

Doubly Aligned Incomplete Multi-view Clustering

Menglei Hu and Songcan Chen*

College of Computer Science & Technology, Nanjing University of Aeronautics & Astronautics
Collaborative Innovation Center of Novel Software Technology and Industrialization
{ml.hu, s.chen}@nuaa.edu.cn

Abstract

Nowadays, multi-view clustering has attracted more and more attention. To date, almost all the previous studies assume that views are complete. However, in reality, it is often the case that each view may contain some missing instances. Such incompleteness makes it impossible to directly use traditional multi-view clustering methods. In this paper, we propose a Doubly Aligned Incomplete Multi-view Clustering algorithm (DAIMC) based on weighted semi-nonnegative matrix factorization (semi-NMF). Specifically, on the one hand, DAIMC utilizes the given instance alignment information to learn a common latent feature matrix for all the views. On the other hand, DAIMC establishes a consensus basis matrix with the help of $L_{2,1}$ -Norm regularized regression for reducing the influence of missing instances. Consequently, compared with existing methods, besides inheriting the strength of semi-NMF with ability to handle negative entries, DAIMC has two unique advantages: 1) solving the incomplete view problem by introducing a respective weight matrix for each view, making it able to easily adapt to the case with more than two views; 2) reducing the influence of view incompleteness on clustering by enforcing the basis matrices of individual views being aligned with the help of regression. Experiments on four real-world datasets demonstrate its advantages.

1 Introduction

Many datasets in real world naturally appear in multiple views or come from multiple sources [Blum and Mitchell, 1998; Schechter *et al.*, 2017], which are called multi-view data. For example, a document can be translated into different languages, and images can be described by different features such as Fourier shape descriptors and K-L expansion coefficients. In multi-view data, these different views often share some consistency and complementary information [Sun, 2013; Zhao *et al.*, 2017]. Such information can be beneficial to learning tasks such as classification and clustering. This leads to a surge of interest in multi-view learning

[Potthast *et al.*, 2018; Jing *et al.*, 2017], whose goal is to integrate information and give a compatible solution across all views. Nowadays, multi-view learning has been widely studied in different areas such as face recognition, image processing and natural language processing [Romero *et al.*, 2017; Xing *et al.*, 2017; Nie *et al.*, 2018].

In all the tasks of multi-view learning, multi-view clustering [Bickel and Scheffer, 2004; Fan *et al.*, 2017; Nie *et al.*, 2017; Chao *et al.*, 2017] has attracted more and more attentions due to exempting the expensive requirement of data labeling. The goal of multi-view clustering is making full use of multi-view data to get a better clustering result than just simply concatenated views. To date, many multi-view clustering methods have been proposed. Among these methods, one of the most widely used techniques is nonnegative matrix factorization (NMF) [Wang *et al.*, 2016b; Li, 2016]. [Lee and Seung, 1999] proposes the NMF, which has received much attention because of its straightforward interpretability for applications. Then, some researchers utilize the NMF for multi-view learning, especially multi-view clustering. A joint NMF process with the consistency constraint is formulated in [Liu *et al.*, 2013], which performs NMF for each view and pushes each view's low dimensional representation towards a common consensus. Besides, some researchers have integrated manifold learning and multi-view learning by imposing the manifold regularization on the objective function of NMF respectively for individual views data [Wang *et al.*, 2016a; Zong *et al.*, 2017].

Most of the previous studies on multi-view clustering make a common assumption that all of the views are complete. However, in real world applications, multi-view data tend to be incomplete. For example, in the camera network, for some reasons, the camera may temporarily fail, or be blocked by some objects, making the instance missing. Another example is in document clustering, different languages of the documents represent multiple views. However, we may not get all the documents translated into each language. All the above-mentioned cases lead to the incompleteness of multi-view data. As a result, the lack of the whole row or column makes the traditional instance imputation methods fail. So how to make full use of the complementary knowledge hidden in different views and reduce the impact of missing instances are the most challenging problems of incomplete multi-view learning. Recently, some incomplete multi-view clustering meth-

ods have been proposed, for example, [Li *et al.*, 2014] proposes PVC by utilizing the information of instance alignment to learn a common latent subspace for aligned instances and a private latent representation for unaligned instances via NMF. Borrowing this idea, [Zhao *et al.*, 2016] proposes IMG by integrating PVC and manifold learning to learn the global structure over the entire data instances across all views. However, PVC and IMG can only deal with the two-view incomplete multi-view clustering, limiting their application scope. A method for clustering more than two incomplete views is proposed in [Shao *et al.*, 2015](MIC) by filling the missing instances with the average feature values in each incomplete view, then handling the problem with the help of weighted NMF and $L_{2,1}$ -Norm regularization [Kong *et al.*, 2011; Wu *et al.*, 2018]. However, such a simple imputation will cause a great deviation when the missing ratio is large. As a result, incomplete multi-view clustering still faces significant challenges.

In this paper, we propose the Doubly Aligned Incomplete Multi-view Clustering (DAIMC) to meet the challenges. By integrating semi-NMF [Ding *et al.*, 2010] and $L_{2,1}$ -Norm regularized regression model, DAIMC tries to learn a common latent feature matrix for all the views from two aspects of instances aligned and the basis matrix aligned.

Compared with the existing methods around the incomplete multi-view clustering, besides inheriting the strength of semi-NMF with ability to handle negative entries, DAIMC mainly has the following advantages:

1. DAIMC extends PVC. Borrowing the idea of the weighted NMF, DAIMC introduces a respective weight matrix for each incomplete view to assign the missing instances zero weights and the presented instances one weights in each view, making it able to be easily and straightforwardly extended to the scenario with more than two incomplete views.
2. Besides the instance alignment, DAIMC considers the global information by enforcing the basis matrices of individual views being aligned with the help of the $L_{2,1}$ -Norm regularized regression model, which further reduces the influence of missing instances on clustering performance.
3. An iterative optimization algorithm for DAIMC with convergent guarantee is proposed. Experimental results on four real-world datasets demonstrate its advantages.

The rest of this paper is organized as follows. In Section 2, we overview some related work on semi-NMF and incomplete multi-view clustering. Section 3 proposes our DAIMC and an efficient iterative algorithm of solving it in detail, respectively. Extensive experimental results and analysis are reported in Section 4. Section 5 concludes this paper with future research directions.

2 Related Work

2.1 Semi-nonnegative Matrix Factorization

The semi-NMF [Ding *et al.*, 2010] is an effective latent factor learning method, which is the extension of NMF. Given an input data matrix $\mathbf{X} \in \mathbb{R}^{M \times N}$, each column of \mathbf{X} is an instance vector. The semi-NMF aims to find a matrix $\mathbf{U} \in \mathbb{R}^{M \times K}$ and a nonnegative matrix $\mathbf{V} \in \mathbb{R}^{N \times K}$ whose product can well approximate the original matrix \mathbf{X} . To fa-

cilitate discussion, we call \mathbf{U} the *basis matrix* and \mathbf{V} the *latent feature matrix*. Thus we can get the following minimization problem

$$\begin{aligned} \min \|\mathbf{X} - \mathbf{UV}^T\|_F^2 \\ \text{s.t. } \mathbf{V} \geq 0 \end{aligned} \quad (1)$$

Similar to NMF, the objective function in (1) is biconvex. Therefore, it is unrealistic to expect an algorithm to find the global minimum. [Ding *et al.*, 2010] proposes an iterative updating algorithm to find the locally optimal solution as follows:

Update \mathbf{U} (while fixing \mathbf{V}) using the rule

$$\mathbf{U} = \mathbf{XV}(\mathbf{V}^T\mathbf{V})^{-1} \quad (2)$$

Update \mathbf{V} (while fixing \mathbf{U}) using

$$\mathbf{V}_{jk} \leftarrow \mathbf{V}_{jk} \sqrt{\frac{(\mathbf{X}^T\mathbf{U})_{jk}^+ + [\mathbf{V}(\mathbf{U}^T\mathbf{U})^-]_{jk}}{(\mathbf{X}^T\mathbf{U})_{jk}^- + [\mathbf{V}(\mathbf{U}^T\mathbf{U})^+]_{jk}}} \quad (3)$$

where we separate the positive and negative parts of a matrix \mathbf{A} as:

$$\mathbf{A}_{jk}^+ = (|\mathbf{A}_{jk}| + \mathbf{A}_{jk})/2, \quad \mathbf{A}_{jk}^- = (|\mathbf{A}_{jk}| - \mathbf{A}_{jk})/2. \quad (4)$$

2.2 Incomplete Multi-view Clustering

Given a dataset with N instances, C categories, n_v views $\{\mathbf{X}^{(i)}, i = 1, 2, \dots, n_v\}$, where $\mathbf{X}^{(v)} \in \mathbb{R}^{d_v \times N}$ is the v -th view of the dataset. An indicator matrix $\mathbf{M} \in \mathbb{R}^{n_v \times N}$ for incomplete multi-view clustering problem is defined as:

$$\mathbf{M}_{ij} = \begin{cases} 1 & \text{if } j\text{-th instance is in the } i\text{-th view} \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

where each row of \mathbf{M} is the instance presence for corresponding view. If every view contains all the instances, then the matrix \mathbf{M} is an all one matrix. And if the v -th view is incomplete, the data matrix $\mathbf{X}^{(v)}$ will have a number of column missing, *i.e.*, $\sum_{j=1}^N \mathbf{M}_{vj} < N$.

The aim of the incomplete multi-view clustering is to integrate all the incomplete views to cluster the N instances into C clusters.

3 Proposed Approach

In this section, we present our Doubly Aligned Incomplete Multi-view Clustering(DAIMC) in detail. We model the DAIMC as a joint weighted semi-NMF problem and use $L_{2,1}$ -norm regularized regression to enforce the basis matrix of individual views being aligned. In the following, we propose our model in two aspects and then give a unified objective function for implementing DAIMC.

3.1 Weighted Semi-NMF for Incomplete Multi-view Data

For the v -th view, similarly to the weighted NMF, the weighted semi-NMF factorizes the data matrix $\mathbf{X}^{(v)} \in \mathbb{R}^{d_v \times N}$ into two matrix $\mathbf{U}^{(v)}$ and $\mathbf{V}^{(v)}$, where $\mathbf{U}^{(v)} \in \mathbb{R}^{d_v \times K}$, $\mathbf{V}^{(v)} \in \mathbb{R}^{N \times K}$, while giving different weights to the reconstruction

errors of different instances. K denotes dimension of subspace. In the experiments of previous works [Shao *et al.*, 2015; Zhao *et al.*, 2016], for multi-view clustering, the K is set to the number of the categories of the data matrix $\mathbf{X}^{(v)}$, i.e., $K = C$. As a result, the weighted semi-NMF optimization problem is formulated as:

$$\begin{aligned} \min & \|(\mathbf{X}^{(v)} - \mathbf{U}^{(v)}\mathbf{V}^{(v)T})\mathbf{W}^{(v)}\|_F^2 \\ \text{s.t. } & \mathbf{V}^{(v)} \geq 0 \end{aligned} \quad (6)$$

where the weight matrix $\mathbf{W}^{(v)} \in \mathbf{R}^{N \times N}$ is a diagonal matrix.

$$\mathbf{W}_{jj}^{(i)} = \begin{cases} 1 & \text{if } j\text{-th instance is in the } i\text{-th view} \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

Note that $\mathbf{W}_{jj}^{(v)}$ indicates the weight of the j -th instance in view v . If the j -th instance is missing, then the loss of this instance will be ignored.

However, (6) only independently decomposes different views without considering their consistency information. To address this issue, we assume that different views have distinct basis matrices $\{\mathbf{U}^{(i)}\}_{i=1}^{n_v}$, but share the same latent feature space \mathbf{V} . As a result, (6) is rewritten as follows:

$$\begin{aligned} \min & \sum_{i=1}^{n_v} \|(\mathbf{X}^{(i)} - \mathbf{U}^{(i)}\mathbf{V}^T)\mathbf{W}^{(i)}\|_F^2 \\ \text{s.t. } & \mathbf{V} \geq 0 \end{aligned} \quad (8)$$

By solving (8), we can obtain a common representation \mathbf{V} for multiple incomplete-view instances.

3.2 $L_{2,1}$ -Norm Regularized Regression with Basis Matrix

To further reduce the influence of the missing instances, DAIMC attempts to incorporate the global information among views. In multi-view data, different views have different representations for a same data matrix. Thus, we can align the different basis matrices $\{\mathbf{U}^{(i)}\}_{i=1}^{n_v}$ of individual views with the help of regression by solving the following problem for the basis matrices intended to be aligned.

$$\min \sum_{i=1}^{n_v} \|\mathbf{B}^{(i)T} \mathbf{U}^{(i)} - \mathbf{U}^{(*)}\|_F^2 + \beta \|\mathbf{B}^{(i)}\|_{2,1} \quad (9)$$

where $\mathbf{B}^{(i)} \in \mathbb{R}^{d_i \times K}$ is the regression coefficient matrix for view i . The $L_{2,1}$ -norm regularization term is here introduced for ensuring $\mathbf{B}^{(i)}$ sparse in rows. In this way, $\mathbf{B}^{(i)}$ performs a feature selection during the alignment process. The matrix $\mathbf{U}^{(*)} \in \mathbb{R}^{P \times K}$ is the same low dimensional representation for the basis matrices of all the views. P denotes dimension of subspace. The value of P will affect the result. Instead of looking for an appropriate P , we simply set the matrix $\mathbf{U}^{(*)}$ equal to a K dimensional identity matrix \mathbf{I}_K , whose columns correspond to the cluster encodings. For such a setting, our experiments later confirm its effectiveness. As a result, (9) is rewritten as:

$$\min \sum_{i=1}^{n_v} \|\mathbf{B}^{(i)T} \mathbf{U}^{(i)} - \mathbf{I}\|_F^2 + \beta \|\mathbf{B}^{(i)}\|_{2,1} \quad (10)$$

where β is the trade-off hyper-parameter between sparsity and accuracy of regression for the i -th view, $\|\cdot\|_{2,1}$ is the $L_{2,1}$ norm and defined as:

$$\|\mathbf{B}^{(v)}\|_{2,1} = \sum_{i=1}^{d_v} \sqrt{\sum_{j=1}^K \mathbf{B}_{ij}^{(v)2}}$$

3.3 Unified Objective Function

Considering the objective for instance alignment information as well as the basis matrix alignment information simultaneously, we minimize the following objective function

$$\begin{aligned} \mathcal{J} &= \sum_{i=1}^{n_v} \left\{ \|(\mathbf{X}^{(i)} - \mathbf{U}^{(i)}\mathbf{V}^T)\mathbf{W}^{(i)}\|_F^2 \right. \\ &\quad \left. + \alpha (\|\mathbf{B}^{(i)T} \mathbf{U}^{(i)} - \mathbf{I}\|_F^2 + \beta \|\mathbf{B}^{(i)}\|_{2,1}) \right\} \\ \text{s.t. } & \mathbf{V} \geq 0 \end{aligned} \quad (11)$$

where α is nonnegative hyper-parameter that controls the trade-off between the aforementioned two objectives.

3.4 Optimization

The objective function in Eq.(11) is not convex over all variables $\{\mathbf{U}^{(i)}\}_{i=1}^{n_v}$, \mathbf{V} , $\{\mathbf{B}^{(i)}\}_{i=1}^{n_v}$, simultaneously. To solve this optimization problem, we propose an alternating iteration procedure.

Subproblem of $\{\mathbf{U}^{(i)}\}_{i=1}^{n_v}$. With $\{\mathbf{B}^{(i)}\}_{i=1}^{n_v}$ and \mathbf{V} fixed, for each $\mathbf{U}^{(i)}$, we need to minimize the following objective function:

$$\mathcal{J}(\mathbf{U}^{(i)}) = \|(\mathbf{X}^{(i)} - \mathbf{U}^{(i)}\mathbf{V}^T)\mathbf{W}^{(i)}\|_F^2 + \alpha \|\mathbf{B}^{(i)T} \mathbf{U}^{(i)} - \mathbf{I}\|_F^2 \quad (12)$$

The partial derivation of $\mathcal{J}(\mathbf{U}^{(i)})$ with respect to $\mathbf{U}^{(i)}$ is

$$\begin{aligned} \frac{\partial \mathcal{J}}{\partial \mathbf{U}^{(i)}} &= 2(\mathbf{U}^{(i)}\mathbf{V}^T - \mathbf{X}^{(i)})\mathbf{W}^{(i)}\mathbf{W}^{(i)T} \mathbf{V} + \\ &\quad 2\alpha \mathbf{B}^{(i)}(\mathbf{B}^{(i)T} \mathbf{U}^{(i)} - \mathbf{I}) \end{aligned} \quad (13)$$

From the definition of $\mathbf{W}^{(i)}$, we can see $\mathbf{W}^{(i)} = \mathbf{W}^{(i)}\mathbf{W}^{(i)T}$. Let $\partial \mathcal{J} / \partial \mathbf{U}^{(i)} = 0$, we get the following equation:

$$(\mathbf{U}^{(i)}\mathbf{V}^T - \mathbf{X}^{(i)})\mathbf{W}^{(i)}\mathbf{V} + \alpha \mathbf{B}^{(i)}(\mathbf{B}^{(i)T} \mathbf{U}^{(i)} - \mathbf{I}) = 0 \quad (14)$$

Eq.(14) is called the *continuous Sylvester equation* with respect to $\mathbf{U}^{(i)}$, which often arises in control theory. When both d_i and K are small, we can solve Eq.(14) via vectorization and get the following updating rule for $\mathbf{U}^{(i)}$:

$$\begin{aligned} \text{vec}(\mathbf{U}^{(i)}) &= [\mathbf{I}_K \otimes (\alpha \mathbf{B}^{(i)} \mathbf{B}^{(i)T}) + (\mathbf{V}^T \mathbf{W}^{(i)} \mathbf{V}) \otimes \mathbf{I}_{d_i}]^{-1} \\ &\quad \text{vec}(\mathbf{X}^{(i)} \mathbf{W}^{(i)} \mathbf{V} + \alpha \mathbf{B}^{(i)}) \end{aligned} \quad (15)$$

And when both d_i and K are large, we instead solve it via *Conjugate Gradient*. In our experiments, we use the *lyap* function of MATLAB to solve Eq.(14).

Subproblem of \mathbf{V} . With $\{\mathbf{B}^{(i)}\}_{i=1}^{n_v}$ and $\{\mathbf{U}^{(i)}\}_{i=1}^{n_v}$ fixed, we need to minimize the following objective function:

$$\begin{aligned} \mathcal{J}(\mathbf{V}) &= \sum_{i=1}^{n_v} \|(\mathbf{X}^{(i)} - \mathbf{U}^{(i)}\mathbf{V}^T)\mathbf{W}^{(i)}\|_F^2 \\ \text{s.t. } & \mathbf{V} \geq 0 \end{aligned} \quad (16)$$

The partial derivation of $\mathcal{J}(\mathbf{V})$ with respect to \mathbf{V} is

$$\frac{\partial \mathcal{J}}{\partial \mathbf{V}} = \sum_{i=1}^{n_v} \{2\mathbf{W}^{(i)}\mathbf{V}\mathbf{U}^{(i)T}\mathbf{U}^{(i)} - 2\mathbf{W}^{(i)}\mathbf{X}^{(i)T}\mathbf{U}^{(i)}\} \quad (17)$$

Similar to the optimization of semi-NMF, using the KKT complementary condition for the nonnegativity of \mathbf{V} , we get

$$\left(\sum_{i=1}^{n_v} \mathbf{W}^{(i)}\mathbf{V}\mathbf{U}^{(i)T}\mathbf{U}^{(i)} - \sum_{i=1}^{n_v} \mathbf{W}^{(i)}\mathbf{X}^{(i)T}\mathbf{U}^{(i)}\right)_{jk} \mathbf{V}_{jk} = 0 \quad (18)$$

Based on this equation, we can write the updating rule for \mathbf{V} as:

$$\mathbf{V}_{jk} \leftarrow \mathbf{V}_{jk}$$

$$\sqrt{\frac{\sum_{i=1}^{n_v} \{[\mathbf{W}^{(i)}(\mathbf{X}^{(i)T}\mathbf{U}^{(i)})+]_{jk} + [\mathbf{W}^{(i)}\mathbf{V}(\mathbf{U}^{(i)T}\mathbf{U}^{(i)})-]_{jk}\}}{\sum_{i=1}^{n_v} \{[\mathbf{W}^{(i)}(\mathbf{X}^{(i)T}\mathbf{U}^{(i)})-]_{jk} + [\mathbf{W}^{(i)}\mathbf{V}(\mathbf{U}^{(i)T}\mathbf{U}^{(i)})+]_{jk}\}}} \quad (19)$$

It is worth to note that if $\{\mathbf{U}^{(i)}\}_{i=1}^{n_v}$ and \mathbf{V} are a solution of Eq.(11), then $\{\mathbf{U}^{(i)}\mathbf{Q}\}_{i=1}^{n_v}$ and $\mathbf{V}\mathbf{Q}^{-1}$ will form another solution for any invertible matrix \mathbf{Q} . With these requirements, the normalization imposed on $\{\mathbf{U}^{(i)}\}_{i=1}^{n_v}$ and \mathbf{V} are achieved by

$$\begin{aligned} \mathbf{V} &\leftarrow \mathbf{V}\mathbf{Q}^{-1} \\ \mathbf{U}^{(i)} &\leftarrow \mathbf{U}^{(i)}\mathbf{Q} \end{aligned} \quad (20)$$

where \mathbf{Q} is a diagonal matrix formally defined as[Wang *et al.*, 2016a]:

$$\mathbf{Q} = \text{diag}\left(\sum_i \mathbf{V}_{i1}, \sum_i \mathbf{V}_{i2}, \dots, \sum_i \mathbf{V}_{iK}\right). \quad (21)$$

Subproblem of $\{\mathbf{B}^{(i)}\}_{i=1}^{n_v}$. With $\{\mathbf{U}^{(i)}\}_{i=1}^{n_v}$ and \mathbf{V} fixed, for each $\mathbf{B}^{(i)}$, we need to minimize the following objective function:

$$\mathcal{J}(\mathbf{B}^{(i)}) = \|\mathbf{B}^{(i)T}\mathbf{U}^{(i)} - \mathbf{I}\|_F^2 + \beta\|\mathbf{B}^{(i)}\|_{2,1} \quad (22)$$

The partial derivation of $\mathcal{J}(\mathbf{B}^{(i)})$ with respect to $\mathbf{B}^{(i)}$ is

$$\frac{\partial \mathcal{J}}{\partial \mathbf{B}^{(i)}} = 2\mathbf{U}^{(i)}(\mathbf{U}^{(i)T}\mathbf{B}^{(i)} - \mathbf{I}) + \beta\mathbf{D}^{(i)}\mathbf{B}^{(i)} \quad (23)$$

Where $\mathbf{D}^{(i)}$ is a diagonal matrix with the j -th diagonal element given by

$$\mathbf{D}_{jj}^{(i)} = \frac{1}{\|\mathbf{B}_{j\cdot}^{(i)}\|_2} \quad (24)$$

where $\mathbf{B}_{j\cdot}^{(i)}$ is the j -th row of matrix $\mathbf{B}^{(i)}$. Let $\partial\mathcal{J}/\partial\mathbf{B}^{(i)} = 0$, we get the following updating rule for $\mathbf{B}^{(i)}$:

$$\mathbf{B}^{(i)} = [\mathbf{U}^{(i)}\mathbf{U}^{(i)T} + 0.5\beta\mathbf{D}^{(i)}]^{-1}\mathbf{U}^{(i)} \quad (25)$$

Generally, in real-world dataset, $d_i \gg K$, thus $\mathbf{U}^{(i)}\mathbf{U}^{(i)T} + 0.5\beta\mathbf{D}^{(i)}$ is close to singular. In order to avoid inaccurate results and reduce the complexity of the algorithm, we use

Algorithm 1 Optimization of DAIMC

Input: Data matrices for incomplete views $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(n_v)}$, weight matrices $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(n_v)}$, hyper-parameters α, β , number of clusters K .

- 1: Initialize $\mathbf{V} \in \mathbb{R}^{N \times K}$, $\mathbf{U}^{(i)} \in \mathbb{R}^{d_i \times K}$ and $\mathbf{B}^{(i)} \in \mathbb{R}^{d_i \times K}$ $i = 1, \dots, n_v$.
- 2: **repeat**
- 3: **for** $i = 1$ to n_v **do**
- 4: Update $\mathbf{U}^{(i)}$ by Eq.(14) and *lyap* function.
- 5: Update $\mathbf{B}^{(i)}$ by Eq.(26).
- 6: **end for**
- 7: **repeat**
- 8: Update \mathbf{V} by Eq.(19).
- 9: **until** Eq.(16) converges
- 10: Normalize $\mathbf{U}^{(i)}$ s and \mathbf{V} by Eq.(20).
- 11: **until** Eq.(11) converges
- 12: Apply K -means on \mathbf{V} to get the clustering results.
- 13: **Output:** Basis matrices $\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(n_v)}$, regression coefficient matrices $\mathbf{B}^{(1)}, \dots, \mathbf{B}^{(n_v)}$, the common latent feature matrix \mathbf{V} and clustering results.

the matrix inverse equation [Bishop and Nasrabadi, 2007], to reformulate the update rules for $\mathbf{B}^{(i)}$:

$$\begin{aligned} \mathbf{B}^{(i)} &= \frac{2}{\beta} [\mathbf{D}^{(i)-1} - \mathbf{D}^{(i)-1}\mathbf{U}^{(i)}(\mathbf{U}^{(i)T}\mathbf{D}^{(i)-1}\mathbf{U}^{(i)} + 0.5\beta\mathbf{I})^{-1} \\ &\quad \mathbf{U}^{(i)T}\mathbf{D}^{(i)-1}] \mathbf{U}^{(i)} \end{aligned} \quad (26)$$

The entire optimization procedure for DAIMC is summarized in Algorithm 1.

3.5 Convergence and Complexity

Convergence Analysis. As shown by Algorithm 1, the optimization of DAIMC can be divided into three subproblems, each of which is convex w.r.t one variable. Thus, by finding the optimal solution for each subproblem alternatively, our algorithm can at least find a locally optimal solution.

Complexity Analysis. The time complexity of DAIMC is dominated by matrix multiplication and inverse operations. In each iteration, the *lyap* function costs $\mathcal{O}(d_i^3)$ and the matrix inversion in Eq.(27) costs $\mathcal{O}(K^3)$. The complexities of multiplication operations in updating $\mathbf{U}^{(i)}$, \mathbf{V} and $\mathbf{B}^{(i)}$ are $\mathcal{O}(d_i^2K + d_iN + KN)$, $\mathcal{O}(M(d_iKN + d_iK^2 + K^2N))$ and $\mathcal{O}(K^3 + d_i^2K)$ respectively, where M is the iteration times of the inner loop. In general, $K \leq d_i$ and N . Suppose L, d_{max} are the iteration times of the outer loop and the largest dimensionality of all the views respectively, thus the time complexity of DAIMC is $\mathcal{O}(n_v L d_{max}^3 + L M d_{max} K N)$.

4 Experiments and Analysis

Datasets: The experiments are conducted on four real-world multi-view datasets. The important statistics of these datasets are given in the Table 1.

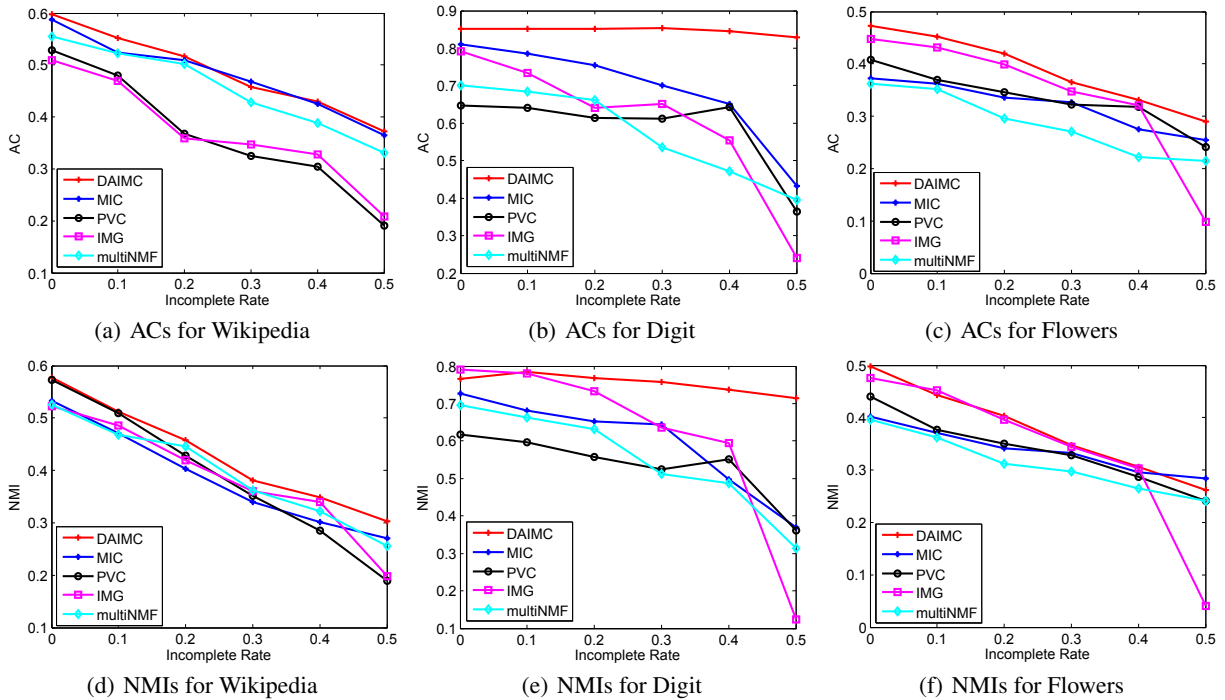


Figure 1: Performance of clustering on Wikipedia, Digit and Flowers.

Dataset	# instances	# views	# clusters
Wikipedia ¹	693	2	10
Digit ²	2000	5	10
3Sources ³	416	3	6
Flowers	1360	3	17

Table 1: Statistics of the datasets

Compared methods: In the experiments, DAIMC is compared with the following state-of-the-art multi-view clustering methods. (1)**MultiNMF**: Multi-view NMF [Liu *et al.*, 2013] seeks a common latent subspace based on joint NMF. However, this method can not deal with the incomplete multi-view data, in our experiment, we therefore first fill the missing instances in each incomplete view with average feature values. (2)**PVC**: Partial multi-view clustering [Li *et al.*, 2014] is one of the state-of-art incomplete multi-view clustering methods, which learns a common latent subspace for the aligned instances and a private latent subspace for the unaligned instances. (3)**IMG**: Incomplete multi-modal visual data grouping [Zhao *et al.*, 2016] integrates PVC and manifold learning, which bridges the connection of missing instance data from different views by learning a complete graph Laplacian term. (4)**MIC**: Multiple incomplete views clustering via weighted NMF [Shao *et al.*, 2015] is a feasible method for incomplete multi-view clustering, which first fills the missing instances in each incomplete view with average feature values, then learns

a common latent subspace with $L_{2,1}$ -norm regularization. All of the hyper-parameters of these methods are selected through grid-search.

For the evaluation metric, we follow [Li *et al.*, 2014], using Normalized Mutual Information (NMI). Besides, precision of clustering result is also reported to give a comprehensive view. Similarly to [Shao *et al.*, 2015], for the complete datasets, we randomly remove some instances from each view to make the views incomplete. The incomplete rate is from 0 (all the views are complete) to 0.5 (all the views have 50% instances missing). It is also worth to note that 3Sources is naturally incomplete. Also since PVC and IMG can only deal with two incomplete views, in order to compare PVC and IMG with other methods, we train these models on all the two-views combinations and report the best result.

4.1 Experimental Results

Table 2 and Figure 1 report the AC and NMI values on image and text datasets with different incomplete rates, respectively. From these table and figures, we can get the following results.

From Figure 1(a) and Figure 1(d), we can see that on Wikipedia dataset, DAIMC raises the performance around 8.65% in NMI with different incomplete rate settings. And in AC, the performance of DAIMC and MIC are close for incomplete rates from 0.2 to 0.5 with the interval of 0.1, the difference between the two methods is just 1%. But when the incomplete rate varies from 0 to 0.1, DAIMC outperforms all the other methods by about 3.53%.

From Figure 1(b) and Figure 1(e), we can see that on Digit dataset, the experimental results of DAIMC are much better than those of other methods. Especially when the incom-

¹<http://www.svcl.ucsd.edu/projects/crossmodal/>

²<http://archive.ics.uci.edu/ml/datasets.html>

³<http://mlg.ucd.ie/datasets/3sources.html>

Methods	BBC-Guardian		BBC-Reuters		Guardian-Reuters		3Sources	
	NMI	AC	NMI	AC	NMI	AC	NMI	AC
DAIMC	0.4264	0.5399	0.4482	0.5641	0.3805	0.5124	0.4733	0.5963
MIC	0.3813	0.4988	0.3814	0.4912	0.3800	0.4612	0.4512	0.5631
PVC	0.2412	0.4334	0.2931	0.4252	0.2488	0.4145	\	\
IMG	0.2614	0.4511	0.3612	0.4624	0.3411	0.4384	\	\
MultiNMF	0.3647	0.4693	0.3687	0.4517	0.3487	0.4281	0.4134	0.4756

Table 2: The NMIs and ACs of different methods on various subsets of 3Sources.

plete rates are large(0.4 and 0.5), DAIMC raises the performance around 60.78% in NMI and 64.67% in AC, respectively. The main reason for this phenomenon is due to that the Digit dataset contains 5 views. DAIMC effectively uses information from different views, reduces the impact of missing samples, and obtains better experimental results.

On Flowers dataset, from Figure 1(c) and Figure 1(f), we can easily see that DAIMC raises the performance around 10.29% in NMI and 20.37% in AC, respectively. Besides, IMG performs very well when the incomplete rate is small (0-0.4), the main reason is due to that Flowers dataset contains a terrible view *D.texturegc*, which plays a negative role in clustering results. Thus, the performances of MIC and multiNMF are bad. In spite of this, DAIMC still performs good with the help of aligning basis matrices.

Table 2 shows the results on the 3Sources dataset, we conduct the experiment on all the two-views combinations and the whole dataset. From Table 2, we can also observe that DAIMC outperforms all the other methods in both NMI and AC.

In summary, when dealing with text data or multi-view data that contains less alignment information, IMG often gets poor result. Meanwhile, although MIC can handle the clustering problem with more than two-views, simply filling the missing instances with the global feature average will lead to a deviation, especially when the incomplete rate is large. By utilizing the information of instance alignment and enforcing the alignment among basis matrices, the proposed DAIMC can get better performances no matter whether it is text dataset or image dataset. Especially when the number of views is large, DAIMC yields more better results.

4.2 Convergence Study

For the convergence study, we conduct an experiment on Digit dataset with the incomplete rate of 0.4 and set the hyper-parameters $\{\alpha, \beta\}$ as $\{1e1, 1e-1\}$ respectively. In Figure 2(a), we show the convergence curve and the NMI values with respect to the number of iterations. The blue curve shows the value of the objective function and the red dashed line indicates the NMI of our method. As can be seen, the algorithm has converged just after 30 iterations.

4.3 View Number Study

In order to demonstrate that the proposed method DAIMC can effectively exploit the information of multiple views, we conduct an experiment on Digit dataset with different view numbers. Similar to convergence study, we set incomplete

rate and hyper-parameters $\{\alpha, \beta\}$ as 0.5 and $\{1e1, 1e0\}$ respectively. The results are shown in Figure 2(b). Obviously, with the increase of the available view number, we get much better result.

4.4 Hyper-parameter Study

The proposed DAIMC method contains two hyper-parameters $\{\alpha, \beta\}$. We conduct the hyper-parameter experiments on Digit dataset. We set the incomplete rate as 0.3 and 0.5 respectively, and report the clustering performance of DAIMC by ranging α and β within the set of $\{1e-4, 1e-3, 1e-2, 1e-1, 1e0, 1e1, 1e2, 1e3\}$. As shown in Figure 2(c) and Figure 2(d), DAIMC obtains a relatively good performance when $\alpha = 1e1$ and $\beta = \{1e-1, 1e0, 1e1\}$.

5 Conclusion

In this paper, we proposed an effective method to deal with incomplete multi-view clustering problem by considering the instance alignment information and enforcing different basis matrices being aligned simultaneously. The experimental results on four real-world multi-view datasets demonstrate the effectiveness of our method. In the future, large scale data will be considered by introducing online learning and incremental learning strategies into our model.

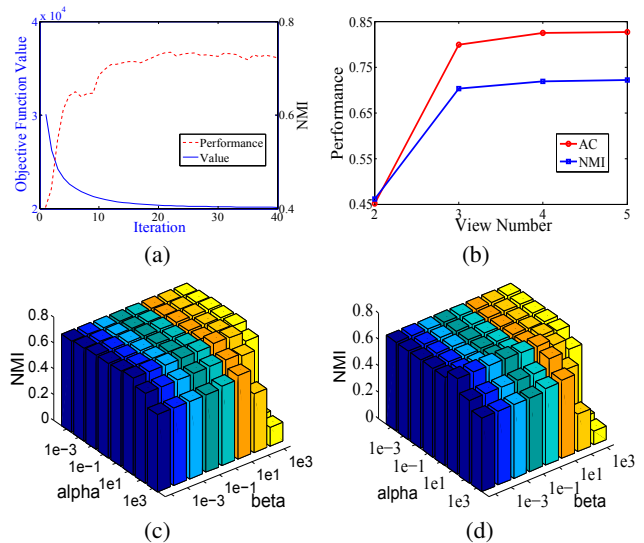


Figure 2: Convergence, view number and hyper-parameter studies on the Digit dataset.

Acknowledgments

This work is supported in part by the NSFC under Grant No. 61672281, and the Key Program of NSFC under Grant No. 61732006

References

- [Bickel and Scheffer, 2004] Steffen Bickel and Tobias Scheffer. Multi-view clustering. In *IEEE International Conference on Data Mining (ICDM)*, pages 19–26. IEEE Computer Society, 2004.
- [Bishop and Nasrabadi, 2007] Christopher M Bishop and Nasser M Nasrabadi. Pattern recognition and machine learning. *Journal of Electronic Imaging*, 16(4):049901, 2007.
- [Blum and Mitchell, 1998] Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Annual Conference on Computational Learning Theory (COLT)*, pages 92–100. ACM, 1998.
- [Chao et al., 2017] Guoqing Chao, Shiliang Sun, and Jinbo Bi. A survey on multi-view clustering. *arXiv preprint arXiv:1712.06246*, 2017.
- [Ding et al., 2010] Chris HQ Ding, Tao Li, and Michael I Jordan. Convex and semi-nonnegative matrix factorizations. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 32(1):45–55, 2010.
- [Fan et al., 2017] Yanbo Fan, Jian Liang, Ran He, Bao-Gang Hu, and Siwei Lyu. Robust localized multi-view subspace clustering. *arXiv preprint arXiv:1705.07777*, 2017.
- [Jing et al., 2017] Xiao-Yuan Jing, Fei Wu, Xiwei Dong, Shiguang Shan, and Songcan Chen. Semi-supervised multi-view correlation feature learning with application to webpage classification. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 1374–1381, 2017.
- [Kong et al., 2011] Deguang Kong, Chris Ding, and Heng Huang. Robust nonnegative matrix factorization using l_{21} -norm. In *ACM International Conference on Information and Knowledge Management (CIKM)*, pages 673–682. ACM, 2011.
- [Lee and Seung, 1999] Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
- [Li et al., 2014] Shao-Yuan Li, Yuan Jiang, and Zhi-Hua Zhou. Partial multi-view clustering. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 1968–1974, 2014.
- [Li, 2016] Yifeng Li. Advances in multi-view matrix factorizations. In *International Joint Conference on Neural Networks (IJCNN)*, pages 3793–3800. IEEE, 2016.
- [Liu et al., 2013] Jialu Liu, Chi Wang, Jing Gao, and Jiawei Han. Multi-view clustering via joint nonnegative matrix factorization. In *SIAM International Conference on Data Mining (SDM)*, pages 252–260. SIAM, 2013.
- [Nie et al., 2017] Feiping Nie, Guohao Cai, and Xuelong Li. Multi-view clustering and semi-supervised classification with adaptive neighbours. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 2408–2414, 2017.
- [Nie et al., 2018] Feiping Nie, Guohao Cai, Jing Li, and Xuelong Li. Auto-weighted multi-view learning for image clustering and semi-supervised classification. *IEEE Transactions on Image Processing*, 27(3):1501–1511, 2018.
- [Potthast et al., 2018] Christian Potthast, Andreas Breitenmoser, Fei Sha, and Gaurav S Sukhatme. Active multi-view object recognition and online feature selection. In *Robotics Research*, pages 471–488. Springer, 2018.
- [Romero et al., 2017] Andrés Romero, Juan León, and Pablo Arbeláez. Multi-view dynamic facial action unit detection. *arXiv preprint arXiv:1704.07863*, 2017.
- [Schechter et al., 2017] Ian Schechter, Tim Wakeling, and Ann M Wollrath. Processing data from multiple sources, March 28 2017. US Patent 9,607,073.
- [Shao et al., 2015] Weixiang Shao, Lifang He, and S Yu Philip. Multiple incomplete views clustering via weighted nonnegative matrix factorization with $l_{2,1}$ regularization. In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD*, pages 318–334, 2015.
- [Sun, 2013] Shiliang Sun. A survey of multi-view machine learning. *Neural Computing and Applications*, pages 2031–2038, 2013.
- [Wang et al., 2016a] Hao Wang, Yan Yang, and Tianrui Li. Multi-view clustering via concept factorization with local manifold regularization. In *IEEE International Conference on Data Mining (ICDM)*, pages 1245–1250, 2016.
- [Wang et al., 2016b] Jing Wang, Xiao Wang, Feng Tian, Chang Hong Liu, Hongchuan Yu, and Yanbei Liu. Adaptive multi-view semi-supervised nonnegative matrix factorization. In *International Conference on Neural Information Processing (ICONIP)*, pages 435–444. Springer, 2016.
- [Wu et al., 2018] Baolei Wu, Enyuan Wang, Zhen Zhu, Wei Chen, and Pengcheng Xiao. Manifold nmf with l_{21} norm for clustering. *Neurocomputing*, 273:78–88, 2018.
- [Xing et al., 2017] Junliang Xing, Zhiheng Niu, Junshi Huang, Weiming Hu, Xi Zhou, and Shuicheng Yan. Towards robust and accurate multi-view and partially-occluded face alignment. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2017.
- [Zhao et al., 2016] Handong Zhao, Hongfu Liu, and Yun Fu. Incomplete multi-modal visual data grouping. In *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 2392–2398, 2016.
- [Zhao et al., 2017] Handong Zhao, Zhengming Ding, and Yun Fu. Multi-view clustering via deep matrix factorization. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 2921–2927, 2017.
- [Zong et al., 2017] Linlin Zong, Xianchao Zhang, Long Zhao, Hong Yu, and Qianli Zhao. Multi-view clustering via multi-manifold regularized non-negative matrix factorization. *Neural Networks*, 88:74–89, 2017.