's or \mathcal{C} 's, of which the extension can be found in the supplementary material.

The Clustering Model based on Triplets

We greedily optimize the proposed model selection reward R_m and the fusion reward R_f in autoSC to simultaneously estimate the number of clusters and generate the segmentation among samples:

$$\begin{aligned} \max_{\mathcal{G}, \widehat{K}} \quad & \sum_{k=1}^{\widehat{K}} R_m(\mathcal{G}_k) + \lambda \sum_{k=1}^{\widehat{K}} R_f(\mathcal{G}_k | \mathbf{X}), \\ \text{s.t.} \quad & \mathcal{G}_k \cap \mathcal{G}_{k' \neq k} = \emptyset, \bigcup_{k=1}^{\widehat{K}} \mathcal{G}_k = [1, \dots, N], \end{aligned} \quad (13)$$

where λ is the trade-off, $\mathcal{G} = \{\mathcal{G}_1, \dots, \mathcal{G}_{\widehat{K}}\}$ denotes the set of the result groups, \widehat{K} is the estimated number of clusters and $[1, \dots, N]$ denotes the universal ordinal set of samples.

We illustrate the proposed autoSC in Algorithm 1. In this section, we interpret the implementation details of the optimization in three steps including: 1) generating the triplet

relationships T from the coefficient matrix C ; 2) estimating the number of clusters \hat{K} and initializing the clusters \mathcal{C} ; 3) assigning the samples $\mathbf{x} \in \mathbf{X}_{out}$ into proper cluster.

The Generation of Triplets The coefficient matrix C reflects the correlations among samples (Elhamifar and Vidal 2009). Larger value indicates stronger belief for the connection between samples. For instance, $c_{ij} > c_{ik}$ indicates a larger probability for \mathbf{x}_i and \mathbf{x}_j being in the same cluster over \mathbf{x}_i and \mathbf{x}_k . Accordingly, we explore the intrinsic local correlations among samples by the proposed triplets derived from C .

Many subspace representations guarantee the mapping invariance via a dense coefficient matrix C . However, the generation of triplets relies only on the strongest connections to avoid the wrong assignment. Therefore, for each column of C , *i.e.*, c_i , we preserve only the top m values which are then modified to 1 for a new binary coefficient matrix C^* .

Then, we extract each triplet from C^* by the following function:

$$\begin{aligned} \tau &= \{\mathbf{x}_{n_1}, \mathbf{x}_{n_2}, \mathbf{x}_{n_3}\} \in T, \\ \text{if and only if: } &c_{n_1 n_2}^* \times c_{n_2 n_3}^* \times c_{n_3 n_1}^* = 1, \end{aligned} \quad (14)$$

where c_{xy}^* denotes the xy -th value of C^* . Note each sample \mathbf{x} can belong to many triplets. Therefore, we consider each τ as a meta-element in the clustering, which improves the robustness due to the complementarity constraints.

The Initialization of the Clusters In the I -th iteration, we first determine the initialized triplet (termed as τ_{ini}^I) from T_{out} to be the basement of the cluster \mathcal{C} . Then, we merge the most correlated samples of τ_{ini}^I into \mathcal{C} . Finally, we detect the repetitiveness between \mathcal{C} and the ‘‘in-cluster’’ samples \mathbf{X}_{in}^I to avoid the redundancy.

Following (Rodriguez and Laio 2014), we initialize a new cluster from τ_{ini}^I with highest local density:

$$\tau_{ini}^I = \arg \max_{\tau} \rho(\tau, \mathbf{X}_{out}^I), \quad (15)$$

where ρ calculate the local density as shown in Definition 5. The high local density of the triplet evokes the most connections between τ_i and other triplets, which induces the most connections between $\mathbf{x}_{n_j}^i$ and other samples in \mathbf{X}_{out}^I .

Once the initialized triplet τ_{ini}^I is determined, we iteratively extend the cluster \mathcal{C} by fusing the most confident triplets. For each triplet τ_i in T_{out} , we calculate the sum of the connection score regarding the samples in \mathcal{C} to greedily determine whether the samples in τ_i should be assigned into \mathcal{C} :

$$\begin{aligned} \tau^* &= \arg \max_{\tau} \sum_{j=1}^3 \sum_{\kappa}^{|m|} s_{n_j m_{\kappa}}, \\ \text{s.t. } &\sum_{j=1}^3 \sum_{\kappa}^{|m|} s_{n_j m_{\kappa}} > 1; \{\mathbf{x}_{n_j}\}_{j=1}^3 \in \tau; \{\mathbf{x}_{m_{\kappa}}\}_{\kappa=1}^{|m|} \in \mathcal{C}, \end{aligned} \quad (16)$$

where \mathbf{n}, \mathbf{m} denote the set of ordinals for the samples in τ and \mathcal{C} , respectively. We iteratively update the auxiliary sets $T_{out}^I, T_{in}^I, \mathbf{X}_{out}^I$ and \mathbf{X}_{in}^I in the iteration.

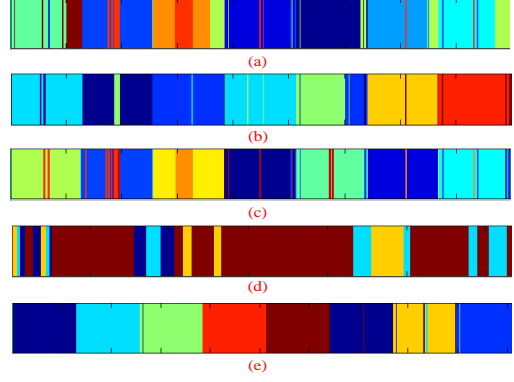


Figure 1: The visualization of the clustering labels for four methods, *i.e.*, (a) SCAMS, (b) DP, (c) SVD, (d) DP-space and (e) the proposed autoSC. The experiments are executed on the extended Yale B dataset with 8 subjects. As shown, SCAMS over-segments the samples, while DP-space assigns the majority into one cluster. The proposed autoSC does not suffer from these problems.

We terminate the process of model selection and get \tilde{K} clusters if and only if $\tau_{ini}^{\tilde{K}+1}$ satisfies:

$$\rho(\tau_{ini}^{\tilde{K}+1}, \mathbf{X}_{out}^{\tilde{K}+1}) \leq \rho(\tau_{ini}^{\tilde{K}+1}, \mathbf{X}_{in}^{\tilde{K}+1}). \quad (17)$$

Specifically, if the samples of $\tau_{ini}^{\tilde{K}+1}$ are of high frequency in $\mathbf{X}_{in}^{\tilde{K}+1}$, *i.e.*, the triplet with the highest local density in $T_{out}^{\tilde{K}+1}$ is already contained in $T_{in}^{\tilde{K}+1}$, we consider the clusters are sufficient for modeling the intrinsic sub-manifolds.

We also introduce an alternative step to check the redundancy among clusters $\{\mathcal{C}_i\}_{i=1}^{\tilde{K}}$. We calculate the connection scores s for small-scale clusters against others, and merge the highly correlated clusters \mathcal{C}_i and \mathcal{C}_j if we have

$$s_{ij} > \min(|\mathcal{C}_i|, |\mathcal{C}_j|), \quad (18)$$

where $|\mathcal{C}|$ denotes the number of samples in \mathcal{C} . We then get the initialized clusters $\{\mathcal{C}_i\}_{i=1}^{\hat{K}}$, where \hat{K} is the estimated number of clusters and $\hat{K} \leq \tilde{K}$.

The Assignment of the Remaining Samples In this stage, we assign each of the remaining samples into clusters which evokes an optimal fusion reward. For \mathbf{x}_j , we find its optimal cluster \mathcal{C}^* by the following equation:

$$\mathcal{C}^* = \arg \max_{\mathcal{C}_i} R_f(\mathcal{C}_i | \mathbf{x}_j), \quad i \in \{1, 2, \dots, \hat{K}\}, \quad (19)$$

where $R_f(\mathcal{C} | \mathbf{x})$ is the fusion reward defined by (10).

Experiment

Setup

In the experiments, we first compare the proposed autoSC with various automatic methods on two benchmark datasets, *i.e.*, the extended Yale B and the COIL-20 dataset. Then, we verify the robustness of the proposed method with combinations to different \mathcal{C} derived from various subspace representations. We also design comprehensive evaluation metrics to

Table 1: Overall Comparison between autoSC and other algorithms on subsets of the extended Yale B and COIL-20 dataset. The coefficient matrix C derived from SMR is utilized as the correlation matrix of DP, SVD and the proposed autoSC. As shown, autoSC achieves the state-of-the-art performance on all reported configurations.

Methods	Metrics	extended Yale B			COIL-20		
		8	15	30	5	10	15
SCAMS	NC_e	9.26	23.60	76.22	8.48	19.72	32.40
	NMI	0.7183	0.7272	0.7266	0.5885	0.6527	0.6668
DP	NC_e	3.06	7.84	24.76	2.22	5.30	9.72
	NMI	0.6196	0.5026	0.2166	0.6864	0.4467	0.3643
SVD	NC_e	2.40	9.06	24.00	0.48	2.58	8.36
	NMI	0.7078	0.4993	0.2808	0.7024	0.7127	0.7224
DP-space	NC_e	2.08	8.96	23.92	0.78	4.78	9.38
	NMI	0.0343	0.0226	0.0406	0.0904	0.0829	0.0718
autoSC	NC_e	0.76	2.08	4.98	0.38	1.18	0.80
	NMI	0.9062	0.8589	0.8287	0.8315	0.7701	0.7266

validate the clustering performance, *e.g.*, the error rate of the number of clusters and the triplets, *etc.* All reported results are the average of 50 trials.

Methods We make comparisons with the following methods: SCAMS (Li et al. 2016), density peak based method (DP) (Rodriguez and Laio 2014), singular value decomposition based method (SVD) (Liu et al. 2013) and DP-space (Wang and Zhu 2015). Besides, we utilize the following subspace representation methods to generate different coefficient matrix C , including LRR (Liu, Latecki, and Yan 2010), CASS (Lu et al. 2015), LSR (Lu et al. 2012), smooth representation (SMR) (Hu et al. 2014) and ORGEN (You et al. 2016). The coefficient matrix C is then used to calculate the triplet relationships for the proposed autoSC.

Metrics To evaluate the performance of the proposed triplets, we define the error rate \mathcal{A} as follows:

$$\mathcal{A} = \frac{1}{n} \sum_{i=1}^n \frac{3 - f(\tau_i | \mathbf{g}_i^*)}{2}, \quad (20)$$

where n denotes the number of the triplets and $f(\tau | \mathbf{g}^*)$ is the counting function on the frequency that $\mathbf{x} \in \mathbf{g}^*$ for all $\mathbf{x} \in \tau$. Here the output of f ranges from 0 to 2. The dynamic set \mathbf{g}_i^* consists of samples in one subspace S according to the ground truth, where S contains as many samples in τ_i as possible.

We introduce the error rate of the number of clusters (NC_e) as the primary evaluation metric for the clustering methods which estimate the number of clusters \hat{K} automatically:

$$NC_e = \frac{1}{M} \sum_{i=1}^M |\hat{K}_i - K|, \quad (21)$$

where K is the real number of clusters, M is the number of trials and \hat{K}_i is the estimated number of clusters in the i -th trial. We also use the standard normalized mutual information (NMI) (Li et al. 2003) to measure the similarity between two clustering distributions, *i.e.*, the prediction and the ground truth. With respect to NMI, the entropy illustrates the non-determinacy of one clustering to the other, and the mutual

Table 2: The evaluation results of \mathcal{A} for the proposed autoSC on the extended Yale B (eYaleB) and COIL-20 dataset. For each column, we utilize the coefficient matrix of a specific subspace representation module. As shown, the autoSC achieves consistency on the calculation of the triplets, which guarantees the performance of the model selection and clustering process.

Datasets	LRR	CASS	LSR	SMR	ORGEN	
eYaleB	8	0.0155	0.0158	0.0144	0.0135	0.0166
	15	0.0147	0.0148	0.0148	0.0149	0.0145
	30	0.0176	0.0157	0.0181	0.0181	0.0177
COIL20	5	0.0185	0.0195	0.0188	0.0175	0.0196
	10	0.0252	0.0198	0.0188	0.0182	0.0210
	15	0.0224	0.0203	0.0212	0.0196	0.0215

information quantifies the amount of information that one variable obtains from the other.

The Comparisons among Automatic Clustering

We conduct experiments on two datasets (the extended Yale B and the COIL-20) with different number of subjects, and compare four methods with the proposed autoSC on the metrics of NC_e and NMI. For DP, SVD and our autoSC, the optimization module in SMR (Hu et al. 2014) is employed to generate the coefficient matrix C . The DP-space method simultaneously estimates the \hat{K} and finds the subspaces without the requirement of coefficient matrix. All parameters of the contrasted methods are tuned to be the best. Table 1 and Figure 1 show the performance.

As shown, the averaged NC_e of autoSC is relatively smaller than others on all experimental configurations, meaning that we can give a close estimation on the number of clusters. For example, the estimated \hat{K} on the extended Yale B with 8 subjects has the deviation which is less than 1, and evokes the NMI higher than 0.9. SVD achieves comparable result on the small-scale configuration of each dataset, but the performance turns to be bad when the number of samples increases. It is mainly because that the largest gap between

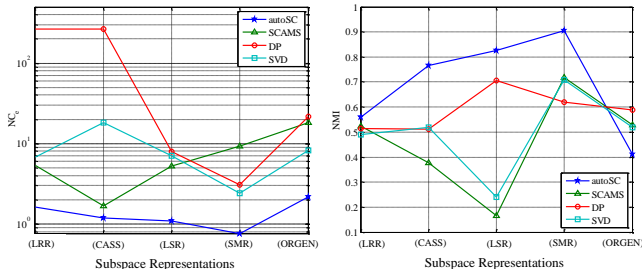


Figure 2: The extension to different subspace representations on the extended Yale B dataset with 8 subjects. The left figure denotes the comparison of four methods using the NC_e metric while the right one is regarding the NMI. Each point in the curves is derived by the combination of clustering method and subspace representation. As shown, the proposed autoSC achieves consistent performance on the evaluation of NC_e .

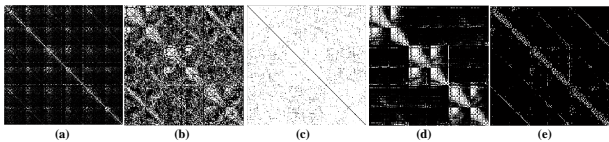


Figure 3: The coefficient matrix C derived from different methods on the extended Yale B dataset with 3 subjects. The white regions denote the locations with non-zero coefficients. Different methods, *i.e.*, (a) LRR, (b) CASS, (c) LSR, (d) SMR and (e) ORGEN, induce C with different characteristics, *e.g.*, (d) derived from SMR is block-diagonal yet (e) derived from ORGEN is sparse.

the pair of singulars decreases when the number of clusters getting larger. SCAMS performs comparably for the NMI on both datasets, however, as is illustrated in Figure 1 (a) and Table 1, it provides a much larger \hat{K} than the ground truth, *e.g.*, $\hat{K} > 100$ when $K = 30$ on the extended Yale B dataset. The NMI is friendly to the situation of over-segmentation, making the metric NC_e be the primary evaluation of the SCAMS method. The DP-space performs well on NC_e , but has bad performance on the NMI. It is because that most samples are assigned into one cluster, and the other clusters are of small scale. As a contrast, our autoSC achieves favorable results with a satisfied generalization ability.

The Robustness to Subspace Representations

As illustrated, the methods including SCAMS (Li et al. 2016), SVD (Liu et al. 2013) and our autoSC require the coefficient matrix C as input. Also, for DP (Rodriguez and Laio 2014), it needs to calculate the distance among samples. We calculate the distance d_{ij} between samples x_i and x_j by $d_{ij} = \frac{1}{c_{ij}}$ rather than the simple Euclidean distance. To verify the robustness of the proposed autoSC regarding various subspace representations, we calculate the coefficient matrix C using 5 subspace representation modules, followed by the combinations with the 4 methods which automatically estimate the number of clusters and segment the samples.

Table 2 shows the evaluation results of \mathcal{A} on both datasets

with the combinations of 5 subspace representations. The experimental results on the extended Yale B dataset with 8 subjects are reported in Figure 2, while other similar experiments with different configurations can be found in the supplementary material. Moreover, we visualize the coefficient matrix C derived from 5 subspace representation modules in Figure 3. We can see from Figure 2 that the SCAMS, DP and SVD methods are sensitive on the choice of the subspace representation module. For example, DP estimates the \hat{K} as a relatively close value to the ground truth when combined with SMR ($NC_e = 3.06$), but generates a totally wrong estimation when combined with LRR ($NC_e = 265.60$). Different subspace representation modules generate coefficient matrices with various intrinsic properties (Vidal 2011), thus the parameter for truncation error ϵ needs to be tuned carefully.

For the proposed autoSC, it is stable on different combinations considering the metric of NC_e and \mathcal{A} , which demonstrates the complementary ability of the proposed method. For all combinations, the error rate of the triplets obtained from (14) is less than 2%, which guarantees the consistency of the proposed autoSC with different kind of C . Furthermore, it shows better performance when combined with CASS, LSR and SMR than other combinations on both metrics in Figure 2. The reason lies on the guarantee of the mapping invariance which is termed as the grouping effect (Lu et al. 2015; Lu et al. 2012; Hu et al. 2014), together with the filtering of weak connections and the self-constraint among samples within triplets. As shown in Figure 3 (b), (c), (d), the coefficient matrices are dense while it shows block-diagonal structure in Figure 3 (d) and each block corresponds to one cluster. Therefore, the nearest neighbors which are used to generate the triplets can be chosen precisely. As shown in Figure 3 (b), (c), (d), the coefficient matrices are dense and block-diagonal which evokes effective triplets. The autoSC can not achieve satisfied performance when combined with the ORGEN as shown in Figure 2. It is because that the C derived from ORGEN is too sparse which contains insufficient localities for constructing effective triplets.

Conclusion

In this paper, we propose the autoSC method to simultaneously estimate the number of clusters and segment the samples according to the coefficient matrix derived by off-the-shelf subspace representations. We consider this problem as an optimization process on two reward functions at the same time, of which the model selection reward constrains the number of clusters and the fusion reward facilitates the segmentation of the samples. These two functions are greedily maximized during the clustering process, utilizing the predefined data structure termed triplet relationship. The triplet is more robust than the pairwise relationships when partitioning the inter-cluster samples near the intersection of two sub-manifolds due to the complementarity of multiple constraints. Besides, it evokes mutual restrictions of neighbors thus depicts a local geometrical structure, by which we can calculate the segmentation greedily. Extensive experiments on the benchmarks demonstrate the effectiveness of our approach.

Acknowledgments

This research was sponsored by NSFC (61620106008, 61572264), CAST (YESS20150117), Huawei Innovation Research Program (HIRP), and IBM Global SUR award.

References

- [Beier, Hamprecht, and Kappes 2015] Beier, T.; Hamprecht, F. A.; and Kappes, J. H. 2015. Fusion moves for correlation clustering. In *CVPR*.
- [Belkin and Niyogi 2001] Belkin, M., and Niyogi, P. 2001. Laplacian eigenmaps and spectral techniques for embedding and clustering. In *NIPS*.
- [Cheng et al. 2016] Cheng, Y.; Wang, Y.; Sznaiar, M.; and Camps, O. 2016. Subspace clustering with priors via sparse quadratically constrained quadratic programming. In *CVPR*.
- [Elhamifar and Vidal 2009] Elhamifar, E., and Vidal, R. 2009. Sparse subspace clustering. In *CVPR*.
- [Elhamifar and Vidal 2013] Elhamifar, E., and Vidal, R. 2013. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [Favaro, Vidal, and Ravichandran 2011] Favaro, P.; Vidal, R.; and Ravichandran, A. 2011. A closed form solution to robust subspace estimation and clustering. In *CVPR*.
- [Feng et al. 2014] Feng, J.; Lin, Z.; Xu, H.; and Yan, S. 2014. Robust subspace segmentation with block-diagonal prior. In *CVPR*.
- [Gao, Tsang, and Chia 2013] Gao, S.; Tsang, I. W.; and Chia, L. T. 2013. Laplacian sparse coding, hypergraph laplacian sparse coding, and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [Guo, Gao, and Li 2013] Guo, Y.; Gao, J.; and Li, F. 2013. Spatial subspace clustering for hyperspectral data segmentation. In *SDIWC*.
- [Hu et al. 2014] Hu, H.; Lin, Z.; Feng, J.; and Zhou, J. 2014. Smooth representation clustering. In *CVPR*.
- [Kim et al. 2014] Kim, S.; Chang, D. Y.; Nowozin, S.; and Kohli, P. 2014. Image segmentation using higher-order correlation clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [Li and Vidal 2015] Li, C. G., and Vidal, R. 2015. Structured sparse subspace clustering: A unified optimization framework. In *CVPR*.
- [Li et al. 2003] Li, M.; Chen, X.; Li, X.; and Ma, B. 2003. Clustering by compression. In *IEEE International Symposium on Information Theory*.
- [Li et al. 2016] Li, Z.; Yang, S.; Cheong, L. F.; and Toh, K. C. 2016. Simultaneous clustering and model selection for tensor affinities. In *CVPR*.
- [Li, Cheong, and Zhou 2014] Li, Z.; Cheong, L. F.; and Zhou, S. Z. 2014. SCAMS: Simultaneous clustering and model selection. In *CVPR*.
- [Liu et al. 2013] Liu, G.; Lin, Z.; Yan, S.; Sun, J.; Yu, Y.; and Ma, Y. 2013. Robust recovery of subspace structures by low-rank representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [Liu, Latecki, and Yan 2010] Liu, H.; Latecki, L.; and Yan, S. 2010. Robust clustering as ensembles of affinity relations. In *NIPS*.
- [Lu et al. 2012] 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 *ECCV*.
- [Lu et al. 2015] Lu, C.; Feng, J.; Lin, Z.; and Yan, S. 2015. Correlation adaptive subspace segmentation by trace lasso. In *CVPR*.
- [Lu et al. 2016] Lu, C.; Feng, J.; Chen, Y.; Liu, W.; Lin, Z.; and Yan, S. 2016. Tensor robust principal component analysis: Exact recovery of corrupted low-rank tensors via convex optimization. In *CVPR*.
- [Nasihatkon and Hartley 2011] Nasihatkon, B., and Hartley, R. 2011. Graph connectivity in sparse subspace clustering. In *CVPR*.
- [Peng et al. 2016] Peng, X.; Yu, Z.; Yi, Z.; and Tang, H. 2016. Constructing the l2-graph for robust subspace learning and subspace clustering. *IEEE Transactions on Cybernetics*.
- [Peng, Zhang, and Yi 2013] Peng, X.; Zhang, L.; and Yi, Z. 2013. Scalable sparse subspace clustering. In *CVPR*.
- [Purkait et al. 2014] Purkait, P.; Chin, T. J.; Ackermann, H.; and Suter, D. 2014. Clustering with hypergraphs: The case for large hyperedges. In *ECCV*.
- [Rodriguez and Laio 2014] Rodriguez, A., and Laio, A. 2014. Clustering by fast search and find of density peaks. *Science*.
- [Schölkopf, Platt, and Hofmann 2006] Schölkopf, B.; Platt, J.; and Hofmann, T. 2006. Learning with hypergraphs: Clustering, classification, and embedding. In *NIPS*.
- [Schroff, Kalenichenko, and Philbin 2015] Schroff, F.; Kalenichenko, D.; and Philbin, J. 2015. Facenet: A unified embedding for face recognition and clustering. In *ICCV*.
- [Shi and Malik 2000] Shi, J., and Malik, J. 2000. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [Vidal 2011] Vidal, R. 2011. Subspace clustering. *IEEE Signal Processing Magazine*.
- [Wang and Zhu 2015] Wang, Y., and Zhu, J. 2015. Dp-space: Bayesian nonparametric subspace clustering with small-variance asymptotics. In *ICML*.
- [Yang et al. 2016] Yang, Y.; Feng, J.; Jovic, N.; Yang, J.; and Huang, T. S. 2016. l_0 -sparse subspace clustering. In *ECCV*.
- [Yin, Gao, and Lin 2016] Yin, M.; Gao, J.; and Lin, Z. 2016. Laplacian regularized low-rank representation and its applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [You et al. 2016] You, C.; Li, C. G.; Robinson, D. P.; and Vidal, R. 2016. Oracle based active set algorithm for scalable elastic net subspace clustering. In *CVPR*.
- [Zhao et al. 2017] Zhao, S.; Yao, H.; Gao, Y.; Ji, R.; and Ding, G. 2017. Continuous probability distribution prediction of image emotions via multitask shared sparse regression. *IEEE Transactions on Multimedia*.