

Fine-Grained Bipartite Concept Factorization for Clustering

Chong Peng¹, Pengfei Zhang¹, Yongyong Chen^{2,*}, Zhao Kang³, Chenglizhao Chen^{4,*}, and Qiang Cheng⁵

¹ College of Computer Science and Technology, Qingdao University

² School of Computer Science and Technology, Harbin Institute of Technology (Shenzhen)

³ School of Computer Science and Engineering, University of Electronic Science and Technology of China

⁴ College of Computer Science and Technology, China University of Petroleum (East China)

⁵ Department of Computer Science, University of Kentucky

{qdcpeng, zpf010108, yongyongchen.cn, sckangz}@gmail.com, cclz123@163.com, qiang.cheng@uky.edu

Abstract

In this paper, we propose a novel concept factorization method that seeks factor matrices using a cross-order positive semi-definite neighbor graph, which provides comprehensive and complementary neighbor information of the data. The factor matrices are learned with bipartite graph partitioning, which exploits explicit cluster structure of the data and is more geared towards clustering application. We develop an effective and efficient optimization algorithm for our method, and provide elegant theoretical results about the convergence. Extensive experimental results confirm the effectiveness of the proposed method.

1. Introduction

Matrix factorization aims at seeking two or more factor matrices, whose product may well approximate the original data [25]. Matrix factorization has been widely studied for dimension reduction [25, 48] and low-dimensional representation [44] of high-dimensional data, which is essentially important in various learning tasks such as clustering [12], classification [24], foreground-background separation in surveillance video [30], community detection [27], link prediction [26], hyperspectral unmixing [8], etc.

Among these learning tasks, nonnegative data are typical due to the fact that they naturally contain only nonnegative elements [15], such as pixels of images [2, 36], connectivity of nodes in a social network [29], etc. For such type of data, the nonnegative matrix factorization (NMF) has been developed for parts-based representation [15, 16], which shows its advantage by mimicking the signal processing mechanism of the human brain [15].

NMF methods have been widely studied during the last decades [19, 39] with various extensions [11, 19, 37]. Among them, the Convex-NMF (CNMF) is quite typical

by restricting the basis vectors to be convex combinations of the samples [6, 50]. Let $\mathbf{X} \in \mathbb{R}^{d \times n}$ be the data matrix with d features, n samples, and c clusters, then the CNMF can be represented as $\mathbf{X} \approx \mathbf{X}\mathbf{W}\mathbf{V}^T$, where $\mathbf{W} \in \mathbb{R}^{n \times c}$ and $\mathbf{V} \in \mathbb{R}^{n \times c}$ denote the so-called score and coefficient matrices [32, 38], respectively. The CNMF is also known as the concept factorization (CF) by revealing the distinct concepts, i.e., the cluster centers of the data, as well as the cluster membership of each sample [50]. The CNMF allows a direct measurement of pairwise similarity of the data [6], which is straightforward to apply the kernel technique for nonlinear relationship learning [6, 32]. Moreover, it has been revealed that the variants of CNMF such as Cluster-NMF, Semi-NMF, and Kernel-NMF, are all soft relaxations of the K-means clustering [6]. With the success of deep learning, deep factorization has been successfully developed for the NMF [28, 51] and CF community [3] from the original one-layer factorization, which helps recognize hierarchical structure of the data [3, 43].

Recently, the CF has been shown closely related with subspace clustering methods in nature [32, 38], where the product of the factor matrices can be treated as the self-expressive representation matrix [32]. These methods integrate the local geometric structure of the data in recovering the score and coefficient factor matrices [32, 38]. It has been shown that local geometric structure is critical and effective in revealing the underlying structure of the data [32], which can be described from the following perspectives: 1) It helps avoid over-fitting issue by recovering local geometric structure of the data [20]; 2) A full kernel matrix does not consider local density of the samples and might be inefficient in the representation capability [17, 21]; 3) The intrinsic information of high-dimensional data needs to be maintained by the low-dimensional representation [45]. For NMF methods, the local geometric structure of the data is often preserved by techniques such as sparse regularization [36] and graph Laplacian [2]. For the former, the lo-

cal relationship of the data is preserved by the sparse representation [46], meanwhile the parts-based representation of NMF naturally leads to sparsity [2]. For the latter, the NMF method recovers the low-dimensional representation on a local manifold that maintains local neighbors [2]. Besides, some other techniques have been developed such as inter-class separability [32] and local centroid learning [4]. All these methods only consider the first-order neighbors, which might be insufficient to fully exploit local geometric structure of the data [13], and more comprehensive information still needs to be exploited.

In fact, high-order neighbors may provide more complicated relationships and reveal latent structure of the data [13,42], which has been rarely considered in the CF or NMF community. In this paper, we fully exploit the high-order neighbor relationships of the data for concept factorization, which provides comprehensive and complementary information of the data. Moreover, we seek factorization with bipartite graph partitioning, which renders the score and coefficient matrices to have clear group structure and thus are more geared towards clustering application.

We summarize the key contributions of this paper as follows: 1) We seek the concept factorization using the fine-grained neighbor graph composed of cross-order neighbors, which provides comprehensive and complementary information of the data; 2) The cluster information is preserved by the bipartite graph partitioning, which renders the coefficient matrix to have a clear group structure; 3) The fine-grained neighbor graph allows the data to have mixed signs, which potentially extends the applicability of NMF methods in broader areas; 4) Efficient optimization algorithm is developed with theoretical convergence guarantee; 5) Extensive experimental results confirm the effectiveness of the proposed method in clustering and data representation.

2. Related Work

We briefly review some closely related methods as follows. Given data set $\mathbf{X} \in \mathbb{R}^{d \times n}$, the CF seeks a factorization with the following minimization problem [50]:

$$\min_{\mathbf{W} \geq 0, \mathbf{V} \geq 0} \|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{V}^T\|_F^2, \quad (1)$$

with $\|\cdot\|_F$ being the Frobenius norm, which allows the NMF methods to handle mixed-signed data [6]. It is seen that (1) is closely related with subspace clustering methods [32,33], which assume self-expressiveness of the data and can be expressed as $\mathbf{X} \approx \mathbf{X}\mathbf{Z}$ with $\mathbf{Z} \in \mathbb{R}^{n \times n}$ being the representation matrix [33,40]. It is straightforward to have nonlinear expansion of (1) using the kernel trick [6,31].

3. Fine-grained Bipartite CF

The CF methods have been widely studied [1, 38, 50] and have been shown to be closely related with the spec-

tral clustering-based subspace clustering [32]:

$$\min_{\mathbf{W}, \mathbf{V}} \frac{1}{2} \|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{V}^T\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}\mathbf{V}^T\|_F^2, \mathbf{W} \geq 0, \mathbf{V} \geq 0, \quad (2)$$

where $\alpha \geq 0$ is a balancing parameter, the factor $\frac{1}{2}$ is used to facilitate the optimization, and $\mathbf{W}\mathbf{V}^T \in \mathbb{R}^{n \times n}$ can be treated as a low-rank self-expressive representation matrix of subspace clustering methods [33,35]. Although the product of $\mathbf{W}\mathbf{V}^T$ has strict rank constraint, it does not necessarily leads to a proper group structure of $\mathbf{W}\mathbf{V}^T$ nor a strong clustering capability of \mathbf{V} . To further enhance the representation capability of $\mathbf{W}\mathbf{V}^T$, as well as \mathbf{V} , in revealing underlying group structure of the data, we first define the following symmetric and nonnegative bipartite graph matrix $\mathbf{S} = \begin{bmatrix} \mathbf{0} & (\mathbf{W}\mathbf{V}^T)^T \\ (\mathbf{W}\mathbf{V}^T) & \mathbf{0} \end{bmatrix}$, and its graph Laplacian matrix $\mathbf{L} = \mathbf{D}_S - \mathbf{S}$, with $\mathbf{0}$ being an $n \times n$ zero matrix, \mathbf{D}_S being diagonal, and $(\mathbf{D}_S)_{ii} = \sum_{j=1}^n \mathbf{S}_{ij}$. For ease of notation, we define an operator $\mathbb{L}(\cdot)$ such that $\mathbf{L} = \mathbb{L}(\mathbf{W}, \mathbf{V})$. Then, it is desired that \mathbf{S} has exactly c connected components that correspond to c clusters of the data, which is ensured by the following Ky-Fan's Theorem [22]:

Theorem 1 ([22]). *For any $\mathbf{S}^T = \mathbf{S} \geq 0$, let \mathbf{L} be its Laplacian matrix. Then the number of zero eigenvalues of \mathbf{L} equals the number of connected components in \mathbf{S} .*

It is straightforward to impose rank constraint of \mathbf{L} in (2) for the desired property. However, the hard rank constraint is usually difficult to solve. Thus, we relax it to minimize the c smallest eigenvalues of \mathbf{L} , which leads to

$$\min_{\mathbf{W}, \mathbf{V}, \mathbf{F}} \frac{1}{2} \|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{V}^T\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}\mathbf{V}^T\|_F^2 + \beta \sum_{i=1}^c \lambda_i(\mathbf{L}) \quad (3)$$

$$s.t. \quad \mathbf{W} \geq 0, \mathbf{V} \geq 0, \mathbf{L} = \mathbb{L}(\mathbf{W}, \mathbf{V}),$$

where $\beta \geq 0$ is a balancing parameter and $\lambda_i(\cdot)$ returns the i -th smallest eigenvalue of the input matrix. Thus, it is guaranteed that $\sum_{i=1}^c \lambda_i(\mathbf{L}) = 0$ with a properly large β . According to [7], we may rewrite (3) as the following:

$$\min_{\mathbf{W}, \mathbf{V}, \mathbf{F}} \frac{1}{2} \|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{V}^T\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}\|_F^2 + \beta \text{Tr}(\mathbf{F}^T \mathbf{L} \mathbf{F}) \quad (4)$$

$$s.t. \quad \mathbf{W} \geq 0, \mathbf{V} \geq 0, \mathbf{F}^T \mathbf{F} = \mathbf{I}_c, \mathbf{L} = \mathbb{L}(\mathbf{W}, \mathbf{V}),$$

where \mathbf{I}_c denotes an identity matrix of size $c \times c$. Moreover, it is desired to make \mathbf{V} a relaxed indicator matrix for intermediate clustering interpretation. Therefore, we impose the constraint $\mathbf{V}^T \mathbf{V} = \mathbf{I}_c$ and our model becomes

$$\min_{\mathbf{W}, \mathbf{V}, \mathbf{F}} \frac{1}{2} \|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{V}^T\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}\|_F^2 + \beta \text{Tr}(\mathbf{F}^T \mathbf{L} \mathbf{F}) \quad (5)$$

$$s.t. \quad \mathbf{W} \geq 0, \mathbf{V} \geq 0, \mathbf{F}^T \mathbf{F} = \mathbf{I}_c, \mathbf{V}^T \mathbf{V} = \mathbf{I}_c, \mathbf{L} = \mathbb{L}(\mathbf{W}, \mathbf{V}).$$

To further enhance the nonlinear learning capability, we expand the above model using the kernel trick [31]. In particular, we may expand the first term of (5) and replace $\mathbf{X}^T \mathbf{X}$ with the kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$, which leads to

$$\begin{aligned} & \min_{\mathbf{W}, \mathbf{V}, \mathbf{F}} \frac{1}{2} \text{Tr}(\mathbf{V} \mathbf{W}^T \mathbf{K} \mathbf{W} \mathbf{V}^T - 2 \mathbf{K} \mathbf{W} \mathbf{V}^T) \\ & + \frac{\alpha}{2} \|\mathbf{W}\|_F^2 + \beta \text{Tr}(\mathbf{F}^T \mathbf{L} \mathbf{F}), \quad s.t. \quad \mathbf{W} \geq 0, \mathbf{V} \geq 0, \quad (6) \\ & \mathbf{F}^T \mathbf{F} = \mathbf{I}_c, \mathbf{V}^T \mathbf{V} = \mathbf{I}_c, \mathbf{L} = \mathbb{L}(\mathbf{W}, \mathbf{V}), \end{aligned}$$

where \mathbf{K} is obtained by adopting a certain kernel function. It has been revealed that local relationships play an essential role in learning the underlying structures of the data [32,38]. However, such information is omitted in the above model. To combat this issue, we construct a local neighbor graph by element as $\mathbf{K}_{ij} = \mathbb{I}_{\{\mathbf{X}_j \in \mathcal{N}_K(\mathbf{X}_i)\}}$, where $\mathbb{I}_{\{\cdot\}}$ is an indicator function that returns 1 if the conditions in the subscript hold and 0 otherwise, \mathbf{X}_i denotes the i -th sample and $\mathcal{N}_K(\cdot)$ denotes the set of K -nearest neighbors of the input. In this paper, for simplicity and yet without loss of generality, we use the binary kernel similarity to measure the pair-wise neighbor relationships of the data [2], with K being set to 5.

Essentially, the weighted graph is closely related with the random walk and Markov transition probability matrix (TPM) [41,49,52], where the neighbor relationships are directly revealed by the probabilities of the one-step random walks. Inspired by [49], we construct the TPM as $\mathbf{P} = \mathbf{D}_K^{-1} \mathbf{K}$, where \mathbf{D}_K is diagonal and $(\mathbf{D}_K)_{ii} = \sum_{j=1}^n \mathbf{K}_{ij}$. We may also treat the probability as the soft neighbor relationship. In practice, there often exist high-order neighbor relationships among the data, which are not directly available in a weighted graph. Intuitively, such relationships can be measured by the probabilities of multi-step random walks [42]. To account for the high-order neighbor information, we define the following a -th order probability matrix of the data as $\mathbf{P}_{[a]} = \mathbf{P}^a = \underbrace{\mathbf{P} \cdot \mathbf{P} \cdot \dots \cdot \mathbf{P}}_{a \text{ times}} \in \mathbb{R}^{n \times n}$,

which measures the neighbor relationships up to an a -step random walk. To fully exploit such comprehensive cross-order neighbor information, we fuse the high-order probability matrices. Moreover, to ensure the symmetry, which is a natural and desired property for neighbor relationships, we define the following fine-grained probability matrix as $\mathbf{P}_A = \sum_{a=1}^A \frac{\mathbf{P}_{[a]} + \mathbf{P}_{[a]}^T}{2} = \sum_{a=1}^A \frac{\mathbf{P}^a + (\mathbf{P}^a)^T}{2}$, with A being the highest order of the neighbor information. Thus, with \mathbf{P}_A , our model can be finally developed as:

$$\begin{aligned} & \min_{\mathbf{W}, \mathbf{V}, \mathbf{F}} \frac{1}{2} \text{Tr}(\mathbf{V} \mathbf{W}^T \mathbf{P}_A \mathbf{W} \mathbf{V}^T - 2 \mathbf{P}_A \mathbf{W} \mathbf{V}^T) \\ & + \frac{\alpha}{2} \|\mathbf{W}\|_F^2 + \beta \text{Tr}(\mathbf{F}^T \mathbf{L} \mathbf{F}), \quad s.t. \quad \mathbf{W} \geq 0, \mathbf{V} \geq 0, \quad (7) \\ & \mathbf{F}^T \mathbf{F} = \mathbf{I}_c, \mathbf{V}^T \mathbf{V} = \mathbf{I}_c, \mathbf{L} = \mathbb{L}(\mathbf{W}, \mathbf{V}). \end{aligned}$$

We name model (7) as the **F**ine-grained **B**ipartite **C**oncept

Factorization (Figer-CF). It is noted that $\mathbf{P}_A \geq 0$, which renders the Figer-CF suitable for general data with mixed signs and potentially applicable in broader applications.

Remark 1. In practice, to ensure that \mathbf{P}_A is positive semi-definite, we add an identity matrix with a scaling factor $\tau \geq 0$ to it. This strategy makes sense because it is natural that one sample has the highest similarity with itself. In practice, we may set the scaling factor as $\tau = |\min\{0, \lambda_n(\mathbf{P}_A)\}|$. Moreover, we may provide a general lower bound of τ for any \mathbf{P}_A , which only depends on the sample size n and order of neighbors A . In rest of this paper, the adjusted $\mathbf{P}_A + \tau \mathbf{I}_n$ is denoted as \mathbf{P}_A for ease and clarity of representation. The lower bound is provided in the following Theorem 2.

Theorem 2. *It is ensured that $\mathbf{P}_A + \tau \mathbf{I}_n \succeq 0$ with $\tau \geq A\sqrt{n}$.*

Proof. According to the definition of \mathbf{P}_A , it is clear that $\mathbf{P}_A + \tau \mathbf{I}_n$ is real and symmetric. Thus, $\mathbf{P}_A + \tau \mathbf{I}_n$ is diagonalizable and we only need to show that all its eigenvalues are nonnegative. It is straightforward that $\|\mathbf{P}_A\|_F \leq \sum_{a=1}^A \|\frac{1}{2}(\mathbf{P}_{[a]} + \mathbf{P}_{[a]}^T)\|_F \leq \sum_{a=1}^A \|\mathbf{P}_{[a]}\|_F$. Moreover, $\|\mathbf{P}_{[a]}\|_F^2 = \sum_{a=1}^n \|\mathbf{P}_{[a]}^{(i)}\|_2^2 \leq \sum_{a=1}^n \|\mathbf{P}_{[a]}^{(i)}\|_1^2 \leq \sum_{a=1}^n 1 = n$. Thus, we have $\|\mathbf{P}_A\|_F \leq \sum_{a=1}^A \sqrt{n} = A\sqrt{n}$. According to Schur's inequality, we have $|\lambda_i(\mathbf{P}_A)| \leq \|\mathbf{P}_A\|_F \leq A\sqrt{n}$ for $i = 1, \dots, n$. Thus, it is straightforward that $\lambda_n(\mathbf{P}_A) \geq -A\sqrt{n}$, and $\lambda_n(\mathbf{P}_A + \tau \mathbf{I}_n) = \lambda_n(\mathbf{P}_A) + \tau \geq 0$, which guarantees that $\mathbf{P}_A + \tau \mathbf{I}_n \succeq 0$. \square

4. Optimization

In this section, we develop an effective iterative optimization algorithm.

4.1. Optimization of \mathbf{F}

Keeping the other variables fixed, the optimization problem of \mathbf{F} is as $\min_{\mathbf{F}^T \mathbf{F} = \mathbf{I}_c} \text{Tr}(\mathbf{F}^T \mathbf{L} \mathbf{F})$, which can be effectively solved by the eigenvalue decomposition. As a result, \mathbf{F} is obtained by computing the eigenvectors of \mathbf{L} that are associated with its top c eigenvalues, which is denoted as:

$$\mathbf{F} = \mathbb{E} \mathbb{V} \mathbb{D}_c(\mathbf{L}). \quad (8)$$

4.2. Optimization of \mathbf{W} and \mathbf{V}

With fixed \mathbf{F} , the subproblem of \mathbf{W} and \mathbf{V} is complicated, since \mathbf{L} is associated with both \mathbf{W} and \mathbf{V} . To facilitate the optimization, we reformulate the term $\text{Tr}(\mathbf{F}^T \mathbf{L} \mathbf{F})$ to explicitly represent it in \mathbf{W} and \mathbf{V} . First, we divide \mathbf{F} into two blocks, i.e., $\mathbf{F}^T = [\mathbf{F}_{(1)}^T, \mathbf{F}_{(2)}^T]$, with $\mathbf{F}_{(1)} \in \mathbb{R}^{n \times c}$ and $\mathbf{F}_{(2)} \in \mathbb{R}^{n \times c}$, respectively. Then, we have $\text{Tr}(\mathbf{F}^T \mathbf{L} \mathbf{F}) = \frac{1}{2} \sum_{i,j=1}^{2n} \|\mathbf{F}^i - \mathbf{F}^j\|_2^2 \mathbf{S}_{ij} = \sum_{i,j=1}^n \|\mathbf{F}_{(1)}^i - \mathbf{F}_{(2)}^j\|_2^2 (\mathbf{W} \mathbf{V}^T)_{ji} = \sum_{i,j=1}^n \mathbf{T}_{ij} (\mathbf{W} \mathbf{V}^T)_{ji} = \text{Tr}(\mathbf{T} \mathbf{W} \mathbf{V}^T)$,

where \mathbf{F}^i , $\mathbf{F}_{(1)}^i$, and $\mathbf{F}_{(2)}^i$ denote the i -th row of \mathbf{F} , $\mathbf{F}_{(1)}$, and $\mathbf{F}_{(2)}$, and $\mathbf{T}_{ij} = \|\mathbf{F}_{(1)}^i - \mathbf{F}_{(2)}^j\|_2^2$ is the (i, j) -th element of matrix $\mathbf{T} \in \mathbb{R}^{n \times n}$, respectively. Thus, the joint minimization problem associated with \mathbf{W} and \mathbf{V} becomes

$$\min_{\mathbf{W}, \mathbf{V}} \frac{1}{2} \text{Tr}(\mathbf{V}\mathbf{W}^T \mathbf{P}_A \mathbf{W}\mathbf{V}^T - 2\mathbf{P}_A \mathbf{W}\mathbf{V}^T) + \frac{\alpha}{2} \|\mathbf{W}\|_F^2 + \beta \text{Tr}(\mathbf{T}\mathbf{W}\mathbf{V}^T), \text{ s.t. } \mathbf{W} \geq 0, \mathbf{V} \geq 0. \quad (9)$$

Given \mathbf{P}_A and \mathbf{T} , we provide the following updating rules of \mathbf{W} and \mathbf{V} for the joint minimization problem (9):

$$\mathbf{W}_{ik} \leftarrow \mathbf{W}_{ik} \sqrt{\frac{(\mathbf{P}_A \mathbf{V})_{ik}}{\beta(\mathbf{T}^T \mathbf{V})_{ik} + (\mathbf{P}_A \mathbf{W}\mathbf{V}^T \mathbf{V})_{ik} + \alpha \mathbf{W}_{ik}}}, \quad (10)$$

$$\mathbf{V}_{ik} \leftarrow \mathbf{V}_{ik} \sqrt{\frac{(\mathbf{P}_A \mathbf{W})_{ik} + \beta(\mathbf{V}\mathbf{V}^T \mathbf{T}\mathbf{W})_{ik}}{(\mathbf{V}\mathbf{V}^T \mathbf{P}_A \mathbf{W})_{ik} + \beta(\mathbf{T}\mathbf{W})_{ik}}}. \quad (11)$$

We denote the objective function as $f(\mathbf{F}, \mathbf{W}, \mathbf{V})$. We repeat (8), (10) and (11) until convergence of objective value sequence $\{f(\mathbf{F}^{(t)}, \mathbf{W}^{(t)}, \mathbf{V}^{(t)})\}$ and obtain the solution $\{\mathbf{F}^{(*)}, \mathbf{W}^{(*)}, \mathbf{V}^{(*)}\}$, where we use the superscript $(\cdot)^{(t)}$ to denote the number of iterations. Then, we perform the K-means to $\mathbf{V}^{(*)}$ to obtain the final clusters of the data.

5. Convergence Analysis

In this section, we provide the main theoretical results, including correctness of the updating rules and convergence of our algorithm.

5.1. Correctness and Convergence of (10)

For the updating rule of (10), we present the following two main results: 1) When convergent, the limiting solution of (10) satisfies the KKT condition. 2) The iteration of (10) converges. We formally establish the above results in Theorems 3 and 4, respectively.

Theorem 3. Fixing \mathbf{V} and \mathbf{T} , the limiting solution of the updating rule in (10) satisfies the KKT condition.

Proof. Fixing \mathbf{V} , the subproblem for \mathbf{W} is

$$\min_{\mathbf{W}} \frac{1}{2} \text{Tr}(-2\mathbf{P}_A \mathbf{W}\mathbf{V}^T + \mathbf{V}\mathbf{W}^T \mathbf{P}_A \mathbf{W}\mathbf{V}^T) + \beta \text{Tr}(\mathbf{W}^T \mathbf{T}^T \mathbf{V}) + \frac{\alpha}{2} \text{Tr}(\mathbf{W}^T \mathbf{W}), \text{ s.t. } \mathbf{W} \geq 0. \quad (12)$$

We introduce an Lagrangian multiplier $\Psi \in \mathbb{R}^{n \times c}$ to the above problem and obtain its Lagrangian function:

$$L_{\mathbf{W}} = \frac{1}{2} \text{Tr}(-2\mathbf{P}_A \mathbf{W}\mathbf{V}^T + \mathbf{V}\mathbf{W}^T \mathbf{P}_A \mathbf{W}\mathbf{V}^T) + \frac{\alpha}{2} \text{Tr}(\mathbf{W}^T \mathbf{W}) + \beta \text{Tr}(\mathbf{W}^T \mathbf{T}^T \mathbf{V}) + \text{Tr}(\Psi \mathbf{W}^T). \quad (13)$$

First, we may obtain the gradient of $L_{\mathbf{W}}$ as

$$\frac{\partial L_{\mathbf{W}}}{\partial \mathbf{W}} = -\mathbf{P}_A \mathbf{V} + \beta \mathbf{T}^T \mathbf{V} + \mathbf{P}_A \mathbf{W}\mathbf{V}^T \mathbf{V} + \alpha \mathbf{W} + \Psi. \quad (14)$$

By the complementary slackness condition, we can see that

$$(-\mathbf{P}_A \mathbf{V} + \beta \mathbf{T}^T \mathbf{V} + \mathbf{P}_A \mathbf{W}\mathbf{V}^T \mathbf{V} + \alpha \mathbf{W})_{ik} \mathbf{W}_{ik} = \Psi_{ik} \mathbf{W}_{ik} = 0, \quad (15)$$

which provides a fixed point condition that the limiting solution of (10) should satisfy. It is clear that (10) can be reduced to the following equation:

$$(-\mathbf{P}_A \mathbf{V} + \beta \mathbf{T}^T \mathbf{V} + \mathbf{P}_A \mathbf{W}\mathbf{V}^T \mathbf{V} + \alpha \mathbf{W})_{ik} \mathbf{W}_{ik}^2 = 0, \quad (16)$$

which is identical to (15) because both of them imply the same conditions that either $\mathbf{W}_{ik} = 0$ holds in both (15) and (16) or $(-\mathbf{P}_A \mathbf{V} + \beta \mathbf{T}^T \mathbf{V} + \mathbf{P}_A \mathbf{W}\mathbf{V}^T \mathbf{V} + \alpha \mathbf{W})_{ik} = 0$ holds in both (15) and (16), which concludes the proof. \square

Before we provide the result about convergence of (10), we introduce the technique of auxiliary function [2, 16] that plays an essential role in the following analysis.

Definition 1. Given functions $\bar{L}(\mathbf{Q}, \mathbf{Q}')$ and $L(\mathbf{Q})$, then $\bar{L}(\mathbf{Q}, \mathbf{Q}')$ is called an auxiliary function of $L(\mathbf{Q})$ if

$$\bar{L}(\mathbf{Q}, \mathbf{Q}') \geq L(\mathbf{Q}) \text{ and } \bar{L}(\mathbf{Q}, \mathbf{Q}) = L(\mathbf{Q}) \quad (17)$$

hold for any \mathbf{Q} and \mathbf{Q}' .

Proposition 1. Given a function $L(\mathbf{Q})$, let $\bar{L}(\mathbf{Q}, \mathbf{Q}')$ be one of its auxiliary functions. Then, the variable sequence $\{\mathbf{Q}^{(t)}\}$, with $t \in \mathbb{N}$ and

$$\mathbf{Q}^{(t+1)} = \arg \min_{\mathbf{Q}} \bar{L}(\mathbf{Q}, \mathbf{Q}^{(t)}), \quad (18)$$

satisfies the following chain of inequalities: $L(\mathbf{Q}^{(t)}) = \bar{L}(\mathbf{Q}^{(t)}, \mathbf{Q}^{(t)}) \geq \bar{L}(\mathbf{Q}^{(t+1)}, \mathbf{Q}^{(t)}) \geq L(\mathbf{Q}^{(t+1)})$, which implies that the value sequence $\{L(\mathbf{Q}^{(t)})\}_{t=0}^{\infty}$ is monotonically decreasing (nonincreasing).

Proposition 2 ([6]). For any matrices $\Gamma \in \mathbb{R}_+^{n \times n}$, $\Omega \in \mathbb{R}_+^{c \times c}$, $\Delta \in \mathbb{R}_+^{n \times c}$, and $\Delta' \in \mathbb{R}_+^{n \times c}$, with Γ and Ω being symmetric, the following inequality holds: $\sum_{i=1}^n \sum_{s=1}^c \frac{(\Gamma \Delta' \Omega)_{is} \Delta_{is}^2}{\Delta'_{is}} \geq \text{Tr}(\Delta^T \Gamma \Delta \Omega)$.

Next, we provide the main result about the convergence of (10) in the following theorem and proof.

Theorem 4. For fixed \mathbf{V} , the objective in (12) is monotonically decreasing under the updating rule in (10).

Proof. We define the subproblem of \mathbf{W} as

$$P(\mathbf{W}) = \frac{\alpha}{2} \text{Tr}(\mathbf{W}^T \mathbf{W}) + \text{Tr}(-\mathbf{W}^T \mathbf{P}_A \mathbf{V}) + \frac{1}{2} \text{Tr}(\mathbf{W}^T \mathbf{P}_A \mathbf{W}\mathbf{V}^T \mathbf{V}) + \beta \text{Tr}(\mathbf{W}^T \mathbf{T}^T \mathbf{V}). \quad (19)$$

Then, according to Proposition 2, as well as $a \leq \frac{a^2+b^2}{2b}$ and $a \geq 1 + \log a$ for $a, b \geq 0$, it is straightforward that the following inequalities hold for any $\mathbf{W}' \in \mathbb{R}^{n \times c} \geq 0$:

$$\text{Tr}(\mathbf{W}^T \mathbf{P}_A \mathbf{W} \mathbf{V} \mathbf{V}^T \mathbf{V}) \leq \sum_{ik} \frac{(\mathbf{P}_A \mathbf{W}' \mathbf{V}^T \mathbf{V})_{ik} \mathbf{W}'_{ik}}{\mathbf{W}'_{ik}}, \quad (20)$$

$$\text{Tr}(\mathbf{W}^T \mathbf{T}^T \mathbf{V}) \leq \sum_{ik} (\mathbf{T}^T \mathbf{V})_{ik} \frac{\mathbf{W}'_{ik} + \mathbf{W}'_{ik}}{2\mathbf{W}'_{ik}}, \quad (21)$$

$$\text{Tr}(\mathbf{W}^T \mathbf{P}_A \mathbf{V}) \geq \sum_{ik} (\mathbf{P}_A \mathbf{V})_{ik} \mathbf{W}'_{ik} \left(1 + \log \frac{\mathbf{W}'_{ik}}{\mathbf{W}'_{ik}}\right), \quad (22)$$

where the equal sign holds when $\mathbf{W} = \mathbf{W}'$. Thus, we may construct an auxiliary function of $P(\mathbf{W})$ as

$$\begin{aligned} \bar{P}(\mathbf{W}, \mathbf{W}') &= - \sum_{ik} (\mathbf{P}_A \mathbf{V})_{ik} \mathbf{W}'_{ik} \left(1 + \log \frac{\mathbf{W}'_{ik}}{\mathbf{W}'_{ik}}\right) \\ &+ \frac{1}{2} \sum_{ik} \frac{(\mathbf{P}_A \mathbf{W}' \mathbf{V}^T \mathbf{V})_{ik} \mathbf{W}'_{ik}}{\mathbf{W}'_{ik}} \\ &+ \beta \sum_{ik} (\mathbf{T}^T \mathbf{V})_{ik} \frac{\mathbf{W}'_{ik} + \mathbf{W}'_{ik}}{2\mathbf{W}'_{ik}} + \frac{\alpha}{2} \sum_{ik} \mathbf{W}'_{ik}, \end{aligned} \quad (23)$$

where $\bar{P}(\mathbf{W}, \mathbf{W}') \geq P(\mathbf{W})$ and $\bar{P}(\mathbf{W}, \mathbf{W}) = P(\mathbf{W})$ hold for any $\mathbf{W}, \mathbf{W}' \geq 0$.

Next, the key is to show that (10) essentially follows (18). It is easy to verify that $\bar{P}(\mathbf{W}, \mathbf{W}')$ is convex in \mathbf{W} . Thus, the global optimum is ensured by the first-order optimality condition, which gives rise to

$$\begin{aligned} \frac{\partial \bar{P}(\mathbf{W}, \mathbf{W}')}{\partial \mathbf{W}_{ik}} &= - \frac{(\mathbf{P}_A \mathbf{V})_{ik} \mathbf{W}'_{ik}}{\mathbf{W}_{ik}} + \alpha \mathbf{W}_{ik} \\ &+ \frac{(\mathbf{P}_A \mathbf{W}' \mathbf{V}^T \mathbf{V})_{ik} \mathbf{W}_{ik}}{\mathbf{W}'_{ik}} + \beta \frac{(\mathbf{T}^T \mathbf{V})_{ik} \mathbf{W}_{ik}}{\mathbf{W}'_{ik}} = 0. \end{aligned} \quad (24)$$

With straight algebra, (24) can be reduced to

$$\mathbf{W}_{ik} = \mathbf{W}'_{ik} \sqrt{\frac{(\mathbf{P}_A \mathbf{V})_{ik}}{(\mathbf{P}_A \mathbf{W}' \mathbf{V}^T \mathbf{V})_{ik} + \beta (\mathbf{T}^T \mathbf{V})_{ik} + \alpha \mathbf{W}'_{ik}}}. \quad (25)$$

Let $\mathbf{W}' = \mathbf{W}^{(t)}$ and $\mathbf{W} = \mathbf{W}^{(t+1)}$, then (25) falls back to (10). Thus, (10) essentially follows (18). According to Definition 1 and Proposition 1, we may conclude the proof. \square

5.2. Correctness and Convergence of (11)

For the updating rule of (11), we present the following two main results: 1) When convergent, the limiting solution of (11) satisfies the KKT condition. 2) The iteration of (11) converges. We formally establish the above results in Theorems 5 and 6, respectively. Before presenting these results, we first introduce the following Lagrangian function.

Fixing \mathbf{W} and \mathbf{T} , we need to solve (26) for \mathbf{V} :

$$\begin{aligned} \min_{\mathbf{V}} \frac{1}{2} \text{Tr}(-2\mathbf{P}_A \mathbf{W} \mathbf{V} \mathbf{V}^T + \mathbf{V} \mathbf{W}^T \mathbf{P}_A \mathbf{W} \mathbf{V} \mathbf{V}^T) \\ + \beta \text{Tr}(\mathbf{W}^T \mathbf{T}^T \mathbf{V}), \quad s.t. \mathbf{V} \geq 0, \mathbf{V}^T \mathbf{V} = \mathbf{I}_c. \end{aligned} \quad (26)$$

We introduce the symmetric Lagrangian multipliers $\Theta \in \mathbb{R}^{c \times c}$ and obtain the Lagrangian function to be minimized:

$$\begin{aligned} L_{\mathbf{V}} &= \frac{1}{2} \text{Tr}(-2\mathbf{P}_A \mathbf{W} \mathbf{V} \mathbf{V}^T + \mathbf{V} \mathbf{W}^T \mathbf{P}_A \mathbf{W} \mathbf{V} \mathbf{V}^T) \\ &+ \beta \text{Tr}(\mathbf{W}^T \mathbf{T}^T \mathbf{V}) + \frac{1}{2} \text{Tr}(\Theta (\mathbf{V}^T \mathbf{V} - \mathbf{I}_c)) \\ &= \text{Tr}(-\mathbf{P}_A \mathbf{W} \mathbf{V} \mathbf{V}^T + \frac{1}{2} \mathbf{V} (\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta) \mathbf{V}^T) \\ &- \frac{1}{2} \text{Tr}(\Theta), \end{aligned} \quad (27)$$

where for an arbitrary input matrix \mathbf{M} , $(\cdot)^+$ and $(\cdot)^-$ return matrices such that $\mathbf{M}^+, \mathbf{M}^- \geq 0$ and $\mathbf{M}^+ - \mathbf{M}^- = \mathbf{M}$. The correctness and convergence of (11) are ensured by the following Theorems 5 and 6.

Theorem 5. Fixing \mathbf{V} and \mathbf{T} , the limiting solution of (28) satisfies the KKT complementarity condition of (27):

$$\mathbf{V}_{ik} \leftarrow \mathbf{V}_{ik} \sqrt{\frac{(\mathbf{P}_A \mathbf{W})_{ik} + (\mathbf{V} (\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^-)_{ik}}{\beta (\mathbf{T} \mathbf{W})_{ik} + (\mathbf{V} (\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^+)_{ik}}}. \quad (28)$$

Proof. The gradient of $L_{\mathbf{V}}$ is

$$\frac{\partial L_{\mathbf{V}}}{\partial \mathbf{V}} = -\mathbf{P}_A \mathbf{W} + \mathbf{V} \mathbf{W}^T \mathbf{P}_A \mathbf{W} + \beta \mathbf{T} \mathbf{W} + \mathbf{V} \Theta. \quad (29)$$

Then the KKT complementarity condition gives

$$(-\mathbf{P}_A \mathbf{W} + \mathbf{V} \mathbf{W}^T \mathbf{P}_A \mathbf{W} + \beta \mathbf{T} \mathbf{W} + \mathbf{V} \Theta)_{ik} \mathbf{V}_{ik} = 0, \quad (30)$$

which provides a fixed point condition that the limiting solution should satisfy. From the updating rule of (28), it is seen that \mathbf{V} satisfies the following condition:

$$\begin{aligned} (-\mathbf{P}_A \mathbf{W} - \mathbf{V} (\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^- \\ + \mathbf{V} (\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^+ + \beta \mathbf{T} \mathbf{W})_{ik} \cdot \mathbf{V}_{ik}^2 = 0. \end{aligned} \quad (31)$$

Moreover, since

$$\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta = (\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^+ - (\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^-,$$

the condition in (31) can be further reduced to

$$(-\mathbf{P}_A \mathbf{W} + \mathbf{V} \mathbf{W}^T \mathbf{P}_A \mathbf{W} + \beta \mathbf{T} \mathbf{W} + \mathbf{V} \Theta)_{ik} \mathbf{V}_{ik}^2 = 0. \quad (32)$$

It is clear that (30) and (32) are identical because they imply the same conditions that either $\mathbf{V}_{ik} = 0$ holds in both (30) and (32) or $(-\mathbf{P}_A \mathbf{W} + \mathbf{V} \mathbf{W}^T \mathbf{P}_A \mathbf{W} + \beta \mathbf{T} \mathbf{W} + \mathbf{V} \Theta)_{ik} = 0$ holds in both (30) and (32), which concludes the proof. \square

Theorem 6. Fixing \mathbf{W} and \mathbf{T} , the limiting solution of (11) satisfies the KKT complementary condition, and the Lagrangian function $L_{\mathbf{V}}$ monotonically decreases under (11).

Proof. For ease of notation, we define $J(\mathbf{V}) = L_{\mathbf{V}}$. Similarly to the proof of Theorem 4, we obtain the following inequalities, which provide some proper upper bounds to the first four terms of $J(\mathbf{V})$ in (27):

$$\text{Tr}(\mathbf{V}^T \mathbf{T} \mathbf{W}) \leq \sum_{ik} (\mathbf{T} \mathbf{W})_{ik} \frac{\mathbf{V}_{ik}^2 + \mathbf{V}'_{ik}{}^2}{2\mathbf{V}'_{ik}}, \quad (33)$$

$$\begin{aligned} \text{Tr}(\mathbf{V}(\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^+ \mathbf{V}^T) \\ \leq \sum_{ik} \frac{(\mathbf{V}'(\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^+)_{ik} \mathbf{V}_{ik}^2}{\mathbf{V}'_{ik}} \end{aligned} \quad (34)$$

$$\text{Tr}(\mathbf{V}^T \mathbf{P}_A \mathbf{W}) \geq \sum_{ik} (\mathbf{P}_A \mathbf{W})_{ik} \mathbf{V}'_{ik} \left(1 + \log \frac{\mathbf{V}_{ik}}{\mathbf{V}'_{ik}}\right), \quad (35)$$

$$\begin{aligned} \text{Tr}(\mathbf{V}(\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^- \mathbf{V}^T) \\ \geq \sum_{ikl} (\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^-_{kl} \mathbf{V}'_{ik} \mathbf{V}'_{il} \left(1 + \log \frac{\mathbf{V}_{ik} \mathbf{V}_{il}}{\mathbf{V}'_{ik} \mathbf{V}'_{il}}\right), \end{aligned} \quad (36)$$

where the equal sign holds when $\mathbf{V} = \mathbf{V}' \geq 0$. Combining the above inequalities with proper weights, we construct the following auxiliary function for $J(\mathbf{V})$:

$$\begin{aligned} \bar{J}(\mathbf{V}, \mathbf{V}') &= -\frac{1}{2} \text{Tr}(\Theta) - \sum_{ik} (\mathbf{P}_A \mathbf{W})_{ik} \mathbf{V}'_{ik} \left(1 + \log \frac{\mathbf{V}_{ik}}{\mathbf{V}'_{ik}}\right) \\ &+ \beta \sum_{ik} (\mathbf{T} \mathbf{W})_{ik} \frac{\mathbf{V}_{ik}^2 + \mathbf{V}'_{ik}{}^2}{2\mathbf{V}'_{ik}} \\ &+ \frac{1}{2} \sum_{ik} \frac{(\mathbf{V}'(\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^+)_{ik} \mathbf{V}_{ik}^2}{\mathbf{V}'_{ik}} \\ &- \frac{1}{2} \sum_{ikl} (\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^-_{kl} \mathbf{V}'_{ik} \mathbf{V}'_{il} \left(1 + \log \frac{\mathbf{V}_{ik} \mathbf{V}_{il}}{\mathbf{V}'_{ik} \mathbf{V}'_{il}}\right), \end{aligned} \quad (37)$$

where $\bar{J}(\mathbf{V}, \mathbf{V}') \geq J(\mathbf{V})$ and $\bar{J}(\mathbf{V}, \mathbf{V}) = J(\mathbf{V})$ hold for any $\mathbf{V}, \mathbf{V}' \geq 0$. It is straightforward to check that $\bar{J}(\mathbf{V}, \mathbf{V}')$ is convex in \mathbf{V} . According to the first-order optimality condition of $\bar{J}(\mathbf{V}, \mathbf{V}')$:

$$\begin{aligned} \frac{\partial \bar{J}(\mathbf{V}, \mathbf{V}')}{\partial \mathbf{V}_{ik}} &= -\frac{(\mathbf{P}_A \mathbf{W})_{ik} \mathbf{V}'_{ik}}{\mathbf{V}_{ik}} + \beta \frac{(\mathbf{T} \mathbf{W})_{ik}}{\mathbf{V}'_{ik}} \mathbf{V}_{ik} \\ &+ \frac{(\mathbf{V}'(\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^+)_{ik} \mathbf{V}_{ik}}{\mathbf{V}'_{ik}} \\ &- \frac{(\mathbf{V}'(\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^-)_{ik} \mathbf{V}'_{ik}}{\mathbf{V}_{ik}} = 0, \end{aligned} \quad (38)$$

we obtain the global minimum of $\bar{J}(\mathbf{V}, \mathbf{V}')$:

$$\mathbf{V}_{ik} = \mathbf{V}'_{ik} \sqrt{\frac{(\mathbf{P}_A \mathbf{W})_{ik} + (\mathbf{V}'(\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^-)_{ik}}{\beta(\mathbf{T} \mathbf{W})_{ik} + (\mathbf{V}'(\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^+)_{ik}}}. \quad (39)$$

Let $\mathbf{V}^{(t)} = \mathbf{V}'$ and $\mathbf{V}^{(t+1)} = \mathbf{V}$, we may recover (28) from (39). By the first-order optimality condition of $L_{\mathbf{V}}$, we have

$$\begin{aligned} \mathbf{V}^T (-\mathbf{P}_A \mathbf{W} + \mathbf{V} \mathbf{W}^T \mathbf{P}_A \mathbf{W} + \beta \mathbf{T} \mathbf{W} + \mathbf{V} \Theta) \\ = -\mathbf{V}^T \mathbf{P}_A \mathbf{W} + \mathbf{W}^T \mathbf{P}_A \mathbf{W} + \beta \mathbf{V}^T \mathbf{T} \mathbf{W} + \Theta = 0, \end{aligned} \quad (40)$$

which implies that $\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta = \mathbf{V}^T \mathbf{P}_A \mathbf{W} - \beta \mathbf{V}^T \mathbf{T} \mathbf{W}$. Thus, it is proper to define $(\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^+ = \mathbf{V}^T \mathbf{P}_A \mathbf{W}$, $(\mathbf{W}^T \mathbf{P}_A \mathbf{W} + \Theta)^- = \beta \mathbf{V}^T \mathbf{T} \mathbf{W}$. Then, it is straightforward to obtain (11) by substituting the above definitions into (28), which concludes the proof. \square

5.3. Convergence of Overall Algorithm

The convergence of the overall algorithm is guaranteed by the following Theorem 7.

Theorem 7. The objective value sequence of (7) is monotonically decreasing under the updating rules of (8), (10) and (11), and thus converges.

Proof. Define the objective function as $f(\mathbf{F}, \mathbf{W}, \mathbf{V})$ and denote the iteration by superscript $(\cdot)^{(t)}$, respectively, then it is clear that the following chain of inequalities hold:

$$\begin{aligned} f(\mathbf{F}^{(t)}, \mathbf{W}^{(t)}, \mathbf{V}^{(t)}) \\ \geq f(\mathbf{F}^{(t+1)}, \mathbf{W}^{(t)}, \mathbf{V}^{(t)}) \geq f(\mathbf{F}^{(t+1)}, \mathbf{W}^{(t+1)}, \mathbf{V}^{(t+1)}). \end{aligned} \quad (41)$$

Thus, the sequence $\{f(\mathbf{F}^{(t)}, \mathbf{W}^{(t)}, \mathbf{V}^{(t)})\}_{t=0}^{\infty}$ is monotonically decreasing. Moreover, since \mathbf{P}_A is symmetric and p.s.d., there exists $\mathbf{C} \in \mathbb{R}^{n \times n}$ such that $\mathbf{P}_A = \mathbf{C}^T \mathbf{C}$. Then, it is clear that

$$\begin{aligned} f(\mathbf{F}^{(t)}, \mathbf{W}^{(t)}, \mathbf{V}^{(t)}) \\ = \frac{1}{2} \|\mathbf{C} - \mathbf{C} \mathbf{W}^{(t)} (\mathbf{V}^{(t)})^T\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}^{(t)}\|_F^2 \\ + \beta \text{Tr}(\mathbf{F}^{(t)T} \mathbf{L}^{(t)} \mathbf{F}^{(t)}) - \text{Tr}(\mathbf{C}^T \mathbf{C}) \geq -\text{Tr}(\mathbf{C}^T \mathbf{C}). \end{aligned} \quad (42)$$

Therefore, the sequence $\{f(\mathbf{F}^{(t)}, \mathbf{W}^{(t)}, \mathbf{V}^{(t)})\}_{t=0}^{\infty}$ is decreasing and lower-bounded, and thus converges. \square

6. Experiments

In this section, we conduct extensive experiments to verify the effectiveness of the Figer-CF. We compare our method with nine state-of-the-art clustering methods, including the WNMf [14], ONMF [5], CNMF [6], KNMF [6], RMNMF [10], OPMC [18], MKKM-SR [23], EWNMF [47], and GLS-MKNMF [32], on four widely used benchmark data sets, including the Jaffe, PIX, Semeion, and COIL20. Three metrics, including the clustering accuracy (ACC), normalized mutual information (NMI), and

best performance in almost all cases, with improvements, in general, by about 2-3% in ACC and PUR and 6-7% in NMI, respectively, compared with the top second method. On the COIL20 data, methods such as the EWNMF and GLS-MKNMF have relatively better performance among the baselines. However, both of them are inferior to the Figer-CF and it is difficult to tell which of them is better. In general, the Semeion and COIL20 data appear to be more difficult, since they contain pairs of clusters such as {"0", "O"} and {"2", "Z"}, which are quite similar and confusing. In this case, it may help to reveal the underlying relationships of the data by seeking the local geometric structures of the data since the global relationships may provide misleading information. All these observations confirm the effectiveness and superiority of the Figer-CF.

6.3. Convergence Study

To empirically verify the convergent property of our method, we show curves of the objective value sequences on the Semeion and COIL20 data sets in Fig. 1, where the parameters are set according to Section 6.1. It is seen that the objective value sequences of both data generally converge within about 50 iterations, which confirms the effectiveness and efficiency of the optimization algorithm.

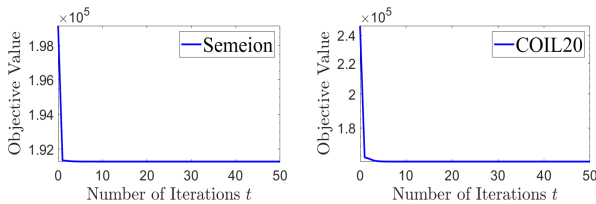


Figure 1. Convergence curve of the objective value sequence.

6.4. Effects of High-order Neighbors

In this test, we empirically show the effects of high-order neighbors using the COIL20 and Semeion data. In particular, we show the performance of the Figer-CF using \mathbf{P}_A of orders from 1 to 3, respectively, and report the performance in Fig. 2, where the other parameters are tuned to the best accordingly. It is observed that the Figer-CF has significantly improved performance in all metrics when a higher-order \mathbf{P}_A is used, which confirms the effectiveness of exploiting comprehensive and complementary neighbor information. If the order tends to be even higher, the improvement tends to be tight. Considering the difficulty in searching a proper A for the Figer-CF in practical sense, it is effective and fairly convincing to suggest $A = 3$ in real-world applications for both effectiveness and efficiency.

6.5. Parameter Sensitivity

We test how the balancing parameters affect the performance of the Figer-CF. Without loss of generality and due

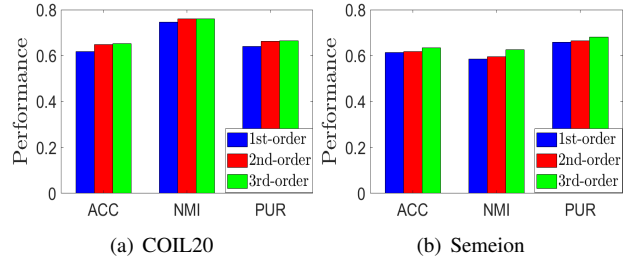


Figure 2. Effects of the high-order neighbors on the clustering capability of Figer-CF.

to space limit, we report the results on the Jaffe data. We report the performance of the Figer-CF with respect to all combinations of $\{\alpha, \beta\} \in \mathbb{S} \times \mathbb{S}$ in Fig. 3. It is seen that the Figer-CF generally has high performance in all metrics with a broad range of parameter values. Similar observations can be found on other data sets as well, which confirms that the Figer-CF is rather insensitive to parameters. This property is desired by unsupervised learning and is essentially important for potential applicability in real applications.

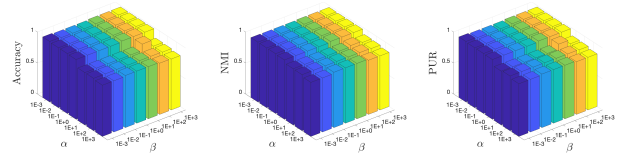


Figure 3. Performance changes with respect to different combinations of parameters on the Jaffe data.

7. Conclusion

In this paper, we propose a novel concept factorization method that uses the cross-order neighbor information of the data to learn score and coefficient matrices with bipartite graph partitioning, which fully exploits comprehensive and complementary neighbor information, as well as explicit cluster structure of the data. Our algorithm admits both elegant theoretical guarantee for convergent property and promising experimental performance in clustering, which confirms the effectiveness and efficiency of the new method.

Acknowledgement

Yongyong Chen and Chenglizhao Chen are corresponding authors. This work is supported in part by the National Natural Science Foundation of China (NSFC) under Grants 62276147, 62172246, and 62106063; in part by the Shandong Province Colleges and Universities Youth Innovation Technology Plan Innovation Team Project under Grants 2022KJ149, 2021KJ062, and 2020KJN011; and in part by the Guangdong Major Project of Basic and Applied Basic Research under Grant 2023B0303000010.

References

- [1] Deng Cai, Xiaofei He, and Jiawei Han. Locally consistent concept factorization for document clustering. *IEEE Transactions on Knowledge & Data Engineering*, 23(6):902–913, 2011. [2](#)
- [2] Deng Cai, Xiaofei He, Jiawei Han, and Thomas S Huang. Graph regularized nonnegative matrix factorization for data representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(8):1548–1560, 2011. [1](#), [2](#), [3](#), [4](#), [7](#)
- [3] Shuai Chang, Jie Hu, Tianrui Li, Hao Wang, and Bo Peng. Multi-view clustering via deep concept factorization. *Knowledge-Based Systems*, 217:106807, 2021. [1](#)
- [4] Mulin Chen and Xuelong Li. Concept factorization with local centroids. *IEEE Transactions on Neural Networks and Learning Systems*, 32(11):5247–5253, 2020. [2](#)
- [5] Chris Ding, Tao Li, Wei Peng, and Haesun Park. Orthogonal nonnegative matrix t-factorizations for clustering. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 126–135. ACM, 2006. [6](#)
- [6] Chris HQ Ding, Tao Li, and Michael I Jordan. Convex and semi-nonnegative matrix factorizations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(1):45–55, 2008. [1](#), [2](#), [4](#), [6](#)
- [7] Ky Fan. On a theorem of weyl concerning eigenvalues of linear transformations i. *Proceedings of the National Academy of Sciences*, 35(11):652–655, 1949. [2](#)
- [8] Xin-Ru Feng, Heng-Chao Li, Rui Wang, Qian Du, Xiuping Jia, and Antonio Plaza. Hyperspectral unmixing based on nonnegative matrix factorization: A comprehensive review. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 15:4414–4436, 2022. [1](#)
- [9] Baoxiang Huang, Linyao Ge, Ge Chen, Milena Radenkovic, Xiaopeng Wang, Jinming Duan, and Zhenkuan Pan. Nonlocal graph theory based transductive learning for hyperspectral image classification. *Pattern Recognition*, 116:107967, 2021.
- [10] Jin Huang, Feiping Nie, Heng Huang, and Chris Ding. Robust manifold nonnegative matrix factorization. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 8(3):11, 2014. [6](#)
- [11] Yuheng Jia, Hui Liu, Junhui Hou, and Sam Kwong. Semisupervised adaptive symmetric non-negative matrix factorization. *IEEE transactions on cybernetics*, 51(5):2550–2562, 2020. [1](#)
- [12] Kehan Kang, Chenglizhao Chen, and Chong Peng. Consensus low-rank multi-view subspace clustering with cross-view diversity preserving. *IEEE Signal Processing Letters*, 30:1512–1516, 2023. [1](#)
- [13] Zhao Kang, Zhanyu Liu, Shirui Pan, and Ling Tian. Fine-grained attributed graph clustering. In *Proceedings of the 2022 SIAM International Conference on Data Mining (SDM)*, pages 370–378. SIAM, 2022. [2](#)
- [14] Yong-Deok Kim and Seungjin Choi. Weighted non-negative matrix factorization. In *Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on*, pages 1541–1544. IEEE, 2009. [6](#)
- [15] Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999. [1](#)
- [16] Daniel D Lee and H Sebastian Seung. Algorithms for non-negative matrix factorization. In *Advances in neural information processing systems*, pages 556–562, 2001. [1](#), [4](#)
- [17] Jiyuan Liu, Xinwang Liu, Jian Xiong, Qing Liao, Sihang Zhou, Siwei Wang, and Yuexiang Yang. Optimal neighborhood multiple kernel clustering with adaptive local kernels. *IEEE Transactions on Knowledge and Data Engineering*, 34(6):2872–2885, 2022. [1](#)
- [18] Jiyuan Liu, Xinwang Liu, Yuexiang Yang, Li Liu, Siqi Wang, Weixuan Liang, and Jiangyong Shi. One-pass multi-view clustering for large-scale data. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 12324–12333, 2021. [6](#)
- [19] Kai Liu, Xiangyu Li, Zhihui Zhu, Lodewijk Brand, and Hua Wang. Factor-bounded nonnegative matrix factorization. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 15(6):1–18, 2021. [1](#)
- [20] Xinwang Liu, Lei Wang, Jian Zhang, Jianping Yin, and Huan Liu. Global and local structure preservation for feature selection. *IEEE Transactions on Neural Networks and Learning Systems*, 25(6):1083–1095, 2014. [1](#)
- [21] Xinwang Liu, Sihang Zhou, Li Liu, Chang Tang, Siwei Wang, Jiyuan Liu, and Yi Zhang. Localized simple multiple kernel k-means. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9273–9281, 2021. [1](#)
- [22] Canyi Lu, Jiashi Feng, Zhouchen Lin, Tao Mei, and Shuicheng Yan. Subspace clustering by block diagonal representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(2):487–501, 2019. [2](#)
- [23] Jitao Lu, Yihang Lu, Rong Wang, Feiping Nie, and Xuelong Li. Multiple kernel k-means clustering with

- simultaneous spectral rotation. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4143–4147. IEEE, 2022. 6, 7
- [24] Yuwu Lu, Chun Yuan, Wenwu Zhu, and Xuelong Li. Structurally incoherent low-rank nonnegative matrix factorization for image classification. *IEEE Transactions on Image Processing*, 27(11):5248–5260, 2018. 1
- [25] Jiaqi Ma, Yipeng Zhang, and Lefei Zhang. Discriminative subspace matrix factorization for multiview data clustering. *Pattern Recognition*, 111:107676, 2021. 1
- [26] Xiaoke Ma, Penggang Sun, and Guimin Qin. Non-negative matrix factorization algorithms for link prediction in temporal networks using graph communicability. *Pattern Recognition*, 71:361 – 374, 2017. 1
- [27] Xiaoke Ma, Benhui Zhang, Changzhou Ma, and Zhiyu Ma. Co-regularized nonnegative matrix factorization for evolving community detection in dynamic networks. *Information Sciences*, 528:265–279, 2020. 1
- [28] Yang Meng, Ronghua Shang, Fanhua Shang, Licheng Jiao, Shuyuan Yang, and Rustam Stolkin. Semi-supervised graph regularized deep nmf with bi-orthogonal constraints for data representation. *IEEE transactions on neural networks and learning systems*, 31(9):3245–3258, 2019. 1
- [29] Yulong Pei, Nilanjan Chakraborty, and Katia Sycara. Nonnegative matrix tri-factorization with graph regularization for community detection in social networks. In *Proceedings of the 24th International Conference on Artificial Intelligence*, pages 2083–2089, 2015. 1
- [30] Chong Peng, Chenglizhao Chen, Zhao Kang, Jianbo Li, and Qiang Cheng. Res-pca: A scalable approach to recovering low-rank matrices. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7309–7317, 2019. 1
- [31] Chong Peng and Qiang Cheng. Discriminative ridge machine: A classifier for high-dimensional data or imbalanced data. *IEEE Transactions on Neural Networks and Learning Systems*, 32(6):2595–2609, 2021. 2, 3
- [32] Chong Peng, Xingrong Hou, Yongyong Chen, Zhao Kang, Chenglizhao Chen, and Qiang Cheng. Global and local similarity learning in multi-kernel space for nonnegative matrix factorization. *Knowledge-Based Systems*, 279:110946, 2023. 1, 2, 3, 6, 7
- [33] Chong Peng, Zhao Kang, Huiqing Li, and Qiang Cheng. Subspace clustering using log-determinant rank approximation. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 925–934. ACM, 2015. 2
- [34] Chong Peng, Yang Liu, Kehan Kang, Yongyong Chen, Xinxing Wu, Andrew Cheng, Zhao Kang, Chenglizhao Chen, and Qiang Cheng. Hyperspectral image denoising using nonconvex local low-rank and sparse separation with spatial-spectral total variation regularization. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–17, 2022.
- [35] Chong Peng, Qian Zhang, Zhao Kang, Chenglizhao Chen, and Qiang Cheng. Kernel two-dimensional ridge regression for subspace clustering. *Pattern Recognition*, 113:107749, 2021. 2
- [36] Chong Peng, Yiqun Zhang, Yongyong Chen, Zhao Kang, Chenglizhao Chen, and Qiang Cheng. Log-based sparse nonnegative matrix factorization for data representation. *Knowledge-Based Systems*, 251:109127, 2022. 1, 2, 7
- [37] Chong Peng, Zhilu Zhang, Chenglizhao Chen, Zhao Kang, and Qiang Cheng. Two-dimensional semi-nonnegative matrix factorization for clustering. *Information Sciences*, 590:106–141, 2022. 1
- [38] Chong Peng, Zhilu Zhang, Zhao Kang, Chenglizhao Chen, and Qiang Cheng. Nonnegative matrix factorization with local similarity learning. *Information Sciences*, page 325346, Jul 2021. 1, 2, 3
- [39] Siyuan Peng, Wee Ser, Badong Chen, and Zhiping Lin. Robust semi-supervised nonnegative matrix factorization for image clustering. *Pattern Recognition*, 111:107683, 2021. 1
- [40] Xi Peng, Zhang Yi, and Huajin Tang. Robust subspace clustering via thresholding ridge regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 2015. 2
- [41] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 22(8):888–905, 2000. 3
- [42] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-scale information network embedding. In *Proceedings of the 24th international conference on world wide web*, pages 1067–1077, 2015. 2, 3
- [43] Ming Tong, Yiran Chen, Mengao Zhao, Haili Bu, and Shengnan Xi. A deep discriminative and robust nonnegative matrix factorization network method with soft label constraint. *Neural Computing and Applications*, 31:7447–7475, 2019. 1
- [44] Fusheng Wang, Chenglizhao Chen, and Chong Peng. Essential low-rank sample learning for group-aware

- subspace clustering. *IEEE Signal Processing Letters*, 30:1537–1541, 2023. [1](#)
- [45] Haixian Wang, Sibao Chen, Zilan Hu, and Wenming Zheng. Locality-preserved maximum information projection. *IEEE Transactions on Neural Networks*, 19(4):571–585, 2008. [1](#)
- [46] Yu-Xiang Wang, Huan Xu, and Chenlei Leng. Provable subspace clustering: When lrr meets ssc. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. [2](#)
- [47] Jiao Wei, Can Tong, Bingxue Wu, Qiang He, Shouliang Qi, Yudong Yao, and Yueyang Teng. An entropy weighted nonnegative matrix factorization algorithm for feature representation. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–11, 2022. [6](#)
- [48] Svante Wold, Kim Esbensen, and Paul Geladi. Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3):37–52, 1987. [1](#)
- [49] Jianlong Wu, Zhouchen Lin, and Hongbin Zha. Essential tensor learning for multi-view spectral clustering. *IEEE Transactions on Image Processing*, 28(12):5910–5922, 2019. [3](#)
- [50] Wei Xu and Yihong Gong. Document clustering by concept factorization. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 202–209, 2004. [1](#), [2](#)
- [51] Yang Zhao, Huiyang Wang, and Jihong Pei. Deep non-negative matrix factorization architecture based on underlying basis images learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(6):1897–1913, 2021. [1](#)
- [52] Dengyong Zhou, Jiayuan Huang, and Bernhard Schölkopf. Learning from labeled and unlabeled data on a directed graph. In *Proceedings of the 22nd international conference on Machine learning*, pages 1036–1043, 2005. [3](#)