# SA-Net: A deep spectral analysis network for image clustering

Jinghua Wang and Jianmin Jiang*

*Research Institute for Future Media Computing, College of Computer Science & Software Engineering, and Guangdong Laboratory of Artificial Intelligence & Digital Economy (SZ), Shenzhen University, China*
*wangjh2012@foxmail.com; jianmin.jiang@szu.edu.cn*

**Abstract**

Although supervised deep representation learning has attracted enormous attentions across areas of pattern recognition and computer vision, little progress has been made towards unsupervised deep representation learning for image clustering. In this paper, we propose a deep spectral analysis network for unsupervised representation learning and image clustering. While spectral analysis is established with solid theoretical foundations and has been widely applied to unsupervised data mining, its essential weakness lies in the fact that it is difficult to construct a proper affinity matrix and determine the involving Laplacian matrix for a given dataset. In this paper, we propose a SA-Net to overcome these weaknesses and achieve improved image clustering by extending the spectral analysis procedure into a deep learning framework with multiple layers. The SA-Net has the capability to learn deep representations and reveal deep correlations among data samples. Compared with the existing spectral analysis, the SA-Net achieves two advantages: (i) Given the fact that one spectral analysis procedure can only deal with one subset of the given dataset, our proposed SA-Net elegantly integrates multiple parallel and consecutive spectral analysis procedures together to enable interactive learning across different units towards a coordinated clustering model; (ii) Our SA-Net can identify the local similarities among different images at patch level and hence achieves a higher level of robustness against occlusions. Extensive experiments on a number of popular datasets support that our proposed SA-Net outperforms

---

*Corresponding author

11 benchmarks across a number of image clustering applications.

## 1. Introduction

As clustering is one of the most fundamental tasks in machine learning and data mining [1, 2, 3], its main goal is to reveal the meaningful structure of a dataset by categorizing the data samples into a number of clusters, where similar samples are grouped together. It is extensively studied and a large number of methods have been reported, including subspace clustering [4], partitional clustering [5], hierarchical clustering [6, 7], and density-based clustering [8]. Over the past decades, clustering has found a wide range of applications, such as video retrieval [9], text analysis [10], as well as large scale data analysis [11, 12].

Spectral analysis is one of the most promising clustering methods [13, 14], and has been successfully applied in various computer vision tasks [15, 16]. Spectral analysis first derives a Laplacian matrix from the pairwise similarities among the data samples, and then embeds the data samples into an eigenspace of the Laplacian matrix, before the $k$-means is applied to complete the final categorization of all the data samples. Theoretically, the numerical embeddings or the spectral features of the data samples are taken as the relaxation of binary cluster labels [17]. Thus, these spectral features can improve both the intra-cluster compactness and the inter-cluster separability. Spectral analysis has three appealing properties, including (i) it can produce the embeddings analytically via an eigen-decomposition procedure; (ii) it has solid interpretations and can be derived from the theory of random walk, where the diffusion distance between a pair of data samples is equal to the distance between their embeddings [18]; (iii) spectral analysis is effective in revealing the non-convex data structure [19].

In general, spectral analysis has two unsolved problems. The first problem is the fact that it is still unclear how the affinity graph can influence the clustering performance. To construct the affinity graph, there exist three popular strategies, including, $k$-nearest-neighborhood, $\epsilon$-nearest-neighborhood, and fully connected graph. While each of these three strategies has its own advantages and disadvantages, how to choose

a specific strategy and how to determine its optimal parameter still remains to be an open issue. The second problem is that no agreement has ever been reached in the choice of Laplacian matrix for eigenvector decomposition. Both of the two popular Laplacian matrices, i.e. symmetric normalized Laplacian matrix [13] and left normalized Laplacian matrix [20, 21], have their own advantages and disadvantages.

Being the input of clustering, representations of data samples are also important for achieving good performances. Out of the popularity of deep learning, an increasing number of researchers use convolutional neural network (CNN) to learn deep representations that are feasible for clustering [22, 23, 24, 25, 26, 27, 28, 29]. Compared with low-level or handcrafted representations, the deep representations show overwhelming strengths in dealing with complicated data sample distributions [25, 30, 31].

Motivated by the significant success of deep learning, we extend the spectral analysis into multiple layers and propose a new spectral analysis network (SA-Net). Our SA-Net learns deep representations based on both parallel and consecutive spectral analysis procedures and shows its strength in various image clustering tasks. The proposed SA-Net consists of four different types of layers, i.e. spectral analysis layer, pooling layer, binarization layer, and coding layer. While parallel spectral analysis procedures reveal the intrinsic structure of differently distributed data samples, the consecutive spectral analysis procedures inside the SA-Net learn deep features to further improve the clustering friendliness of spectral features. By taking the image patches as the input, a spectral analysis layer learns a patch-level representation space in such a way that similar patches are made close to each other and dissimilar ones are made far away from each other. This procedure can implicitly associate similar patches across different images and identify the local similarity among them. While the spectral analysis layer stacks the patch representations to produce the representation of an image for further processing by other layers, the pooling layer reduces the size of the feature by summarization, the binarization layer binarizes the spectral features, and the coding layer transforms the binary features into numerical feature maps.

Compared with the existing spectral analysis techniques applied to clustering, our SA-Net has the following three advantages.

- While the existing approaches are dominated by one single spectral analysis procedure to learn clustering-friendly features, our SA-Net stacks multiple spectral analysis procedures in both parallel and consecutive manners to identify the best possible features for data samples across various distribution models.

- Our proposed SA-Net elegantly integrates three types of affinity graphs as well as two different normalized Laplacian matrices rather than relying on a single empirically determined spectral analysis procedure. In this way, different spectral analysis procedures can collaborate together in dealing with different data sample distributions. Thus, our network can achieve enhanced clustering performances in dealing with the variety of input datasets.

- While existing spectral analysis can only assess the similarity between image pairs holistically, our proposed SA-Net can reveal the local similarity at patch level via learning with multiple and multi-type layers. As a result, the proposed method is more robust in identifying local similarities among images, especially against occlusions.

The rest of this paper is organized as follows. Section 2 reviews the existing spectral analysis clustering and deep representation learning, which are related to our work. Section 3 presents the details of the proposed spectral analysis network (SA-Net). Section 4 reports the experiments and finally Section 5 provides concluding remarks.

## 2. Related Work

### 2.1. Spectral Analysis Clustering

Given a dataset of $N$ samples, i.e. $I = \{I_1, I_2, \cdots, I_N\}$, a clustering task aims to partition it into $k$ clusters. To achieve this, spectral analysis methods first build a non-directed graph $G = \{I, W\}$, with $W \in R^{N \times N}$ as the affinity matrix. In the graph $G$, each node $I_i (1 \leq i \leq N)$ corresponds to a data sample and the element $w_{ij} (1 \leq i, j \leq N)$ represents the affinity between a pair of nodes $I_i$ and $I_j$. Spectral analysis partitions the graph $G$ into a number of subgraphs based on a graph cut criterion [20], and thus

4

produces a set of clusters. Mathematically, spectral clustering solves the following minimization problem:

$$\min_{\hat{Y}} tr(\hat{Y}^T W \hat{Y}) \tag{1}$$

where the binary assignment matrix $\hat{Y} \in \{0, 1\}^{N \times k}$ satisfies $\hat{Y} \mathbb{1}_k = \mathbb{1}_N$, i.e. each sample belongs to one and only one cluster. The element $\hat{y}_{ic} = 1$ if and only if the data sample $I_i$ is assigned to the $c$th cluster.

It is known that the problem in Eq. (1) is NP-hard. In order to solve this problem numerically, researchers relax it by the spectral graph theory [32]. By allowing the assignment matrix to have continuous values, we obtain the following relaxation for Eq. (1):

$$\min_{Y} tr(Y^T W Y) \quad s.t. \quad Y^T Y = E_k \tag{2}$$

where $Y \in R^{N \times k}$ is the relaxed continuous clustering assignment matrix with orthogonal constraint and $E_k \in R^{k \times k}$ is the identity matrix. Based on a normalized cut criterion, the spectral clustering can also be formulated as [17]:

$$\max_{Y^T D Y = E_k} tr(Y^T W Y) \tag{3}$$

where $Y = \hat{Y}(\hat{Y}^T D \hat{Y})^{-1/2}$ and $D \in R^{N \times N}$ is a diagonal matrix with $D_{ii} = \sum_{j=1}^{N} w_{ij}$. We can also formulate spectral clustering as [13]

$$\min_{Y^T D Y = E_k} tr(Y^T L Y) \tag{4}$$

where $L = D - W$ is the Laplacian matrix. With the solution $Y$ of Eq. (2), (3), or (4), we can obtain the final cluster labels by simply conducting an additional $k$-means procedure. For the convenience of description, we consider $y_i$ as the spectral feature of the data sample $I_i (1 \leq i \leq N)$.

A typical spectral analysis clustering procedure is shown in Algorithm 1, which mainly consists of four steps, i.e. affinity matrix construction, Laplacian matrix computation, matrix eigen-decomposition, and $k$-means clustering. While spectral analysis clustering has the advantage in revealing the intrinsic data distribution [13, 21], it also has a number of deficiencies, compared with other clustering approaches, which can be highlighted as follows.

5

**Algorithm 1** Spectral Clustering

**Input:** data points $I = \{I_1, I_2, \cdots, I_N\}$ and number of clusters $k$;

**Output:** $k$ clusters

1: Construct an affinity matrix $W \in R^{N \times N}$ between data points, where $w_{ij}$ measures the similarity between $I_i$ and $I_j$;

2: Compute the Laplacian matrix $L = D - W$, where $D \in R^{N \times N}$ is the degree matrix with $d_{ii} = \sum_{j=1}^{N} w_{ij}$;

3: Compute the $k$ eigenvectors $q_i (1 \leq i \leq k)$ for Laplacian matrix $L$, corresponding to the $k$ smallest eigenvalues, and denote them by: $Q = [q_i, q_2, \cdots, q_k] \in R^{N \times k}$;

4: For $1 \leq i \leq N$, let $y_i \in R^k$ be the $i$th row of the matrix $Q$, and apply $k$-means to cluster the points $y_i (1 \leq i \leq N)$ to obtain the $k$ clusters: $Cluster_j (1 \leq j \leq k)$.

Firstly, there is little theoretical analysis that could lead us to a proper affinity matrix $W$ for a given dataset, although it is extensively studied [33, 34]. While three different similarity measurements are popularly used to construct the affinity matrix, including $k$-nearest-neighborhood, $\epsilon$-nearest-neighborhood, and the fully connected graph, each of these three affinity matrices can only deal with some but not all types of data sets. While the $k$-nearest-neighborhood strategy might break a connected component into several components, the $\epsilon$-nearest-neighborhood strategy can not handle datasets with varying densities, and the fully connected affinity method suffers from high computational complexity. In addition, the clustering results are also sensitive to the parameter variations of these similarity measurements.

Secondly, the relevant research communities have not reached consensus on how to choose between different Laplacian matrices. The unnormalized Laplacian matrix $L$ has two popular extensions, including symmetric normalization and left normalization, and details of these two normalization matrices are described in the following two equations:

$$L_{sym} = D^{-1/2} L D^{-1/2} = E_N - D^{-1/2} W D^{-1/2} \tag{5}$$

$$L_{rw} = D^{-1} L = E_N - D^{-1} W \tag{6}$$

6

where $E_N \in R^{N \times N}$ is the identity matrix. Both the normalized matrices are positive semi-definite and thus have $N$ real-valued eigenvalues. Let $v_{rw}$ be an eigenvector of $L_{rw}$ corresponding to eigenvalue $\lambda_{rw}$, i.e. $L_{rw}v_{rw} = \lambda_{rw}v_{rw}$. Then, we have the following equations to describe their relationships.

$$
\begin{aligned}
&E_N v_{rw} - D^{-1}W v_{rw} = \lambda_{rw}v_{rw} \Leftrightarrow \\
&v_{rw} - D^{-1/2}D^{-1/2}WD^{-1/2}D^{1/2}v_{rw} = \lambda_{rw}v_{rw} \Leftrightarrow \\
&D^{1/2}v_{rw} - (D^{-1/2}WD^{-1/2})D^{1/2}v_{rw} = \lambda_{rw}D^{1/2}v_{rw} \Leftrightarrow \\
&L_{sym}(D^{1/2}v_{rw}) = \lambda_{rw}(D^{1/2}v_{rw})
\end{aligned}
\tag{7}
$$

Thus, $D^{1/2}v_{rw}$ is the eigenvector of $L_{sym}$ corresponding to the eigenvalue $\lambda_{rw}$. This means that $L_{sym}$ and $L_{rw}$ have the same set of eigenvalues and their eigenvectors differ by a scaling of $D^{1/2}$. While Ng [13] adopted the symmetrically normalized Laplacian matrix and claimed superior performances, Shi [20] and Luxburg [21] recommended the left normalized matrix. Both normalizations have their individual advantages and disadvantages.

Thirdly, spectral clustering is computationally expensive. In general, Algorithm 1 has the computational complexity of $O(N^2)$ in terms of space and $O(N^3)$ in terms of time, and thus many efforts have been reported to reduce the computational complexity. Dhillon et al. [35] eliminate the need for eigen-decomposition by proposing a multilevel algorithm to optimize weighted graph cuts. Yan et al. [36] propose to conduct spectral clustering based on the representative centroids. In a similar manner, Zhang et al. [37] minimize the quantization error of samples and thus improve Nystrom spectral clustering. Based on the Nystrom method, Fowlkes et al. [38] approximate the affinity matrix using a subset of data samples. Wang et al. [39] select a subset of data points based on data-dependent nonuniform probability distribution and use them to construct a low-rank approximation for the affinity matrix. In 2017, Han and Filippone [40] propose to recover the Laplacian spectrum via mini-batch-based stochastic gradient optimization on Stiefel manifolds. In contrast, our proposed approach is computationally efficient and converges to critical points, even with a data set as large as $580K$, providing the potential of allowing us to conduct spectral analysis on large datasets.

## 2.2. Deep Representation Learning for Clustering

To achieve satisfactory clustering performances, data representations are also of vital importance in addition to the clustering techniques. Along with the popularity of deep learning [41], existing research has increasingly focused on deep representations, leading to significant improvement of clustering performances. Based on the reconstruction task, Hinton and Salakhutdinov propose an autoencoder to learn deep representations [42]. Tian et al. [25] correlate spectral clustering with autoencoder, and thus put forward a so-called sparse autoencoder for deep representation learning. Chen [43] takes the nonparametric maximum margin clustering results as the supervision information and learns the deep data representations using a DBN (deep belief network). As a recurrent framework for agglomerative clustering, JULE [28] integrates CNN-based representation learning with the cluster assignment learning, and through such an integration, these two learning procedures can boost each other. Both DCN [27] and DBC [23] propose a well-designed objective function in order to learn deep representations which are suitable for $k$-means clustering. DEC [26] proposes soft assignments for data samples based on the representation distribution and refines them iteratively. DEPICT [22] designs a new network structure by stacking a soft-max layer on top of a multilayer autoencoder and trains it by minimizing relative entropy. Shaham et al. [24] train a SpectralNet that can learn the data embeddings as well as the cluster assignments at the same time. Such spectralNet can deal with out-of-sample problem as well as large data set. Recently, Aljalbout et al. [30] present a systematic taxonomy for clustering with deep learning.

All of the above mentioned methods learn deep representations based on convolutional neural networks. In this work, we extend [44] and propose a new framework for deep representation learning based on spectral analysis with multiple layers.

## 3. The Proposed SA-Net

### 3.1. Main Idea

At present, the dominating technique in deep representation learning is CNN [45], consisting of multiple layers, such as convolutional layers, pooling layers, and softmax

8

layers. In a CNN, the convolution operation is the key to produce the representative and discriminative representations. Researchers also attempt to extend other techniques, such as $k$-means [46], to a network structure for deep representation learning. In this paper, we propose a new representation learning method via expansion of the concept in spectral analysis, and to the best of our knowledge, we are the first to build a deep learning network by stacking spectral analysis procedures both consecutively and parallelly.



(a) Data samples    (b) One spectral analysis procedure    (c) Two spectral analysis procedures
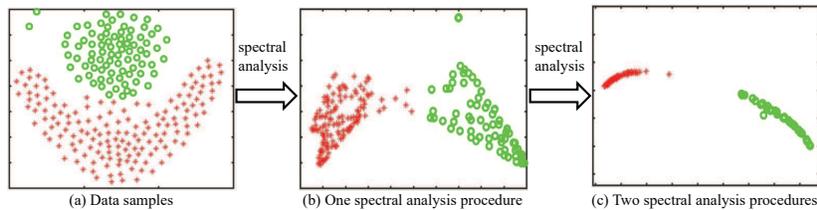
Figure 1: Illustration of samples in a two-cluster dataset and their corresponding spectral features, where (a) the data samples; (b) the shallow spectral features of the samples obtained from one spectral analysis procedure; and (c) the deep spectral features of the samples obtained from two consecutive spectral analysis procedures.

It is widely recognized that, as relaxation of the cluster assignment vectors $\hat{Y} = \{\hat{y}_i\}_{i=1}^N$, the spectral features $Y = \{y_i\}_{i=1}^N$ in algorithm 1 are more suitable for clustering than the original data points $\{I_i\}_{i=1}^N$ [21]. In other words, a spectral analysis procedure can enhance the intra-cluster similarity as well as the inter-cluster separability. Inspired by such an observation, we consider to conduct multiple spectral analysis procedures **consecutively** to learn deep representations and hence create more spaces for research towards improved clustering performances. An intuitive example is shown in Fig. 1, where part-(a) shows a two-cluster data set and part-(b) shows the spectral features obtained from a single spectral analysis procedure. By taking the features in part-(b) as the input, the second spectral analysis procedure produces a set of more clustering-friendly features (as shown in part-(c)). Specifically, the deep features in part-(c) associate with a higher Calinski-Harabasz(CH) score than the shallow features in part-(b).

As previously mentioned, there are three different methods to construct the affin-

9

ity matrix, involving two different types of Laplacian matrices. Since each of them has its own advantage, and it remains difficult to determine which one to use for a given dataset, yet the variation of a parameter can influence the clustering results significantly. Fig. 2 shows an example on two-cluster dataset and visualizes the clustering results of different spectral analysis procedures with a symmetrically normalized matrix. In Fig. 2, part-(a) and part-(b) adopt a fully connected graph (i.e. $w_{ij} = exp(-\|x_i - x_j\|^2/(2\sigma^2)))$ to construct the affinity matrix with $\sigma = 0.2$ and $\sigma = 0.5$, respectively, and part-(c) adopts knn to construct the affinity matrix with the parameter $k = 3$. As seen, all these three spectral analysis procedures do not have optimal parameters and consequently, each of them mis-cluster a portion of data points. In order to improve the clustering performance, we propose to conduct multiple spectral analysis procedures **parallelly** and integrate them together by simply concatenating their spectral features. After applying a $k$-means clustering on the concatenated features, we present our clustering results in part-(d). As seen, the integration of three spectral analysis procedures can cluster all the data points correctly, though none of them can achieve this individually. This indicates that we can improve the clustering performance by integrating multiple spectral analysis procedures with sub-optimal affinity matrices. Similarly, different Laplacian matrices can also collaborate with each other to boost the clustering performance. In this way, multiple parallel spectral analysis procedures can be elegantly integrated to work together collectively and collaboratively, and thus deep correlations across the data samples can be effectively identified and exploited for improved clustering performances.

### 3.2. SA-Net Structure

In this subsection, we propose a new network structure that can extract deep features or representations for the task of image clustering. Fig. 3 provides an overview of the proposed network structure. As seen, the proposed network consists of five layers, i.e. two spectral analysis layers, one pooling layer, one binarization layer, and one coding layer. Note that we can easily build deeper networks by adding more layers out of the four different types.

In order to enhance the clustering friendliness of its input, a spectral analysis layer

10

(a) fully connected graph, $\sigma = 0.2$      (b) fully connected graph, $\sigma = 0.5$
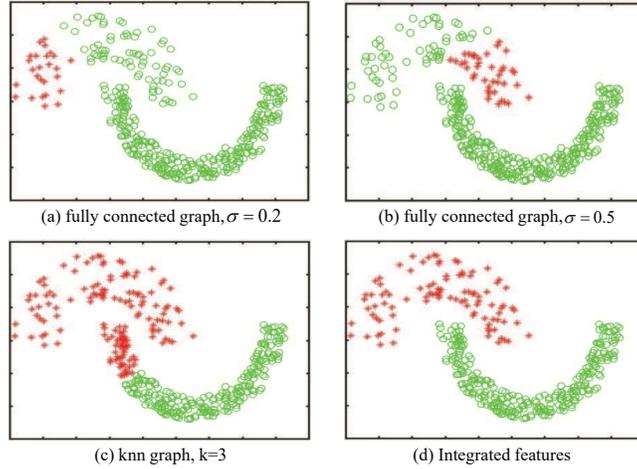
(c) knn graph, k=3      (d) Integrated features

Figure 2: Illustration of a two-cluster dataset and the clustering results by different spectral analysis procedures, where a), b), and c) use a single type of spectral features, d) concatenates all of the three spectral features before conducting $k$-means.

produces the concatenation of spectral features obtained from multiple parallel spectral analysis procedures. For the convenience of description, the output features of a spectral analysis layer is referred to as spectral features. In Algorithm 1, $y_i$ denotes the spectral feature of sample $I_i$. In general, spectral features are theoretically more intra-cluster compact and inter-cluster separate [21]. In practice, however, different spectral analysis procedures have their own capabilities in dealing with data samples with various distributions. In order to deal with data samples that follow different distributions, we extract a range of spectral features from multiple spectral analysis procedures and stack them together. The first spectral analysis layer takes the image patches as the input and then its extracted spectral features are fed into the second layer for further processing to increase the discriminative power of the first layer spectral features.

The pooling layer summarizes the neighboring spectral features, and takes the strongest response to represent the visual appearance. The binarization layer binarizes the spectral features, and the coding layer transforms the binary feature maps into numerical feature maps, which provides better suitability and more friendliness for $k$-means clustering.

11

**(a) The network structure**



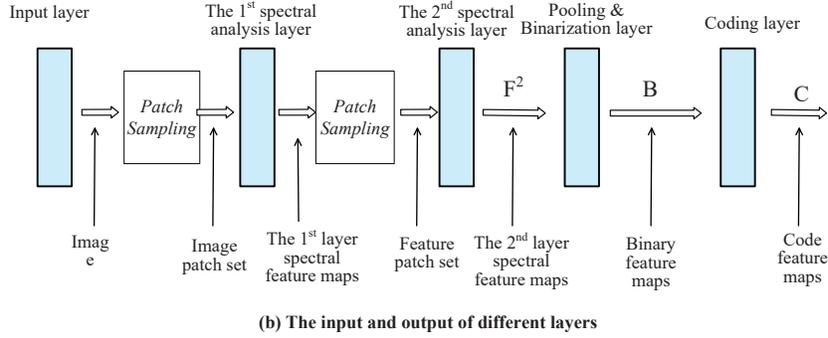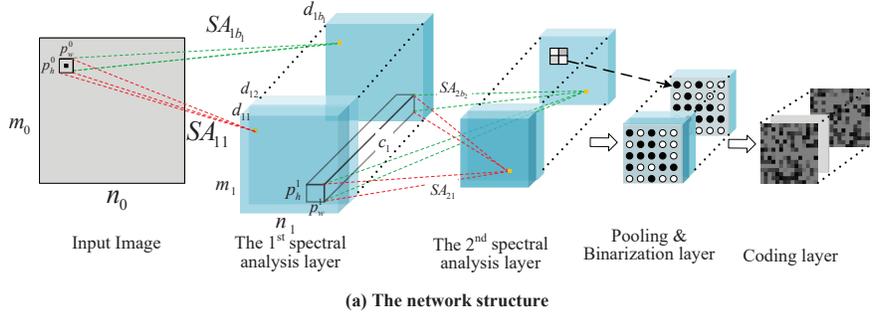**(b) The input and output of different layers**

Figure 3: Illustration of the structure of the proposed SA-Net. In addition to the input and output layers, there are five layers in the proposed network, including two spectral analysis layers, one pooling layer, one binarization layer, and one coding layer. (a) shows the details of the whole network structure and (b) shows the input and output of different layers

Detailed descriptions and discussions of all these layers are provided as follows. For an image clustering task, we assume that the image dataset $I = \{I_i | 1 \leq i \leq N\}$ consists of $N$ samples, and the image size is $m_0 \times n_0$ with $c_0$ channels, i.e. $I_i \in R^{m_0 \times n_0 \times c_0}$.

(1) The first spectral analysis layer

The first layer conducts spectral analysis based on the densely sampled image patches. In the patch sampling procedure, we pad the image to include the border information. Around each of a subset (or all) pixels, we crop an image patch of size $p_h^0 \times p_w^0 \times c_0$. With a stride of $s_0$ in patch sampling, we obtain $n_p^0 = m_1 \times n_1 =$

12

$\lceil m_0/s_0 \rceil \times \lceil n_0/s_0 \rceil$ patches in total from the image. In order to address the problem of illumination variation, a normalization procedure subtracts the mean from each image patch. For an image $I_i$, we obtain a set of normalized image patches $X_i^0 = \{x_{ij}^0 | 1 \leq j \leq n_p^0\}$. The patch set is $X^0 = \{X_1^0, X_2^0, \cdots, X_N^0\}$ with size of $N \times n_p^0$.

Given the image patch set $X^0$, the $t$th $(1 \leq t \leq b_1)$ spectral feature in the 1st layer, i.e. $SA_{1t}$, extracts the spectral feature $f_{ij}^{1t} \in R^{1 \times d_{1t}}$ for the patch $x_{ij}^0$. The parameter $b_1$ counts different spectral analysis procedures in the first layer and $d_{1t}$ denotes the dimension of the spectral feature. Taking the normalized Laplacian matrix $L_{sym}$ in Eq. (5) as an example, we can formulate the $t$th spectral analysis procedure by the following objective function:

$$\max_{F_{1t}^T D_t^1 F_{1t} = E_{d_{1t}}} tr(F_{1t}^T W_t^1 F_{1t}) \tag{8}$$

where $W_t^1$ is the affinity matrix, $D_t^1$ is a diagonal matrix with $(D_t^1)_{ii} = \sum_{j=1}^N (W_t^1)_{ij}$, and $F_{1t} = \left[ (f_{11}^{1t})^T \quad (f_{12}^{1t})^T \cdots (f_{Nn_p^0}^{1t})^T \right]^T \in R^{(N*n_p^0) \times d_{1t}}$. The following integrates the $b_1$ spectral analysis procedures into a single objective function:

$$\max_{\substack{F_{1t}^T D_t^1 F_{1t} = E_{d_{1t}} \\ t=1,2,\dots,b_1}} tr \left\{ \begin{bmatrix} F_{11} \\ F_{12} \\ \vdots \\ F_{1b_1} \end{bmatrix}^T \begin{bmatrix} W_1^1 & 0 & 0 & 0 \\ 0 & W_2^1 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & W_{b_1}^1 \end{bmatrix} \begin{bmatrix} F_{11} \\ F_{12} \\ \vdots \\ F_{1b_1} \end{bmatrix} \right\} \tag{9}$$

After obtaining $F_{1t} \in R^{(N*n_p^0) \times d_{1t}}$, we stack the spectral features belonging to the same image into $d_{1t}$ feature maps with the size of $m_1 \times n_1$, instead of clustering them directly. Note that, any one of the spectral analysis procedures in Eq. (9) can also be formulated by the left normalized Laplacian matrix $L_{rw}$. In general, the $b_1$ different spectral analysis procedures can be different from each other in one or more of the following three aspects, i.e. the affinity matrix, the Laplacian matrix, and the computing method to produce the spectral features (further details are given in Sec. 3.3). With $b_1$ different spectral analysis procedures, we obtain the spectral features of an image with the dimensionality of $m_1 \times n_1 \times c_1$, where $c_1 = \sum_{t=1}^{b_1} d_{1t}$ sums the dimensionality of $b_1$ different sets of spectral features. Let $F_i^1 \in R^{m_1 \times n_1 \times c_1}$ denote the first layer spectral features of the $i$th image, and $F^1 = \{F_i^1 | 1 \leq i \leq N\}$.

(2) The second spectral analysis layer

This layer has two operational steps. The first step is to sample the feature patches on the output of the first layer, i.e. $F^1$. Let the feature patch set be $X^1$, the second step is to conduct spectral analysis procedures on $X^1$ and produce the second layer spectral features $F^2$.

Let the size of the feature patches be $p_h^1 \times p_w^1 \times c_1$ as shown in Fig. 3, each feature patch carries the information learned from a larger patch with the size of $(p_h^1 + p_h^0 - 1) \times (p_w^1 + p_w^0 - 1) \times c_0$ in the original image. This indicates that, while the first layer deals with the correlations between small image patches, the second layer discovers the correlations among larger image areas from the original image. In addition, a feature patch also integrates the discriminative information learned by different spectral analysis procedures, which are suitable for the clustering of data samples with various distributions. This layer provides an elegant manner to integrate different spectral analysis procedures together in the feature-level, in order to enable them to work collaboratively.

(3) The pooling layer

The pooling layer summarizes the neighboring spectral features within the same spectral map, which can be conducted on the spectral features of the first or the second spectral analysis layer. To illustrate the specific operation process, we take spectral features of the second layer as an example. The spectral analysis $SC_{2t}$ produces $d_{2t}$ different feature maps with the size $m_2 \times n_2$ for each image, and each feature map is associated with an eigenvalue. The pooling operation is conducted inside each feature map. For a $s_p \times s_p$ grid centered at a point, the pooling operation only keeps the strongest response in terms of an absolute value, which can be mathematically expressed as:

$$pooling(G) = g_{kl} \qquad where \quad |g_{kl}| = \max_{ij} |g_{ij}| \tag{10}$$

where $g_{ij}$ denotes the feature in the $i$th row and the $j$th column of the target feature grid. The pooling grids can be overlapped in our algorithm.

(4) The binarization layer

From the graph cut point of view, the sign of spectral features (positive or nega-

tive) carries the cluster information [21]. For a clustering task with only two clusters, specifically, we can simply take a single eigenvector in step 3 of Algorithm 1, and cluster the data $I_i$ into the first cluster if $y_i > 0$ and the second cluster if $y_i < 0$. This observation explicitly shows the importance of the sign of the spectral features in the clustering process. Following the pooling layer, correspondingly, we use a binarization layer to binarize the spectral features, in which $B_{ij}$ denotes the $j$th binary feature maps associating with the $i$th image, and $B = \{B_{ij}\}$ denotes the binary feature map set.
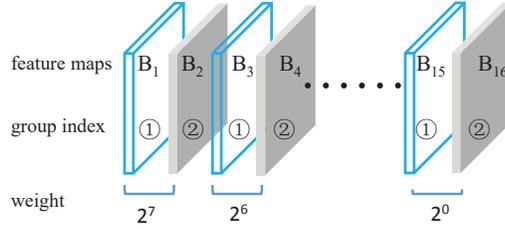


Figure 4: Illustration of partitioning the feature maps into groups, where the 16 feature maps are permuted based on their corresponding eigenvalues in a non-descending order, and we assign larger weights to those the feature maps with smaller eigenvalues.

(5) The coding layer

Following the binarization layer, a coding layer [31] is added to transform the binary code into decimal numbers and thus make it feasible for the following clustering operation. In this layer, we first partition the binary features of each image into different groups. Each group consists of $L$ binary feature maps, and $L$ is normally set to be 8. With $n_b$ binary feature maps, the coding layer produces $\lceil n_b/L \rceil$ decimal feature maps. Let $B_{ij}^k$ be the $j$th ($1 \leq j \leq L$) binary feature map in the $k$th group for the $i$th image. At the position $(u, v)$, we take the $L$ binary features $B_{ij}^k(u, v)$ and convert them into decimal by:

$$C_i^k(u, v) = \sum_{j=1}^{L} 2^{L-j} B_{ij}^k(u, v) \tag{11}$$

In this way, we obtain the $k$th decimal feature map $C_i^k$ for the $i$th image. Note, we assign large weights to the spectral features with small eigenvalues, due to their strong discriminative capabilities. By setting $L$ to be 8, we can produce the gray maps in the

coding layer as shown in Fig. 3. The final clustering results are obtained by applying a simple $k$-means to the output of the coding layer. Fig. 4 shows an example of partitioning 16 feature maps into two groups and assigning a weight to each feature map. In this figure, the feature maps are permuted based on their corresponding eigenvalues in a non-descending order, i.e. $v_i \leq v_{i+1}$, where $v_i$ is the eigenvalue associated with $B_i$.

*3.3. Operational set-up across different spectral analysis procedures*

In our proposed SA-Net, the spectral analysis procedures are different from each other in terms of affinity matrix, Laplacian matrix, and computing method. Tab. 1 presents an overview of all these methods, from which it can be seen that we have three different affinity matrices and three types of Laplacian matrices to design the spectral analysis procedures. In addition, we have four different methods to compute the spectral features based on a given Laplacian matrix. The last column of Tab. 1 provide the details of whether we adopt the corresponding item or how we set the parameters.

To construct an affinity matrix with $k$-nearest-neighborhood, we choose a parameter $k$ so that the affinity graph is connected. To determine the parameter $\epsilon$ for $\epsilon$-nearest-neighborhood affinity matrix construction, we first obtain a minimal spanning tree from the fully connected graph, then set the parameter $\epsilon$ to be $0.5\eta$, $\eta$, or $2\eta$, where $\eta$ denotes the longest edge in the minimal spanning tree. For a fully connected graph, we use empirical parameters or the self-tuning method [33] to determine an appropriate parameter value for each sample point.

There are three different types of Laplacian matrices, including one unnormalized and two normalized matrices. We adopt the two types of normalized Laplacian matrices, i.e. $L_{sym}$ and $L_{rw}$, for our proposed SA-Net, as they have already shown their advantages in clustering. Theoretically, these two normalized Laplacian matrices implement the essential objective in maximizing both the inter-cluster separability and the intra-cluster similarity. As the unnormalized matrix (i.e. $L$) only takes the first half of the objective into consideration, it is removed from our consideration.

Due to its high computational complexity, we do not apply the traditional eigen-decomposition procedure to our patch-based spectral analysis procedures. As the Lanc-

Table 1: Overview of spectral analysis procedures.

| Stage | Methods | Settings |
|---|---|---|
| | $k$-nearest-neighborhood | $k = 5, 9, 17, 21$ |
| Affinity Matrix | $\epsilon$-nearest-neighborhood | $\epsilon = 0.5\eta, \eta, 2\eta^{*}$ |
| | Fully connected matrix | $\sigma = 10^{-1}, 10^{-2}, 10^{-3}$, or self-tunning |
| | L=D-W | No$^{**}$ |
| Laplacian Matrix | $L_{sym} = E - D^{-1/2}WD^{-1/2}$ | Yes |
| | $L_{rw} = E - D^{-1}W$ | Yes |
| | Eigen-decomposition | No |
| Computing Method | Lanczos method [47] | Yes |
| | Nystrom approximation [38] | Yes |
| | Mini-batch analysis [40] | Yes |

$^{*}$ $\eta$ denotes the longest edge in the minimal spanning tree.

$^{**}$ *Yes* or *No* denotes whether the corresponding term is used in this paper.

zos method [47] has shown its advantages in decomposition of sparse matrices, we adopt this method to produce the spectral features from the sparse affinity matrices constructed by $k$-nearest-neighborhood and $\epsilon$-nearest-neighborhood. The computational complexity of Lanczos method is $O(N_{matrix} \cdot n_{eig} \cdot n_{iter})$, where $N_{matrix}$ denotes the width or height of the Laplacian matrix, $n_{eig}$ denotes the number of eigenvectors to produce, and $n_{iter}$ denotes the number of iterations. For image clustering reported in this paper, the parameter $N_{matrix}$ is a multiple of the number of images $N$. For example, we have $N_{matrix} = N \times n_p^0$ in the first spectral analysis layer, where $n_p^0$ denotes the number of patches in each image. As given in each experiment of sec. 4,

the parameter $n_{eig}$ varies in different spectral analysis procedures and is no larger than 64. We set the parameter $n_{iter}$ to be 1000 in this paper. For the dense affinity matrix, we apply the Nystrom approximation-based method [38] or mini-batch spectral clustering (MBSC) [40] to produce the spectral features. The computational complexity of the Nystrom method is $O(n_{eig}\ell_{col}^2 + n_{eig}\ell_{col}N_{matrix})$, where $\ell_{col}$ denotes the number of representative columns sampled from the Laplacian matrix. In the implementation of Nystrom method, we set $\ell_{col} = logN_{matrix}$ and adopt sparse matrix greedy approximation (SMGA) sampling method [48] to select the columns in a greedy manner. The number of eigenvectors $n_{eig}$, which equals to the target rank, is given in each experiment. Taking the face image clustering as an example, we set $n_{eig} = 64$ in the first spectral layer and $n_{eig} = 16$ in the second spectral layer. In the MBSC, we set the size of mini-batch to be $N_{matrix}^{\frac{1}{2}}$, leading to the computational complexity of $O(N_{matrix}n_{eig}^2 + N_{matrix}^2 n_{eig} + n_{eig}^3)$.

## 4. Experiments

To evaluate the proposed SA-Net, we carry out extensive experiments on a range of image clustering tasks, including handwritten digit image clustering, face image clustering, natural image clustering, and fashion image clustering. We show the robustness of our method against parameter variation in digit image clustering. We also show the performance variations of our method when the number of spectral analysis procedures changes in face image clustering. Four popular standard metrics are adopted for measuring the clustering performances, which include accuracy (ACC), nomarlized mutual information (NMI) , adjust rand index (ARI), and F1-score (FS).

To benchmark our proposed SA-Net, we compare our proposed with 11 existing clustering algorithms, which cover almost all the representative existing state of the arts in image clustering. These include: k-Means (KM), normalized cuts (N-Cuts) [20], self-tuning spectral clustering (SC-ST) [33], large-scale spectral clusteirng (SC-LS) [12], agglomerative clustering via path integral (AC-PIC) [7], spectral embedded clustering (SEC) [15], local discriminant models and global integration (LDMGI) [2], NMF with deep model (NMF-D) [3], deep embedding clustering (DEC) [26], joint

18

supervised learning (JULE) [28], and Deep embeded regularized clustering (DEPICT) [22].

### 4.1. Handwritten Digit Image Clustering

For the convenience of validation, we conduct the experiments on two popular handwritten digit image datasets, i.e. MNIST [49] and USPS[1]. The USPS dataset is a handwritten digits dataset produced by the USPS postal service. There are $11,000$ images in this dataset, each belonging to $10$ different classes (i.e. from $0$ to $9$), and the image size is $16 \times 16$. The MNIST dataset is one of the most popular image datasets, widely used for deep learning based research. In total, this dataset consists of $70,000$ images, $60,000$ for training and $10,000$ for testing. Each image in the dataset represents a handwritten digit, from $0$ to $9$. We use all the data samples in our experiments, and the images are centered with the size $28 \times 28$.

We take MNIST dataset as the example to show the implementation details, and it is similar for the USPS dataset. In the first layer, we sample $11 \times 11$ image patches with a stride of $5$ both vertically and horizontally. With padding in the border, we sample $6 \times 6 = 36$ image patches from an $28 \times 28$ image.

For $k$-nearest-neighborhood affinity matrix construction, we set the parameter $k$ to be 9 and 17. As mentioned previously, we have three different settings for the value of $\epsilon$ in $\epsilon$-nearest-neighborhood affinity matrix construction, i.e. $0.5\eta$, $\eta$, and $2\eta$, where $\eta$ denotes the longest edge in the minimal spanning tree. We also construct three different dense affinity matrices. While one dense matrix is determined by the self-tunning method [33], the other two are constructed based on the Gaussian function $w_{ij} = exp(-||x_i - x_j||^2/(2\sigma^2))$ with the parameter $\sigma$ equals to 0.1 and 0.01, respectively. Thus, we have 5 different sparse affinity matrices and 3 different dense affinity matrices.

A symmetric Laplacian matrix is computed from each of the affinity matrices to yield spectral features. We apply Lanczos to obtain spectral features from sparse affinity matrices, and mini-batch anlaysis to derive spectral features from dense affinity

---

[1]https://cs.nyu.edu/roweis/data.html

19

matrices. For each of the Laplacian matrix, we calculate $n_{eig} = 64$ eigenvectors, i.e. the dimensionality of spectral features is 64. In summary, the spectral features produced by the first layer is of the size $6 \times 6 \times 512$ for each image, with 64 dimensional spectral features for every 8 different Laplacian matrices.



Figure 5: The typical visual patterns in the MNIST dataset

To show that the first layer can learn the typical visual patterns in the image dataset, we conduct a $k$-means clustering based on the first layer spectral features and visualize 40 cluster centers in Fig. 5. As seen, while the first three rows represent lines in different angles and positions, the last two rows represent different curve shapes in the digit images. A proper combination of these visual patterns can produce a digit image.

In the second layer, we sample feature patches with the size of $4 \times 4 \times 512$. With a stride of 1 and no feature padding, we obtain $3 \times 3$ feature patches for all the $6 \times 6 \times 512$ spectral features at the first layer. In other words, each image is associated with 9 feature patches with the dimensionality of $4 \times 4 \times 512$. Similarly, the second layer also uses both sparse affinity matrices and dense affinity matrices in the spectral analysis procedures. By setting the parameter $k$ to be 9 and 17, we use $k$-nearest-neighborhood strategy to construct two sparse affinity matrices. To maximize the strength of the multiple spectral analysis procedures, we also construct two dense affinity matrices by self-tuning method and the Gaussian distance with $\sigma = 0.1$, respectively. With a symmetric Laplacian matrix employed, therefore, this layer has 4 different spectral analysis procedures altogether. As each spectral analysis produces $n_{eig} = 16$ dimensional features, the second spectral analysis layer produces 64 feature maps for every image and

Table 2: Comparative results on handwritten digit image datasets

| | MNIST | | | | USPS | | | |
|---|---|---|---|---|---|---|---|---|
| | ACC | ARI | NMI | FS | ACC | ARI | NMI | FS |
| KM | 0.534 | 0.408 | 0.500 | 0.347 | 0.460 | 0.430 | 0.451 | 0.392 |
| N-Cuts | 0.327 | 0.311 | 0.411 | 0.301 | 0.314 | 0.449 | 0.675 | 0.462 |
| SC-ST | 0.311 | 0.291 | 0.416 | 0.289 | 0.308 | 0.572 | 0.726 | 0.491 |
| SC-LS | 0.714 | 0.627 | 0.706 | 0.637 | 0.659 | 0.599 | 0.681 | 0.614 |
| SEC | 0.804 | 0.700 | 0.779 | 0.766 | 0.544 | 0.509 | 0.511 | 0.552 |
| AC-PIC | 0.115 | 0.095 | 0.017 | 0.154 | 0.855 | 0.729 | 0.840 | 0.754 |
| LDMGI | 0.842 | 0.795 | 0.802 | 0.717 | 0.580 | 0.538 | 0.563 | 0.525 |
| NMF-D | 0.175 | 0.159 | 0.152 | 0.212 | 0.382 | 0.334 | 0.287 | 0.251 |
| DEC | 0.844 | 0.795 | 0.816 | 0.716 | 0.619 | 0.554 | 0.586 | 0.635 |
| JULE | 0.959 | 0.832 | 0.906 | 0.814 | 0.922 | 0.749 | 0.858 | 0.801 |
| DEPICT | 0.965 | 0.808 | 0.917 | 0.851 | 0.964 | 0.846 | 0.927 | 0.809 |
| SA-Net | 0.970 | 0.838 | 0.923 | 0.888 | 0.976 | 0.857 | 0.936 | 0.900 |

each feature map is of the size $3 \times 3$. Following the binarization and coding layer (with $L = 8$), we now have $64/8 = 8$ coding feature maps, each of which has the size of $3 \times 3$. In other words, the dimensionality of the features for the final $k$-means procedure is 72.

Tab. 2 lists the experimental results of our proposed SA-Net in comparison with the benchmarks. Across all the four assessment metrics, our proposed SA-Net achieves the best performances indicating that the proposed network can learn feasible deep features for the clustering task.

In order to show the robustness of our method against parameter variations, we conduct further experiments on the MNIST testing subset and list the results in table 3. In both of the two spectral analysis layers, we change the parameter $k$ in k-nearest-neighborhood affinity matrix and parameter $\sigma$ in the dense affinity matrix. The first layer has two k-nearest-neighborhood affinity matrices and two dense matrices determined by $\sigma$. Thus, we have two values for $k$ and two values for $\sigma$ in the first layer.

Similarly, we have two values for $k$ and one value for $\sigma$ in the second layer. With three different settings for these parameters, the accuracy of our method only varies slightly, which are 0.974, 0.966, and 0.969, validating that our proposed SA-Net does achieve a good level of robustness against the parameter variations.

Table 3: The accuracies achieved by SA-Net with different parameters. For simplicity, the following abbreviations are adopted: SAL: spectral analysis layer; PL: pooling layer; BL: binarization layer; CL: coding layer. We also use ✓ and ✗ to indicate whether the layer is included or not in our experiments.

| The 1st SAL | The 2nd SAL | PL | BL | CL | ACC |
|---|---|---|---|---|---|
| | | ✓ | ✓ | ✓ | 0.974 |
| | | ✓ | ✓ | ✗ | 0.957 |
| $k$=5,17;$\sigma$=0.1,0.05 | $k$=5,17;$\sigma$=0.1 | ✓ | ✗ | ✗ | 0.891 |
| | | ✗ | ✗ | ✗ | 0.884 |
| | | ✓ | ✓ | ✓ | 0.966 |
| | | ✓ | ✓ | ✗ | 0.942 |
| $k$=9,21;$\sigma$=0.05,0.01 | $k$=9,21;$\sigma$=0.01 | ✓ | ✗ | ✗ | 0.903 |
| | | ✗ | ✗ | ✗ | 0.876 |
| | | ✓ | ✓ | ✓ | 0.969 |
| | | ✓ | ✓ | ✗ | 0.934 |
| $k$=5,21;$\sigma$=0.1,0.05 | $k$=5,21;$\sigma$=0.01 | ✓ | ✗ | ✗ | 0.912 |
| | | ✗ | ✗ | ✗ | 0.906 |

The experimental results in table 3 also empirically show the effectiveness of the pooling layer, the binarization layer, and the coding layer. When we remove the coding layer, as seen, the accuracy drops about $2\%$ in all of the three settings. This is due to the fact that the coding layer assigns higher weights to the feature maps in order to increase the discriminating power, yet all of the feature maps have equal weights without the coding layer. Although we may lose information in the operations of pooling and binarization, they can indeed improve the accuracy by reducing the intra-cluster distance and improving the cluster compactness.

*4.2. Face Image Clustering against occlusions*

To have more comprehensive evaluations upon our proposed SA-Net, we carry out another phase of experiments by testing our proposed against occlusions via three publicly available face image data sets AR [50], YaleB [51], and CMU PIE [52]. While the occlusions in YaleB-O and CMU PIE-O are introduced by ourselves, the occlusions in AR are real ones generated by scarf and glasses.

The AR dataset [50] consists of more than $4,000$ frontal face images from $126$ people ($70$ men and $56$ women). The face images were captured under different conditions introduced by facial expression, illumination variation, and disguises (sunglasses and scarf). The images were captured in two sessions (with an interval of two weeks). In our experiment, we use a subset of the AR dataset, consisting of $14$ non-occluded images and $12$ occluded images for each of the $120$ persons. Fig. 6 shows three groups of sample images out of 3 different persons.



Figure 6: Illustration of example images from the AR dataset.

The YaleB facial image dataset [51] has around $64$ near frontal images from $38$ individuals. The images are captured under different illuminations. Totally, this dataset has $2414$ face images. The CMU PIE dataset [52] is collected by Carnegie Mellon University. The face images in this dataset are captured with different poses, illumination conditions and expressions. We use $2,856$ images from $68$ persons.

Compared with the AR dataset, neither YaleB nor CMU PIE has occluded face

images. To test the robustness of our proposed SA-Net, we create two datasets with occlusions, i.e. CMU PIE occlusion (CMU PIE-O) and YaleB occlusion (YaleB-O). In these two occluded datasets, we simulate contiguous occlusions by hiding $20\%$ of pixels at a randomly selected location with a block out of another irrelevant image, some samples of which are shown in Fig. 7.
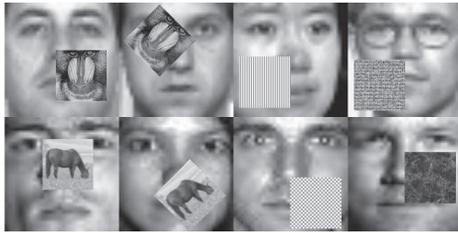


Figure 7: Illustration of eight example face images from YaleB and CMU PIE whose twenty percent pixels are occluded. The occlusions are not related to the face images and are randomly located.

Table 4: Comparative results on AR face image dataset

|        | ACC   | ARI   | NMI   | FS    |        | ACC   | ARI   | NMI   | FS    |
|--------|-------|-------|-------|-------|--------|-------|-------|-------|-------|
| KM     | 0.534 | 0.408 | 0.500 | 0.347 | LDMGI  | 0.460 | 0.430 | 0.451 | 0.392 |
| N-Cuts | 0.327 | 0.311 | 0.411 | 0.301 | NMF-D  | 0.314 | 0.449 | 0.675 | 0.462 |
| SC-ST  | 0.311 | 0.291 | 0.416 | 0.289 | DEC    | 0.308 | 0.572 | 0.726 | 0.491 |
| SC-LS  | 0.714 | 0.627 | 0.706 | 0.637 | JULE   | 0.659 | 0.599 | 0.681 | 0.614 |
| SEC    | 0.804 | 0.700 | 0.779 | 0.766 | DEPICT | 0.544 | 0.509 | 0.511 | 0.552 |
| AC-PIC | 0.115 | 0.095 | 0.017 | 0.154 | SA-Net | 0.855 | 0.729 | 0.840 | 0.754 |

Following the same design as for the clustering of digit images, we apply the proposed SA-Net shown in Fig. 3 to face image clustering. While each spectral analysis procedure produces $n_{eig} = 64$ eigenvectors in the first layer, the second layer produces $n_{eig} = 16$ eigenvectors. However, some of the implementation details need to be adapted correspondingly. Firstly, for face image clustering, we sample patches of size $15 \times 15$ with stride of 7. The patches are larger than the ones used in the handwritten image dataset, in order to allow the typical visual patterns to cover meaningful parts

Table 5: Comparative results on CMU PIE and PIE-O face image datasets

| | CMU PIE | | | | CMU PIE-O | | | |
|---|---|---|---|---|---|---|---|---|
| | ACC | ARI | NMI | FS | ACC | ARI | NMI | FS |
| KM | 0.196 | 0.190 | 0.266 | 0.153 | 0.158 | 0.154 | 0.258 | 0.123 |
| N-Cuts | 0.129 | 0.156 | 0.223 | 0.111 | 0.105 | 0.128 | 0.211 | 0.081 |
| SC-ST | 0.218 | 0.198 | 0.295 | 0.250 | 0.185 | 0.170 | 0.270 | 0.224 |
| SC-LS | 0.282 | 0.191 | 0.277 | 0.265 | 0.246 | 0.169 | 0.189 | 0.242 |
| SEC | 0.112 | 0.122 | 0.126 | 0.149 | 0.079 | 0.086 | 0.105 | 0.086 |
| AC-PIC | 0.244 | 0.320 | 0.221 | 0.189 | 0.224 | 0.283 | 0.194 | 0.169 |
| LDMGI | 0.256 | 0.194 | 0.242 | 0.257 | 0.198 | 0.158 | 0.214 | 0.223 |
| NMF-D | 0.310 | 0.357 | 0.380 | 0.292 | 0.280 | 0.324 | 0.266 | 0.266 |
| DEC | 0.421 | 0.389 | 0.477 | 0.348 | 0.403 | 0.334 | 0.347 | 0.291 |
| JULE | 0.550 | 0.437 | 0.521 | 0.497 | 0.522 | 0.421 | 0.311 | 0.423 |
| DEPICT | 0.535 | 0.484 | 0.488 | 0.453 | 0.516 | 0.374 | 0.320 | 0.347 |
| SA-Net | 0.610 | 0.497 | 0.605 | 0.537 | 0.569 | 0.423 | 0.567 | 0.518 |

of the faces. Secondly, we use the Nystrom approximation method (a different method from the digit image clustering experiment) to compute the spectral features from the dense affinity matrices. It is observed in our experiments that the computing method has little impact upon the clustering performances.

Table 4, 5, and 6 summarize all the experimental results, from which it can be seen that our proposed SA-Net outperforms all the 11 benchmarks selected out of the existing clustering algorithms. Further examinations of the experimental results also reveal that, compared with all other spectral analysis-based clustering methods, the proposed SA-Net achieves additional advantages in dealing with the occluded face images, due to the fact that the network allows us to identify the local similarity between the face images at patch-level. In contrast, the existing spectral analysis-based methods only consider the global similarity between face images, yet the occlusion can significantly reduce the similarity between two images, even when they associate with the same person. This explains why the performances of the existing methods drop significantly on

Table 6: Comparative results on YaleB and YaleB-O face image datasets

| | YaleB | | | | YaleB-O | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | ACC | ARI | NMI | FS | ACC | ARI | NMI | FS |
| KM | 0.234 | 0.208 | 0.287 | 0.216 | 0.117 | 0.070 | 0.111 | 0.083 |
| N-Cuts | 0.163 | 0.139 | 0.219 | 0.160 | 0.093 | 0.113 | 0.163 | 0.150 |
| SC-ST | 0.308 | 0.276 | 0.268 | 0.321 | 0.216 | 0.170 | 0.236 | 0.197 |
| SC-LS | 0.577 | 0.502 | 0.642 | 0.547 | 0.167 | 0.108 | 0.194 | 0.125 |
| SEC | 0.401 | 0.438 | 0.431 | 0.452 | 0.230 | 0.158 | 0.157 | 0.152 |
| AC-PIC | 0.472 | 0.482 | 0.412 | 0.428 | 0.264 | 0.208 | 0.242 | 0.241 |
| LDMGI | 0.574 | 0.597 | 0.615 | 0.509 | 0.230 | 0.176 | 0.172 | 0.223 |
| NMF-D | 0.578 | 0.520 | 0.637 | 0.613 | 0.393 | 0.296 | 0.310 | 0.242 |
| DEC | 0.571 | 0.597 | 0.569 | 0.606 | 0.402 | 0.245 | 0.469 | 0.441 |
| JULE | 0.610 | 0.573 | 0.697 | 0.649 | 0.471 | 0.460 | 0.496 | 0.358 |
| DEPICT | 0.678 | 0.619 | 0.650 | 0.617 | 0.418 | 0.393 | 0.441 | 0.380 |
| SA-Net | 0.766 | 0.648 | 0.707 | 0.681 | 0.621 | 0.526 | 0.580 | 0.497 |

both Yale-O and CMU PIE-O.

In order to show that every spectral analysis procedure contributes positively to the clustering task, we adopt different number of spectral analysis procedures in our experiments and show the average accuracy in Fig. 8. Specifically, we keep all of the spectral analysis procedures in one layer and remove one or more spectral analysis procedures in the other layer. In Fig. 8 (a), we use 1 to 8 spectral analysis procedures in the first layer and 4 spectral analysis procedures in the second layer. In Fig. 8 (b), we use 8 spectral analysis procedures in the first layer and 1 to 4 in the second layer. As seen, the average accuracy increases when we use more spectral analysis procedures. This means that, as more spectral analysis procedures are added with our proposed SA-Net, more discriminative spectral features and subspaces are brought in to increase the discriminating power of our proposed method. In other words, the additional spectral analysis procedures bring beneficial subspaces for clustering, indicating that the parallel spectral analysis procedures are indeed working together collaboratively and
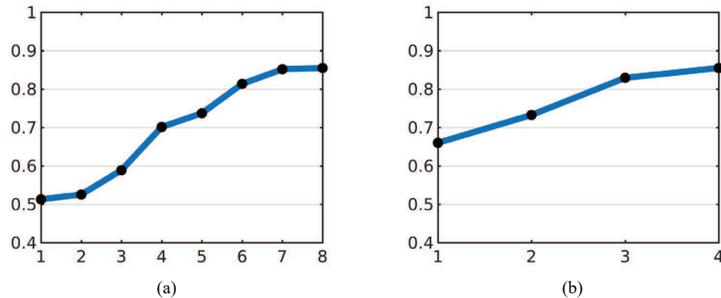
Figure 8: The average accuracy (vertical) versus the number of spectral analysis procedures (horizontal) in the first layer (a) and the second layer (b)

collectively. In principle, we can adopt more spectral analysis procedures, but the increase on clustering accuracies remains trivial, revealing that our choice of 8 spectral analysis procedures in the first layer and 4 in the second layer is sufficient not only in discovering clustering-friendly representations, but also in representing an appropriate balance between the computing cost and the effectiveness.

### 4.3. Natural Image Clustering

To assess how our proposed SA-Net performs on natural image clustering, we further conduct experiments on three more image datasets, STL-10, CIFAR-10, and CIFAR-100. They respectively have 13k, 60k, and 60k images, and the number of clusters for them are respectively 10, 10, and 20. For training purposes, we resize the images from STL-10 to $32 \times 32 \times 3$, and let them be the same size as the images from the other two datasets.

We compare the proposed SA-Net with four representative deep learning methods, i.e. SEC [15], DEC [26], JULE [28], and DEPICT [22]. Based on the four types of layers (i.e. SAL-Spectral Analysis Layer, PL-Pooling Layer, BL-Binarization Layer, and CL-Coding Layer), we build three different network structures, including SA-Net-2 (*SAL-SAL-PL-BL-CL*), SA-Net-3 (*SAL-SAL-PL-SAL-PL-BL-CL*), and SA-Net-5 (*SAL-SAL-PL-SAL-PL-SAL-PL-SAL-PL-BL-CL*).

In the first spectral layer, we sample $11 \times 11 \times 3$ image patches with a stride of $2$ pixels both vertically and horizontally. For each image of size $32 \times 32 \times 3$, we obtain

27

Table 7: Results of SA-Net and baselines on STL-10, CIFAR-10, and CIFAR-100

|  |  | SEC | DEC | JULE | DEPICT | SA-Net-2 | SA-Net-3 | SA-Net-5 |
|---|---|---|---|---|---|---|---|---|
| STL-10 | ACC | 0.148 | 0.276 | 0.182 | 0.195 | 0.317 | **0.331** | 0.328 |
|  | ARI | 0.125 | 0.184 | 0.207 | 0.231 | 0.274 | **0.304** | 0.295 |
|  | NMI | 0.212 | 0.359 | 0.277 | 0.264 | **0.398** | 0.351 | 0.386 |
|  | FS | 0.216 | 0.276 | 0.237 | 0.245 | 0.327 | 0.330 | **0.353** |
| CIFAR-10 | ACC | 0.168 | 0.257 | 0.192 | 0.244 | 0.322 | 0.297 | **0.336** |
|  | ARI | 0.175 | 0.206 | 0.291 | 0.273 | 0.316 | **0.342** | 0.318 |
|  | NMI | 0.204 | 0.301 | 0.272 | 0.295 | 0.341 | 0.356 | **0.374** |
|  | FS | 0.237 | 0.221 | 0.268 | 0.217 | 0.283 | 0.316 | **0.323** |
| CIFAR-100 | ACC | 0.094 | 0.136 | 0.103 | 0.187 | 0.208 | 0.223 | **0.254** |
|  | ARI | 0.141 | 0.154 | 0.169 | 0.184 | 0.194 | 0.222 | **0.240** |
|  | NMI | 0.123 | 0.185 | 0.137 | 0.240 | 0.274 | 0.276 | **0.293** |
|  | FS | 0 106 | 0.162 | 0.173 | 0.200 | 0.207 | 0.246 | **0.273** |

$16 \times 16 = 256$ image patches. To deal with these difficult image datasets, we use 16 different spectral analysis procedures, involving 8 different affinity matrices and 2 normalized Laplacian matrices. The 8 different affinity matrices are: three k-nearest-neighborhood affinity matrices with $k = 9, 17, 21$; three $\epsilon$-nearest-neighborhood affinity matrices with $\epsilon = 0.5\eta, \eta, 2\eta$; and two fully connected affinity matrices determined by a self-tunning method and $\sigma = 0.1$. Each spectral analysis procedure produces spectral features with dimensionality of $n_{eig} = 16$, and thus the dimensionality of the first layer spectral feature for an image is $16 \times 16 \times 256$.

In the second spectral analysis layer, we sample $5 \times 5 \times 256$ feature patches and set the stride to be 2. With 8 different spectral analysis procedures each contributing $n_{eig} = 8$ spectral features, we obtain the second layer spectral feature with a dimensionality of $8 \times 8 \times 64$. In each of the following (i.e. the third, forth, and fifth) spectral analysis layers, we set the patch size to be $3 \times 3$ and adopt 8 different spectral analysis procedures. After pooling, binarization and coding, the dimensionality of the final feature map for an image is $4 \times 4 \times 8 = 128$.

Tab. 7 summarizes the experimental results for both the benchmarks and our meth-

ods, where the best performances are highlighted in bold. As seen, the proposed methods outperform the benchmarks in terms of ACC, ARI, NMI, and FS. While SA-Net-5 achieves the best performance in 8 cases, SA-Net-3 and SA-Net-2 achieve the best in 3 cases and 1 case, respectively. In the most difficult CIFAR-100 dataset, the deepest SA-Net-5 achieves the best performances across all of the four evaluating metrics. Taking FS as the example, SA-Net-5 beats SA-Net-2 by a margin of $0.066$ on CIFAR-100. In terms of NMI, even the shallowest network SA-Net-2 can outperform the existing benchmarks by a margin larger than $0.03$ on all of the three datasets. In terms of ACC (or FS), both SA-Net-3 and SA-Net-5 outperform the existing benchmarks by a margin larger than $0.035$ (or $0.045$) on all of the three datasets.



Figure 9: The STL-10 images which are far away from their associating cluster centers

To show the difficulty of these clustering tasks, Fig. 9 illustrates some sample images from STL-10. As seen, the samples are indeed far away from their associating cluster centers. While the *planes* are captured under different views in the first column, the *cats* are heavily occluded in the fourth column and the *horses* are captured under quite different backgrounds in the eighth column.

*4.4. Fashion image clustering*

To test our proposed SA-Net for its capability in clustering variety of images with different styles, we carry out one more phase of experiments to cluster the fashion images into different styles on the dataset HipsterWars [53]. This dataset consists of $1,893$ fashion images, each associating with one of five style categories, including hipster, bohemian, pinup, preppy, and goth. The numbers of images in these five categories

are 376, 462, 191, 437, and 427 respectively. For the convenience of implementation without losing generality, we resize all the images into $600 \times 400 \times 3$.

As in Sec. 4.3, we also adopt three different network structures, i.e. SA-Net-2, SA-Net-3, and SA-Net-5. We compare our method with three existing state of the art benchmarks, including StyleNet [54], ResNet [55], and PolyLDA [56]. The StyleNet is a network for clothing that is trained from the Fashion 144K dataset [57], and the ResNet is a popular network for image classification. We extract features from these two networks and obtain the clusters by $k$-means. PolyLDA (polylingual Latent Dirichlet Allocation) is a Bayesian nonparametric model to characterize the styles by discovering the compositions of lower-level visual cues.

In this experiment, the first spectral analysis layer samples image patches of size $32 \times 32$. The main goal of the first layer is to discover the typical visual patterns that appears in many fashion images, where we use the 8 different spectral analysis procedures as in the digit image clustering to learn the spectral features. In the second and the subsequent spectral analysis layers, we only use sparse affinity matrix constructed by the $k$-nearest-neighborhood, with the parameter $k$ equals to $5, 9, 17$ and $21$, respectively. We adopt both the symmetric normalized matrix and the left normalized matrix, and apply the Lanczos method for Laplacian matrix decomposition to produce spectral features. A spectral analysis procedure produces $n_{eig} = 64$ eigenvectors in the first and $n_{eig} = 16$ eigenvectors in the second or the subsequent spectral analysis layers.

Table 8: Comparative results on HipsterWars dataset

|      | StyleNet | ResNet | PolyLDA | SA-Net-2 | SA-Net-3 | SA-Net-5 |
|------|----------|--------|---------|----------|----------|----------|
| ACC  | 0.39     | 0.30   | 0.50    | 0.54     | 0.54     | 0.55     |
| ARI  | 0.14     | 0.12   | 0.18    | 0.22     | 0.24     | 0.25     |
| NMI  | 0.20     | 0.16   | 0.21    | 0.20     | 0.20     | 0.21     |
| FS   | 0.30     | 0.28   | 0.33    | 0.39     | 0.41     | 0.41     |

As seen from Tab. 8, our method performs better than the existing benchmarks in the four evaluating metrics. Specifically, our proposed method can improve the ACC by a margin of $5\%$, the ARI by a margin of $7\%$, the FS by a margin of $8\%$.

Table 9: The confusion matrix of the proposed SA-Net-2 on the Hipsterwars dataset

|          | Hipster | Bohemian | Pinup | Preppy | Goth |
|----------|---------|----------|-------|--------|------|
| Hipster  | 140     | 83       | 21    | 81     | 51   |
| Bohemian | 47      | 328      | 18    | 28     | 41   |
| Pinup    | 30      | 25       | 72    | 23     | 41   |
| Preppy   | 92      | 49       | 35    | 202    | 59   |
| Goth     | 64      | 20       | 28    | 32     | 283  |



Figure 10: Illustration of the fashion image samples nearest to the cluster centers of the five different styles

As seen, there exist no significant differences among the performances of our three network structures, indicating that two spectral analysis layers are sufficient for this small dataset.

For the convenience of further examination and analysis, Tab. 9 illustrates the values of our confusion matrix for SA-Net-2, and Fig. 10 illustrates some image samples that are nearest to the cluster centers.

In addition, further examination reveals that all the compared methods fail to achieve good clustering results on this dataset. This is mainly due to two reasons. Firstly, the images are obtained from on-line, not captured under any controlled environment, and

as a result, they are significantly different from each other in terms of the capturing view, illumination, and background (see Fig. 10). Secondly, a style (and the resulting cluster) normally represents the coherent latent appearance between different parts of the fashion images, not simply a composition of several visual components. In other words, the foregrounds of two images can be quite different even though they are from the same cluster.

## 5. Conclusions

In this paper, we have described a new deep learning network SA-Net for image clustering based on the technique of spectral analysis. This provides one more method for deep representation learning, in addition to the popular convolutional neural network. Our proposed network structure has four type of layers, including spectral analysis layer, binarization layer, coding layer, and pooling layer. Compared with the existing spectral analysis clustering methods, SA-Net achieves three advantages. Firstly, while the existing spectral clustering methods learn representations by a single spectral analysis procedure, our proposed SA-Net conducts multiple procedures in both consecutive and parallel manner to learn more clustering-friendly representations. Secondly, SA-Net can elegantly integrate different spectral analysis procedures and thus capable of dealing with different data sample sets. Thirdly, by conducting spectral analysis procedures on image patches, SA-Net can discover the local similarity among images at patch-level, and hence it is more robust against occlusions than the existing spectral clustering methods. Extensive experiments validate the effectiveness of the proposed SA-Net on a range of different image clustering tasks, including handwritten digit image clustering, face image clustering, natural image clustering, and fashion image clustering.

## References

## References

[1] R. Xu, D. Wunsch, Survey of clustering algorithms, IEEE Transactions on Neural Networks 16 (3) (2005) 645–678.

[2] Y. Yang, D. Xu, F. Nie, S. Yan, Y. Zhuang, Image clustering using local discriminant models and global integration, TIP 19 (10) (2010) 2761–2773.

[3] G. Trigeorgis, K. Bousmalis, S. Zafeiriou, B. W. Schuller, A deep semi-nmf model for learning hidden representations, in: ICML, 2014, pp. 1692–1700.

[4] H. Chen, W. Wang, X. Feng, R. He, Discriminative and coherent subspace clustering, Neurocomputing 284 (2018) 177 – 186.

[5] C. H. Wu, C. S. Ouyang, L. W. Chen, L. W. Lu, A new fuzzy clustering validity index with a median factor for centroid-based clustering, IEEE Transactions on Fuzzy Systems 23 (3) (2015) 701–718.

[6] A. A. Liu, Y. T. Su, W. Z. Nie, M. Kankanhalli, Hierarchical clustering multi-task learning for joint human action grouping and recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence 39 (1) (2017) 102–114.

[7] W. Zhang, D. Zhao, X. Wang, Agglomerative clustering via maximum incremental path integral, Pattern Recognition 46 (11) (2013) 3056–3065.

[8] Q. Tong, X. Li, B. Yuan, Efficient distributed clustering using boundary information, Neurocomputing 275 (2018) 2355 – 2366.

[9] Y. Cai, Y. Jiao, W. Zhuge, H. Tao, C. Hou, Partial multi-view spectral clustering, Neurocomputing 311 (2018) 316 – 324.

[10] Y. Lee, J. Im, S. Cho, J. Choi, Applying convolution filter to matrix of word-clustering based document representation, Neurocomputing 315 (2018) 210 – 220.

[11] Y. Zhao, Y. Yuan, F. Nie, Q. Wang, Spectral clustering based on iterative optimization for large-scale and high-dimensional data, Neurocomputing 318 (2018) 227 – 235.

[12] X. Chen, D. Cai, Large scale spectral clustering with landmark-based representation, in: AAAI, 2011, pp. 313–318.

[13] A. Y. Ng, M. I. Jordan, Y. Weiss, On spectral clustering: Analysis and an algorithm.

[14] J. Li, Y. Xia, Z. Shan, Y. Liu, Scalable constrained spectral clustering, IEEE Transactions on Knowledge & Data Engineering 27 (2) (2015) 589–593.

[15] F. Nie, Z. Zeng, I. W. Tsang, D. Xu, C. Zhang, Spectral embedded clustering: A framework for in-sample and out-of-sample spectral clustering, IEEE Transactions on Neural Networks 22 (11) (2011) 1796–1808.

[16] H. Liu, J. Wu, T. Liu, D. Tao, Y. Fu, Spectral ensemble clustering via weighted k-means: Theoretical and practical evidence, IEEE Transactions on Knowledge & Data Engineering 29 (5) (2017) 1129–1143.

[17] S. X. Yu, J. Shi, Multiclass spectral clustering, in: ICCV, 2003, pp. 313–319 vol.1.

[18] B. Nadler, S. Lafon, R. R. Coifman, I. G. Kevrekidis, Diffusion maps, spectral clustering and eigenfunctions of fokker-planck operators, in: NIPS, NIPS'05, 2005, pp. 955–962.

[19] Z. Kang, C. Peng, Q. Cheng, Z. Xu, Unified spectral clustering with optimal graph, CoRR abs/1711.04258 (2017). `arXiv:1711.04258`.

[20] J. Shi, J. Malik, Normalized cuts and image segmentation, IEEE Transactions on Pattern Analysis and Machine Intelligence 22 (8) (2000) 888–905.

[21] U. von Luxburg, A tutorial on spectral clustering, Statistics and Computing 17 (4) (2007) 395–416.

[22] K. G. Dizaji, A. Herandi, C. Deng, W. Cai, H. Huang, Deep clustering via joint convolutional autoencoder embedding and relative entropy minimization, in: ICCV, 2017, pp. 5747–5756.

[23] F. Li, H. Qiao, B. Zhang, X. Xi, Discriminatively boosted image clustering with fully convolutional auto-encoders, CoRR abs/1703.07980 (2017).

[24] U. Shaham, K. P. Stanton, H. Li, B. Nadler, R. Basri, Y. Kluger, Spectralnet: Spectral clustering using deep neural networks, CoRR abs/1801.01587 (2018).

[25] F. Tian, B. Gao, Q. Cui, E. Chen, T.-Y. Liu, Learning deep representations for graph clustering, in: AAAI, AAAI'14, 2014.

[26] J. Xie, R. Girshick, A. Farhadi, Unsupervised deep embedding for clustering analysis, in: ICML, 2016, pp. 478–487.

[27] B. Yang, X. Fu, N. D. Sidiropoulos, M. Hong, Towards k-means-friendly spaces: Simultaneous deep learning and clustering (2016).

[28] J. Yang, D. Parikh, D. Batra, Joint unsupervised learning of deep representations and image clusters, in: CVPR, 2016, pp. 5147–5156.

[29] J. Wang, J. Jiang, An unsupervised deep learning framework via integrated optimization of representation learning and gmm-based modeling, in: ACCV, 2018, pp. 249–265.

[30] E. Aljalbout, V. Golkov, Y. Siddiqui, D. Cremers, Clustering with deep learning: Taxonomy and new methods, CoRR abs/1801.07648 (2018).

[31] T. H. Chan, K. Jia, S. Gao, J. Lu, Z. Zeng, Y. Ma, Pcanet: A simple deep learning baseline for image classification?, IEEE Transactions on Image Processing 24 (12) (2015) 5017–5032.

[32] F. R. K. Chung, Spectral Graph Theory, CBMS Regional Conference Series in Mathematics, 1997.

[33] L. Zelnik-Manor, P. Perona, Self-tuning spectral clustering, in: NIPS, NIPS'04, 2004, pp. 1601–1608.

[34] X. Zhu, C. C. Loy, S. Gong, Constructing robust affinity graphs for spectral clustering, in: CVPR, 2014, pp. 1450–1457.

[35] I. S. Dhillon, Y. Guan, B. Kulis, Weighted graph cuts without eigenvectors a multilevel approach, IEEE Trans. Pattern Anal. Mach. Intell. 29 (11) (2007) 1944–1957.

[36] D. Yan, L. Huang, M. I. Jordan, Fast approximate spectral clustering, in: SIGKDD, KDD '09, 2009, pp. 907–916.

[37] K. Zhang, I. W. Tsang, J. T. Kwok, Improved nystrÖm low-rank approximation and error analysis, in: ICML, ICML '08, 2008, pp. 1232–1239.

[38] C. Fowlkes, S. Belongie, F. Chung, J. Malik, Spectral grouping using the nystrom method, IEEE Transactions on Pattern Analysis and Machine Intelligence 26 (2) (2004) 214–225.

[39] L. Wang, M. Dong, A. Kotov, Multi-level approximate spectral clustering, in: 2015 IEEE International Conference on Data Mining, 2015, pp. 439–448.

[40] Y. Han, M. Filippone, Mini-batch spectral clustering, in: IJCNN, 2017, pp. 3888–3895.

[41] J. Wang, G. Wang, Hierarchical spatial sum product networks for action recognition in still images, IEEE Transactions on Circuits and Systems for Video Technology 28 (1) (2018) 90–100.

[42] G. Hinton, R. Salakhutdinov, Reducing the dimensionality of data with neural networks, Science 313 (5786) (2006) 504 – 507.

[43] G. Chen, Deep learning with nonparametric clustering, CoRR abs/1501.03084
(2015). `arXiv:1501.03084`.
URL `http://arxiv.org/abs/1501.03084`

[44] J. Wang, A. Hilton, J. Jiang, Spectral analysis network for deep representation
learning and image clustering, in: 2019 IEEE International Conference on Multi-
media and Expo (ICME), 2019, pp. 1540–1545.

[45] J. Wang, Z. Wang, D. Tao, S. See, G. Wang, Learning common and specific fea-
tures for rgb-d semantic segmentation with deconvolutional networks, in: ECCV,
2016.

[46] A. Coates, A. Y. Ng, Learning Feature Representations with K-Means, Springer
Berlin Heidelberg, 2012.

[47] Y. Saad, Numerical methods for large eigenvalue problems, Manchester Univer-
sity Press, 2011.

[48] A. J. Smola, B. Schökopf, Sparse greedy matrix approximation for machine learn-
ing, in: ICML, 2000, pp. 911–918.

[49] Y. Lecun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied to
document recognition, Proceedings of the IEEE 86 (11) (1998) 2278–2324.

[50] A. Martinez, R. Benavente., The ar face database, CVC Technical Report 24 (June
1998).

[51] A. S. Georghiades, P. N. Belhumeur, D. J. Kriegman, From few to many: Illumi-
nation cone models for face recognition under variable lighting and pose, IEEE
Transactions on Pattern Analysis and Machine Intelligence 23 (6) (2002) 643–
660.

[52] T. Sim, S. Baker, M. Bsat, The cmu pose, illumination, and expression (pie)
database, in: Proceedings of Fifth IEEE International Conference on Automatic
Face Gesture Recognition, 2002, pp. 46–51.

[53] M. H. Kiapour, K. Yamaguchi, A. C. Berg, T. L. Berg, Hipster wars: Discovering elements of fashion styles, in: ECCV, Springer International Publishing, 2014, pp. 472–488.

[54] E. Simo-Serra, H. Ishikawa, Fashion style in 128 floats: Joint ranking and classification using weak data for feature extraction, in: CVPR, 2016, pp. 298–307.

[55] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: CVPR, 2016, pp. 770–778.

[56] W. L. Hsiao, K. Grauman, Learning the latent "look": Unsupervised discovery of a style-coherent embedding from fashion images, in: ICCV, 2017, pp. 4213–4222.

[57] E. Simo-Serra, S. Fidler, F. Moreno-Noguer, R. Urtasun, Neuroaesthetics in Fashion: Modeling the Perception of Fashionability, in: CVPR, 2015.