

Fair Incomplete Multi-View Clustering via Distribution Alignment

Qianqian Wang¹, Haiming Xu^{1*}, Meiling Liu¹, Wei Feng² and Xiangdong Zhang¹

¹School of Telecommunications Engineering, Xidian University, Xi'an, China

²College of Information Engineering, Northwest A&F University, Yangling, China

qqwang@xidian.edu.cn, 24011211044@stu.xidian.edu.cn, lml1710031211@163.com, wei.feng@nwfau.edu.cn, xdchen@mail.xidian.edu.cn

Abstract

Incomplete multi-view clustering (IMVC) extracts consistent and complementary information from multi-view data with missing views, aiming to partition the data into different clusters. It can effectively address the problem of unsupervised multi-view data analysis in complex environments and has gained considerable attention. However, the fairness of IMVC remains underexplored, particularly when data contains sensitive features (e.g., gender, marital status, and age). To tackle the problem, this work presents a novel Fair Incomplete Multi-View Clustering via Distribution Alignment (FIMVC-DA) method. The proposed FIMVC-DA introduces fairness constraints to ensure clustering results are independent of sensitive features. Additionally, it learns consensus representations to enhance clustering performance by maximizing mutual information and aligning the distributions of different views. Experimental results on three datasets containing sensitive features demonstrate that our method improves the fairness of clustering results while outperforming state-of-the-art IMVC methods in clustering performance.

1 Introduction

Multi-view data has become increasingly prevalent, which incorporates multiple perspectives or viewpoints on a particular object [Xu *et al.*, 2013]. Compared with single-view data, multi-view data provides more details that help enhance the effectiveness of clustering tasks [Tao *et al.*, 2019; Ding and Fu, 2018]. Extensive research has been conducted on multi-view clustering (MVC) [Zhan *et al.*, 2019; Chao *et al.*, 2021; Tao *et al.*, 2020], which has been proven as an effective tool for analyzing unlabeled data [Chen *et al.*, 2022]. However, the effectiveness of existing MVC approaches relies on the completeness of the data. In the presence of missing multi-view data, the performance of most MVC methods decreases drastically with an increasing missing rate [Wen *et al.*, 2023]. Consequently, it motivates the investigation of the incomplete multi-view clustering (IMVC)

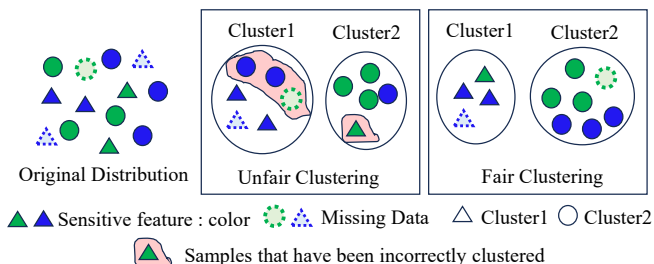


Figure 1: This figure depicts the paper’s motivation. Triangles and circles represent data categories, while blue and green indicate sensitive features like gender. Unfair clustering, as seen in Cluster 1, arises when missing data (e.g., females) leads to overrepresentation of other groups (e.g., males). In Cluster 2, sensitive features themselves may contribute to bias. This paper aims to ensure fair clustering by regulating sensitive feature distribution.

problem[Li *et al.*, 2023].

Incomplete multi-view data [Li *et al.*, 2014a] is ubiquitous due to environmental noise, sensor failure, transmission interference, etc. The simplest IMVC technique involves excluding samples with missing views or imputing the missing view with zeros or the mean, followed by applying an MVC method [Li *et al.*, 2023]. However, the exclusion of samples can lead to the loss of important information [Zhang *et al.*, 2023], and imputing missing values with zeros or means may result in poor clustering results due to the deviation from genuine samples [Xia *et al.*, 2022]. To overcome the difficulty, researchers have extensively investigated the IMVC problem, yielding three main types of IMVC methods: 1) Non-negative matrix factorization (NMF) based methods [Li *et al.*, 2014a]: These methods decompose the data matrix to learn a complementary low-dimensional representation. However, their application is limited to non-negative matrices [Li *et al.*, 2022]. 2) Subspace learning based methods [Yin *et al.*, 2015]: They solve the IMVC problem by learning the completed subspace for clustering; 3) Generation-based methods. Wang *et al.* [2018] utilized the powerful generative capabilities of generative adversarial networks (GANs) to impute missing data, thereby ensuring that the imputed data are closer to the distribution of the original data.

Although the aforementioned methods improve IMVC performance, they treat sensitive features as sample differences,

*corresponding autor

resulting in unfair clustering results [Kenfack *et al.*, 2023]. That is, samples with sensitive features from the same class are always divided into different clusters, as shown in Figure 1. For example, when clustering results are used for targeted advertising or job recommendations, discriminated groups may not receive equal opportunities compared to other groups if the clustering results are not fair. To tackle this problem, several fair MVC have been proposed [Hardt *et al.*, 2016; Lee *et al.*, 2021]. Chierichetti *et al.* [2018] computed the clusters with a fair distribution of sensitive features first and also established evaluation metrics of fairness. The method proposed by Kleindessner *et al.* [2019a] forces groups of different sensitive attributes to have similar similarity distributions in feature space by modifying the spectral clustering objective function, which centers on incorporating fairness constraints into the Laplace matrix construction process, but faces the problem of computational complexity growing with the square of the number of views. The deep fair clustering algorithm proposed by Li *et al.* [2020a] uses a combination of deep representation learning and adversarial training to generate sensitive attribute-independent embeddings via a view-sharing encoder, but the adversarial training model may lead to clustering center drift.

Despite the growing attention towards fairness in MVC, few works consider fairness in IMVC tasks. Nevertheless, the unfairness problem becomes more serious when encountering IMVC problem because IMVC may lose crucial discriminative features, leading to a situation where sensitive features dominate the clustering process. As illustrated in Figure 1, the missing data leads to unfair clustering for cluster 1. Moreover, these sensitive features may be misinterpreted as categorical information, contributing to unfair clustering outcomes for cluster 2.

To address this issue, this paper proposes a novel fair IMVC method that takes into account the fairness of IMVC and reduces the impact of sensitive features on the clustering results. It first maximizes the mutual information to eliminate redundant information from the complete view and extract the minimally sufficient common representation. Then, it aligns the feature distribution between the incomplete view and the complete view, improving the consistency of the common representation. Additionally, the model imposes fairness constraints on the soft clustering allocation, ensuring that the distribution of sensitive features in each cluster closely matches the true distribution. Specifically, our contribution can be summarized as follows:

- We develop a novel fair IMVC method that makes clustering results independent of sensitive features, alleviating the unfairness problem in IMVC. The fairness of clustering is maintained by ensuring that the distribution of sensitive features within each cluster closely aligns with the true distribution.
- We maximize mutual information to mine common features and align complete and incomplete view data distributions to learn the optimal clustering structure.
- Extensive experiments on several datasets demonstrate that the proposed method guarantees the fairness of the

clustering results, and its clustering performance is superior to existing IMVC methods.

2 Related Work

2.1 Incomplete Multi-view Clustering

Incomplete multi-view clustering (IMVC) tackles the challenge of missing views, which impedes effective extraction of cross-view information. Existing IMVC strategies focus on imputing low-dimensional subspace representations or missing views. For instance, matrix factorization methods [Li *et al.*, 2014b; Hu and Chen, 2019] recover non-negative representations from available views, with the work [Yin and Sun, 2022] preserving manifold structures through cosine similarity. Kernel-based techniques similarly impute incomplete matrices using consensus information, including anchor-based strategies [Guo and Ye, 2019] and multiple kernel fusion approaches [Liu *et al.*, 2020]. Generative approaches utilize GANs [Wang *et al.*, 2021; Xu *et al.*, 2019] to synthesize missing views. DIMC-net [Wen *et al.*, 2020] employs graph-guided imputation, while Completer [Lin *et al.*, 2021] leverages dual prediction mechanisms for accurate reconstruction. Despite these advances, two fundamental limitations persist: imputation methods might significantly degrade clustering performance at high missing rates, and the distribution gap between representations of complete and incomplete data negatively influences the imputation validity.

2.2 Fair Incomplete Multi-view Clustering

The integration of fairness and IMVC is an emerging research direction that aims to develop algorithms capable of handling both missing views and clustering bias. Balancing incompleteness and fairness introduces significant technical challenges. DFMVC [2024] uses contrast constraints to align sensitive attribute distributions with the target cluster distribution, but it does not handle the missing-view problem. In contrast, Fair-MVC [2023a] integrates group fairness constraints into IMVC, ensuring that protected groups with sensitive attributes are evenly distributed in each cluster. However, Fair-MVC overlooks the real distribution of sensitive attributes. Besides, a key limitation of both methods is that they do not well model the interaction between missing patterns and sensitive attributes, potentially exacerbating bias. To address these issues, we propose a novel framework based on information bottleneck theory and mutual information maximization that aligns distributions between complete and incomplete views while enforcing fairness constraints, enabling fair clustering on incomplete data without sacrificing clustering performance.

3 Fair Incomplete Multi-View Clustering via Distribution Alignment

Problem setting & Notations: Given a multi-view data matrix $\mathbf{X} = \{\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^v, \dots, \mathbf{X}^V\}$, where $\mathbf{X}^v \in \mathbb{R}^{N \times d_v}$ ($v = 1, 2, \dots, V$), V is the number of views, N is the number of samples, d_v is the dimension of \mathbf{X}^v . Since our model is designed for incomplete multi-view clustering, we divide the data with multiple views \mathbf{X} into two parts: 1)

Notation	Description
\mathbf{X}^v	Input data for the v -th view
$\widehat{\mathbf{X}}^v$	Reconstructed data for the v -th view
\mathbf{X}_p^v	Paired data for the v -th view
\mathbf{X}_u^v	Unpaired data for the v -th view
\mathbf{Z}^v	Subspace feature of the filled v -th view
\mathbf{Z}_p^v	Subspace feature of paired data
\mathbf{Z}_u^v	Subspace feature of unpaired data
\mathbf{Z}_{p+u}^v	Subspace feature of the v -th view
\mathbf{Z}_p	Common feature of paired data
\mathbf{Z}	Shared representation
\mathbf{d}_j^v	Distance from the cluster center.
\mathbf{s}_j	The sensitive features in the j -th cluster.
α, β, γ	Hyperparameters.

Table 1: Descriptions of important notations used in this paper.

paired data of which all the views are complete, denoted as \mathbf{X}_p^v ; 2) unpaired data of which some data is missing \mathbf{X}_u^v . We denote $\widehat{\mathbf{X}}^v$ as the reconstructed data of v -th view.

Definition 1. Fair clustering: It is assumed that there are m cases where the sensitive features $\mathbf{R} \in \mathbb{R}^{N \times m}$ take values, i.e., there are m sensitive groups (based on sensitive features such as gender, race, etc.), and the data is to be clustered into k categories. Let P_s denote the proportion of samples in the entire dataset belonging to sensitive group $s \in \{1, 2, \dots, m\}$. Let $P_{s,j}$ denote the proportion of samples in cluster $j \in \{1, 2, \dots, k\}$ that belong to sensitive group s . We define the fairness ratio $P'_{s,j} = P_{s,j}/P_s$, which measures how the representation of sensitive group s in cluster j compares to its representation in the overall dataset. A clustering result is considered fair if the proportion of each sensitive group within each cluster closely matches the proportion of that sensitive group in the entire dataset, which occurs when $P'_{s,j}$ approaches 1 for all clusters and sensitive groups.

$$\min_{j \in \{1, 2, \dots, k\}, s \in \{1, 2, \dots, m\}} \min\{P'_{s,j}, \frac{1}{P'_{s,j}}\} \rightarrow P'_{s,j} = 1 \quad (1)$$

When $P'_{s,j}$ equals 1, it signifies that the proportion of sensitive groups within each cluster matches the proportion of that sensitive group in the entire dataset, ensuring maximum fairness in clustering. For instance, if the ratio of men to women in a given category is 1:1, clustering should maintain the same ratio within that category.

3.1 Network Architecture

Figure 2 illustrates the overall architecture of FIMVC-DA, which is composed of V encoders/decoders E/D , a common representation learning module, an unpaired data distribution alignment module, and a fairness-aware clustering module. We first send partial multi-view data $[\mathbf{X}_p^v, \mathbf{X}_u^v]$ to the encoder corresponding to each view and obtain latent subspace features $\mathbf{Z}_{p+u}^v = [\mathbf{Z}_p^v, \mathbf{Z}_u^v]$, and the decoders map the latent features to reconstructed samples $\widehat{\mathbf{X}}^v$. Then, the common representation learning module removes the redundant information in the latent subspace feature \mathbf{Z}_p^i and retains the con-

sistent latent representation. Subsequently, we use the unpaired data distribution alignment module to complete the latent representation \mathbf{Z} of data with missing views and guarantee the consistent distribution of complete view data and incomplete view data. Finally, the fairness-aware clustering module achieves the fairness of clustering results.

Encoder/Decoder Module:

The encoder transforms each view’s input data into hidden layer features, which are better suitable for clustering. The subsequent clustering tasks are performed on the feature \mathbf{Z}_{p+u}^v extracted by the encoder. The decoder attempts to reconstruct the original data from the extracted features \mathbf{Z}_{p+u}^v , thereby enforcing the extracted features not to deviate from the original data. Generally, the function of the encoder/decoder module is achieved via a **Reconstruction Loss** as defined below:

$$\mathcal{L}_{rec} = \sum_{v=1}^V \|\mathbf{X}^v - D_v(E_v(\mathbf{X}^v))\|_F^2 = \sum_{v=1}^V \|\mathbf{X}^v - \widehat{\mathbf{X}}^v\|_F^2, \quad (2)$$

where \mathbf{X}^v is the v -th view; E_v represents the encoder of the v -th view; D_v represents the decoder of the v -th view; $\widehat{\mathbf{X}}^v$ represents the reconstructed data of \mathbf{X}^v after undergoing encoding and decoding. A proficiently trained encoder/decoder module should minimize the disparity between \mathbf{X}^v and $\widehat{\mathbf{X}}^v$.

Common Representation Learning Module:

To learn common representations with better clustering structures, we use common features \mathbf{Z}_p as anchors, increasing the common information in paired data from multiple views while removing redundant information from the original data, which is implemented by the following loss:

$$\mathcal{L}_{com} = - \sum_{v=1}^V I(\mathbf{Z}_p^v; \mathbf{Z}_p) + \mathcal{R}(\mathbf{Z}_p^v, \mathbf{Z}_{p+u}^v) \quad (3)$$

where $I(\cdot; \cdot)$ denotes the mutual information between two random variables, which measures the amount of shared information. $\mathcal{R}(\cdot, \cdot)$ is a regularization term whose purpose is to prevent the model from reaching a trivial clustering assignment result, which is formulated as

$$\mathcal{R}(\mathbf{Z}_p^v, \mathbf{Z}_{p+u}^v) = \sum_{v=1}^V I(\mathbf{C}_1^v; \mathbf{C}_2^v), \text{ where } \mathbf{C}_1^v = \sqrt{\mathbf{Z}_p^{vT} \mathbf{Z}_p^v},$$

$$\mathbf{C}_2^v = \sqrt{\mathbf{Z}_{p+u}^{vT} \mathbf{Z}_{p+u}^v}.$$

When the features from multiple views are projected into a shared common feature space, each sample can be represented by features from any view. Building upon these observations, we suggest feature weighting to derive common features for each sample. We denote \mathbf{Z}_p as the learned common feature of the paired data, calculated from the weighted sum of the complete features of all views as

$$\mathbf{Z}_p = \sum_{v=1}^V \mathbf{w}_v \mathbf{Z}_p^v, \quad (4)$$

where \mathbf{w}_v is the weight of the v -th view, which is designed to help \mathbf{Z}_p have more information about the clustering structure.

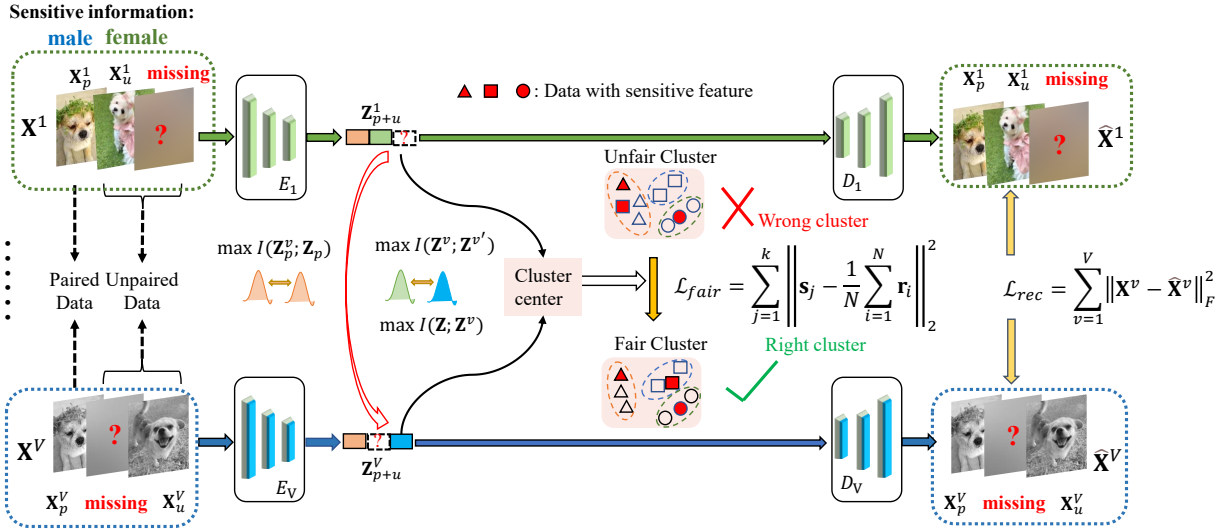


Figure 2: Illustration of the proposed FIMVC-DA model. We maximize $I(\mathbf{Z}_p^v; \mathbf{Z}_p)$ to learn common features in complete views and solve $\max I(\mathbf{Z}^v; \mathbf{Z}^{v'})$ and $\max I(\mathbf{Z}; \mathbf{Z}^v)$ to align the feature distribution between incomplete and complete views. The fair-constrained loss \mathcal{L}_{fair} is used to encourage the fairness of each cluster, and the reconstruction loss \mathcal{L}_{rec} is leveraged to train encoders to ensure the degree of feature restoration to the original data.

It is calculated by the following formula:

$$\mathbf{w}_v = \frac{\sigma(\mathbf{Z}_p^v)}{\sum_{v=1}^V \sigma(\mathbf{Z}_p^v)} \quad (5)$$

where σ denotes the variance, which is used to measure the degree of dispersion of \mathbf{Z}_p^v . The intuition is that a well-defined discrete structure is beneficial for clustering. Therefore, assigning a higher weight to view features with larger variances ensures a better clustering structure for the learned common features \mathbf{Z}_p .

Unpaired Data Distribution Alignment Module:

We aim to learn a comprehensive shared representation from both unpaired and paired data and eliminate differences in \mathbf{Z}_p and \mathbf{Z}_u^v distributions to promote objective Eq. (8). To achieve this goal, we first complete the shared latent representation \mathbf{Z} with unpaired representations provided by other views. This process is conceptually achieved by concatenating the representations from all available data (both paired and unpaired) to form a unified information source. This ensures that the final shared representation \mathbf{Z} for all samples is formed by fusing both the paired and unpaired latent features, i.e., $\mathbf{Z} = \mathbf{Z}_p \oplus \mathbf{Z}_u^v (v = 1, \dots, V)$, where the unpaired features \mathbf{Z}_u^v are implicitly completed by aligning their distributions with the anchor \mathbf{Z}_p . Then, we align the probability distribution of unpaired data with paired data by maximizing the mutual information between two representations of incomplete view data after filling them with each other [Hjelm *et al.*, 2019]. The **Alignment Loss** is as follows:

$$\mathcal{L}_{align} = - \sum_{v=1, v \neq v'}^V (I(\mathbf{Z}^v; \mathbf{Z}^{v'}) + I(\mathbf{Z}; \mathbf{Z}^v) + H(\mathbf{Z}^v)) \quad (6)$$

where v' denotes a view different from v and H denotes entropy. The first mutual information term $I(\mathbf{Z}^v; \mathbf{Z}^{v'})$ aims to

align the distribution of unpaired data, the latter mutual information terms $I(\mathbf{Z}; \mathbf{Z}^v)$ aims to align the distribution between the shared representation \mathbf{Z} (which is informed by all data) and each individual view's representation, and the last term $H(\mathbf{Z}^v)$ is a regularization term to avoid the objective function degenerating into a trivial solution that simply aligns samples into one category, which is calculated as:

$$H(\mathbf{Z}^v) = - \sum_i p(\mathbf{z}_i^v) \log p(\mathbf{z}_i^v) \quad (7)$$

where $p(\mathbf{z}_i^v)$ represents the probability distribution of i -th feature \mathbf{z}_i^v in \mathbf{Z}^v of view v . The entropy term $H(\mathbf{Z}^v)$ serves as a regularization term that maximizes the uncertainty of feature distribution, encouraging the model to explore more diverse feature representations, preventing all samples from being mapped to the same category, and ensuring the diversity and effectiveness of clustering results.

By optimizing this objective function, it aligns effectively unpaired data and avoids trivial solutions. Additionally, it ensures the consistent distribution of complete view data and incomplete view data.

Fairness-aware Clustering Module:

A reasonable strategy for fairness is to ensure that the proportion of sensitive features acquired by clustering is consistent with the original data [Kleindessner *et al.*, 2019b; Cohen *et al.*, 2021]. We first calculate the distance of a sample from each view to each cluster center as:

$$d_{ij}^v = \frac{\exp(\text{sim}(\mathbf{z}_i^v, \mathbf{c}_j))}{\sum_{j'} \exp(\text{sim}(\mathbf{z}_i^v, \mathbf{c}_{j'}))} \quad (8)$$

where $\mathbf{z}_i^v \in \mathbf{Z}^v$ denotes the feature vector of the i -th sample from the v -th view; \mathbf{c}_j is the center of the j -th cluster; $\text{sim}(\cdot, \cdot)$ is the similarity function. d_{ij}^v represents the distance

Algorithm 1 FIMVC-DA

Input: Incomplete multi-view data $\{\mathbf{X}^v\}, v \in \{1, 2 \dots V\}$, Epoch T
Initialize: Extract initial features for all samples from each view.
Parameter: Objective function weights α, β, γ
Output: Clustering results on \mathbf{Z}

- 1: **while** not reach T **do**
- 2: For each view v , extract features for paired samples \mathbf{Z}_p^v and unpaired samples \mathbf{Z}_u^v .
- 3: Calculate the weight of v -th view w_v by Eq. (5).
- 4: Fuse paired features to obtain common representation: $\mathbf{Z}_p = \sum_{v=1}^V w_v \mathbf{Z}_p^v$.
- 5: Concatenate \mathbf{Z}_p and all \mathbf{Z}_u^v to form the unified feature matrix \mathbf{Z} .
- 6: Calculate the distance d_{ij}^v between each sample and each cluster center by Eq. (8).
- 7: Update the distribution of sensitive features \mathbf{s}_j by Eq. (9).
- 8: Adapt the distribution of sensitive features in clustering results by Eq. (10).
- 9: Optimize the total loss by Eq. (11).
- 10: **end while**
- 11: Perform k -means clustering on \mathbf{Z} .

between the i -th sample from the v -th view and the center of the j -th cluster. Suppose \mathbf{r}_i is the sensitive feature of the i -th sample, we calculate the weighted mean of each sensitive feature in j -th cluster as:

$$\mathbf{s}_j = \frac{\sum_{i=1}^N \sum_{v=1}^V d_{ij}^v \mathbf{r}_i}{\sum_{i=1}^N \sum_{v=1}^V d_{ij}^v}, \quad (9)$$

Finally, the **Fairness Constraint Loss** can be defined as:

$$\mathcal{L}_{fair} = \sum_{j=1}^k \left\| \mathbf{s}_j - \frac{1}{N} \sum_{i=1}^N \mathbf{r}_i \right\|_2^2 \quad (10)$$

where k is the number of clusters; $\frac{1}{N} \sum_{i=1}^N \mathbf{r}_i$ is the average value of sensitive features in the entire dataset.

Optimizing this loss ensures that the proportion of sensitive features in each cluster is close to the overall dataset distribution. Thus, this objective function can fix the model's bias problem on sensitive features.

3.2 Objective Function

In summary, the objective loss function of our model is as follows:

$$\mathcal{L} = \mathcal{L}_{rec} + \alpha \mathcal{L}_{com} + \beta \mathcal{L}_{align} + \gamma \mathcal{L}_{fair} \quad (11)$$

where \mathcal{L}_{rec} is the reconstruction loss of encoder and decoder. \mathcal{L}_{com} is a common representation learning loss used to discover common representations across views. \mathcal{L}_{align} is the alignment loss to align the feature distribution between incomplete and complete view data, which makes unpaired data

Dataset	Number	Sensitive feature	Cluster
Credit Card	5000	Gender	5
Zafar	10000	Binary value	2
Bank	5000	Marital status	2

Table 2: Information about the datasets

nearer to its cluster in the feature subspace. \mathcal{L}_{fair} is the fairness-aware clustering loss. α, β , and γ are employed to control \mathcal{L}_{com} , \mathcal{L}_{align} , and \mathcal{L}_{fair} . An excessive β can negatively impact the feature resolution, potentially reducing the clustering to a single class. Conversely, a small β can cause incomplete view data to deviate from the cluster it belongs to. \mathcal{L}_{fair} is the fair constraint loss that ensures that the distribution of sensitive features in each cluster is close to the genuine distribution. We summarize the optimization algorithm in **Algorithm 1**.

4 Experiments

4.1 Experimental Setup

Our experiments were all run on Windows 10 systems with Python 3.7 and Cuda 11.5. The hidden layer dimension for all codec-related algorithms is set to 200 and every random seed has been set to 8. In all experiments, we set the learning rate to 0.0001. We chose Adam as the underlying optimizer. For the Credit Card, and Bank Marketing data set, we use two non-linear functions [McCulloch and Pitts, 1943] (e.g., Sigmoid and Relu) to generate two views.

Datasets We evaluated our model on three fairness datasets: **Credit Card**, **Zafar**, and **Bank** [Zafar *et al.*, 2017]. The Credit Card dataset contains 30,000 samples with gender as the sensitive attribute; Zafar is a widely used synthetic dataset with a binary sensitive attribute; and the Bank dataset includes 40,000 samples with marital status as the sensitive attribute. All datasets were preprocessed to ensure balanced class and sensitive attribute distributions. Table 2 summarizes the main characteristics and sensitive attributes of each dataset.

Comparison Algorithms We compare our method with the following state-of-the-art approaches: Single-view clustering: k -means [Macqueen, 1967], DEC [Xie *et al.*, 2016], CC [Li *et al.*, 2020b], Multi-view clustering: MvDSCN [Wang *et al.*, 2022], Incomplete multi-view clustering: DCP [Lin *et al.*, 2023], APADC [Xu *et al.*, 2023], Fair multi-view clustering: Fair-MVC [Zheng *et al.*, 2023b].

Evaluation Criteria We demonstrated two indicators of the model in clustering tasks: (1) Normalized Mutual Information (NMI); (2) Balance Score $\in [0, 1]$ [Chierichetti *et al.*, 2018]. The tightness of clustering is evaluated using NMI. The fairness of clustering is assessed by the balance score, which is defined as follows:

$$Balance = \min_j \frac{\min_m |C_j \cap s_m|}{|C_j|} \quad (12)$$

where $C_j \in [0, 1]$ represents the j -th cluster, s_m represents m -th protected subgroup. The distribution of sensitive

Datasets	Missing Rate Metrics(%)	0		0.25		0.5		0.75	
		NMI	Balance	NMI	Balance	NMI	Balance	NMI	Balance
Credit Card	<i>k</i> -means[MacQueen, 1967]	20.94 ± 1.14	35.53 ± 0.37	15.67 ± 1.48	36.02 ± 0.60	13.56 ± 0.63	36.32 ± 0.38	9.24 ± 0.72	36.54 ± 0.28
	DEC[Xie <i>et al.</i> , 2016]	21.03 ± 2.09	35.96 ± 0.60	20.05 ± 0.79	36.40 ± 0.78	15.67 ± 1.21	36.26 ± 0.40	10.43 ± 1.53	36.42 ± 0.62
	MvDSCN[Wang <i>et al.</i> , 2022]	21.92 ± 1.53	35.82 ± 0.41	20.34 ± 1.59	35.94 ± 0.48	16.34 ± 1.83	36.63 ± 0.69	11.56 ± 1.62	36.72 ± 0.70
	CC[Li <i>et al.</i> , 2020b]	23.87 ± 1.28	35.74 ± 0.47	20.95 ± 0.68	36.16 ± 0.71	17.62 ± 1.01	36.80 ± 0.97	11.74 ± 0.96	37.18 ± 0.76
	DCP[Lin <i>et al.</i> , 2023]	26.73 ± 0.26	24.19 ± 1.05	23.18 ± 0.34	30.42 ± 2.23	20.04 ± 1.49	34.97 ± 0.47	16.12 ± 1.32	39.18 ± 0.58
	APADC[Xu <i>et al.</i> , 2023]	23.07 ± 0.45	26.32 ± 0.22	23.15 ± 0.24	32.27 ± 0.72	16.30 ± 1.52	32.27 ± 0.71	11.17 ± 1.62	33.29 ± 0.41
	Fair-MVC[Zheng <i>et al.</i> , 2023b]	24.95 ± 0.41	41.98 ± 0.32	22.09 ± 0.75	38.20 ± 0.46	18.71 ± 0.79	39.07 ± 0.76	14.28 ± 0.59	39.82 ± 0.92
FIMVC-DA	26.89 ± 0.53	45.23 ± 1.54	23.26 ± 0.62	43.84 ± 1.53	20.36 ± 0.54	42.31 ± 1.56	17.16 ± 0.56	43.62 ± 1.48	
Zafar	<i>k</i> -means[MacQueen, 1967]	70.32 ± 0.78	17.06 ± 0.76	64.83 ± 1.21	16.35 ± 0.44	60.24 ± 0.90	16.23 ± 1.11	55.65 ± 1.02	16.32 ± 0.86
	DEC[Xie <i>et al.</i> , 2016]	72.55 ± 1.92	16.85 ± 0.73	70.89 ± 1.33	17.11 ± 0.84	64.93 ± 1.28	16.96 ± 0.93	58.96 ± 1.16	17.04 ± 0.86
	MvDSCN[Wang <i>et al.</i> , 2022]	76.91 ± 0.42	17.13 ± 0.65	74.98 ± 1.01	18.32 ± 0.88	69.58 ± 1.66	16.87 ± 0.99	64.26 ± 1.34	17.08 ± 1.17
	CC[Li <i>et al.</i> , 2020b]	78.95 ± 0.68	17.01 ± 0.71	75.22 ± 0.87	18.41 ± 0.67	72.13 ± 0.74	17.96 ± 0.83	68.05 ± 0.72	17.87 ± 1.06
	DCP[Lin <i>et al.</i> , 2023]	81.57 ± 1.57	21.65 ± 1.23	73.79 ± 2.44	22.54 ± 1.85	65.54 ± 1.91	19.89 ± 2.45	57.15 ± 2.18	21.27 ± 2.28
	APADC[Xu <i>et al.</i> , 2023]	72.38 ± 0.80	21.21 ± 0.57	66.37 ± 0.11	21.70 ± 0.31	61.28 ± 0.94	22.02 ± 1.28	55.29 ± 0.82	21.47 ± 0.96
	Fair-MVC[Zheng <i>et al.</i> , 2023b]	81.61 ± 0.57	28.96 ± 0.59	79.89 ± 0.55	30.64 ± 0.89	76.97 ± 0.79	29.33 ± 0.80	73.16 ± 0.92	29.58 ± 0.86
FIMVC-DA	83.23 ± 0.81	30.32 ± 1.67	80.94 ± 0.66	32.32 ± 0.91	78.22 ± 0.73	30.34 ± 0.85	75.84 ± 0.85	30.28 ± 0.92	
Bank	<i>k</i> -means[MacQueen, 1967]	28.67 ± 1.44	37.65 ± 0.66	24.77 ± 1.71	36.66 ± 0.45	21.08 ± 1.17	36.32 ± 0.89	16.83 ± 1.46	36.25 ± 0.92
	DEC[Xie <i>et al.</i> , 2016]	30.93 ± 1.15	37.60 ± 0.96	28.97 ± 1.80	36.42 ± 0.69	24.37 ± 1.02	37.09 ± 0.78	19.43 ± 1.22	36.68 ± 0.85
	MvDSCN[Wang <i>et al.</i> , 2022]	36.24 ± 0.55	37.59 ± 0.67	35.02 ± 1.21	36.11 ± 0.57	31.33 ± 1.76	36.96 ± 0.77	27.15 ± 1.54	37.24 ± 0.86
	CC[Li <i>et al.</i> , 2020b]	36.23 ± 1.01	37.46 ± 0.97	34.88 ± 1.33	37.69 ± 0.95	31.13 ± 1.63	37.80 ± 0.85	27.12 ± 1.54	37.62 ± 0.94
	DCP[Lin <i>et al.</i> , 2023]	39.93 ± 1.84	26.75 ± 2.06	34.20 ± 2.87	20.49 ± 1.54	28.12 ± 2.31	27.29 ± 1.24	21.54 ± 2.46	26.58 ± 1.36
	APADC[Xu <i>et al.</i> , 2023]	40.62 ± 0.25	27.79 ± 2.59	32.21 ± 2.79	28.41 ± 2.11	33.14 ± 2.77	27.82 ± 2.28	26.32 ± 2.56	27.92 ± 2.12
	Fair-MVC[Zheng <i>et al.</i> , 2023b]	38.99 ± 0.91	42.40 ± 0.75	36.66 ± 1.01	41.76 ± 0.83	32.90 ± 1.03	41.61 ± 0.67	28.34 ± 0.98	41.82 ± 0.78
FIMVC-DA	41.12 ± 0.80	44.35 ± 0.67	38.89 ± 1.05	43.30 ± 0.93	34.43 ± 0.91	43.29 ± 0.76	30.25 ± 0.87	43.58 ± 0.89	

Table 3: Clustering NMI and balance scores on the three datasets contain sensitive features. The best results are marked in bold. Suboptimal results are represented in blue.

L_1	L_2	L_3	L_4	Credit Card		Bank		Zafar	
				NMI	Balance	NMI	Balance	NMI	Balance
✓				14.76	14.47	15.53	13.21	22.03	17.71
✓	✓			17.31	23.22	20.00	16.52	39.87	17.12
✓	✓	✓		20.12	26.17	32.15	20.32	72.24	21.47
✓	✓	✓	✓	20.36	42.31	34.43	43.29	78.22	30.34

Table 4: Ablation study results across three datasets. The best and suboptimal results are marked in bold and blue respectively.

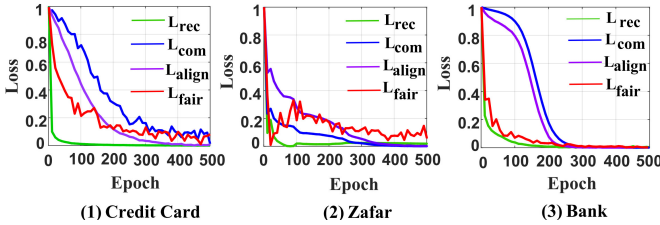


Figure 3: Convergence curves on the three datasets.

traits typically determines the upper limit of balance, and the higher the balance value, the fairer the outcome.

4.2 Experimental Results

Definition 2 Missing Rate If X comprises N samples across V views, the total number of samples is $N * V$. If some views of certain data are missing, let the total number of missing samples be N_{miss} . The missing rate is then calculated as $N_{miss} / (N * V)$.

Clustering Results and Balance Scores on Three Datasets

Table 3 shows that our method outperforms existing clustering methods on the three datasets, including the most recent clustering algorithm. Despite a certain constraint between the balance score and clustering accuracy, the results show that our algorithm achieves higher clustering accuracy. This is

consistent with our alignment loss, which aligns the distribution of partial and complete view data, allowing the learned features to be more beneficial for clustering tasks. Our model exceeds the second-best result in clustering NMI by 2% on average in the three datasets, and the fairness score of clustering clusters is 3% higher on average. This result shows that our model achieves better feature extraction and fusion capabilities, which helps to improve the degree of balance while maintaining clustering accuracy.

Sensitive Features Visualization Figure 4 presents t-SNE visualizations of the Credit Card dataset, illustrating the effect of fairness constraints on clustering results. In the raw data (a), male (blue) and female (red) samples are evenly distributed, reflecting the balanced gender ratio in the original dataset. Without fairness constraints (b), some clusters are dominated by a single gender, indicating the presence of bias in the clustering results. In contrast, with fairness constraints applied (c), the gender distribution within each cluster closely matches that of the original data. These results demonstrate that our method effectively mitigates bias and achieves fair clustering outcomes that are independent of sensitive attributes.

Ablation Study In this part, we conduct an ablation study to analyze the influence of the reconstruction loss \mathcal{L}_{rec} (L_1), the common representation learning loss \mathcal{L}_{com} (L_2), the alignment loss \mathcal{L}_{align} (L_3), and the fairness-aware clustering loss \mathcal{L}_{fair} (L_4). We report experimental results on the three datasets in Table 4. An obvious finding is that the fairness constraint loss significantly improves the fairness of the clustering results. The incomplete view alignment scheme based on information theory effectively improves the clustering performance. As can be seen from the results in the last two rows, the fairness constraint also helps to promote correct clustering to some extent. In summary, the ablation study

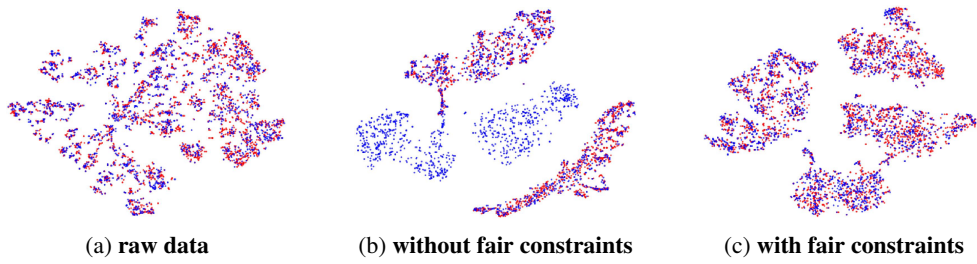


Figure 4: Visualize sensitive features in Credit Card dataset using t-SNE. Blue represents the male feature, while red represents the female.

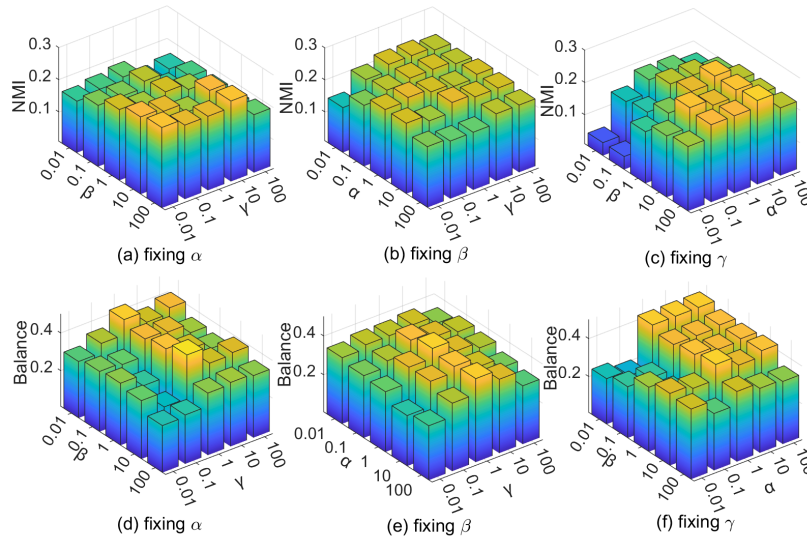


Figure 5: NMI and Balance scores with different hyperparameters at a missing rate of 0.5 on the credit card dataset.

demonstrates that each component of the proposed FIMVC-DA effectively improves clustering performance and fairness.

Convergence Analysis Figure 3 depicts the convergence curves of the model’s objective functions on three datasets. We normalize the loss function values of the four parts because they differ in magnitude. After training 500 batches, each portion of the loss function on the four data sets is reduced to a constant. This suggests that our model exhibits an elevated level of convergence.

Parametric Analysis From Figure 5, it can be seen that when fixing α to be 1, the other two hyperparameters do not have much effect on NMI. However, it is seen that the fairness score is higher when the hyperparameter γ is larger. When fixing β to 10, the other two hyperparameters have little effect on the NMI results as well as on the fairness score. When fixing γ to 1, the other two parameters controlling the consistency representation of the learning view do not differ much in the NMI values when taking values between 1 and 10. The clustering fairness is also the best.

5 Conclusion

This study addresses fair clustering in the context of missing multi-view data, presenting a novel approach named Fair Incomplete Multi-View Clustering via Distribution Alignment

(FIMVC-DA). Initially, we tackle the challenge of integrating incomplete view data and refining common features by maximizing mutual information from complete views. Additionally, we align the distribution of incomplete views to enhance coordination between the data. Fairness constraints are then employed to ensure proximity between sensitive feature distributions in clusters and those in real categories, thereby mitigating model bias risks. Empirical evaluations on three datasets confirm the superior clustering performance of our method, with experiments on sensitive feature visualization highlighting its effectiveness in ensuring clustering fairness.

Acknowledgments

This work is supported by the National Natural Science Foundation of China under Grant 62176203, the Fundamental Research Funds for the Central Universities (ZYTS25267, QTZX25004), and the Science and Technology Project of Xi’an (Grant 2022JH-JSYF-0009), Open Project of Anhui Provincial Key Laboratory of Multimodal Cognitive Computation, Anhui University (No. MMC202416), Selected Support Project for Scientific and Technological Activities of Returned Overseas Chinese Scholars in Shaanxi Province 2023-02, and the Xidian Innovation Fund (Project NoYJSJ25007).

References

- [Chao *et al.*, 2021] Guoqing Chao, Shiliang Sun, and Jinbo Bi. A survey on multiview clustering. *IEEE Transactions on Artificial Intelligence*, 2(2):146–168, 2021.
- [Chen *et al.*, 2022] Man-Sheng Chen, Tuo Liu, Chang-Dong Wang, Dong Huang, and Jian-Huang Lai. Adaptively-weighted integral space for fast multiview clustering. In *Proceedings of the 30th ACM International Conference on Multimedia*, MM ’22, page 3774–3782, New York, NY, USA, 2022. Association for Computing Machinery.
- [Chierichetti *et al.*, 2018] Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, and Sergei Vassilvitskii. Fair clustering through fairlets, 2018.
- [Cohen *et al.*, 2021] Maxime C. Cohen, Adam N. Elmachtoub, and Xiao Lei. Price discrimination with fairness constraints. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’21, page 2, New York, NY, USA, 2021. Association for Computing Machinery.
- [Ding and Fu, 2018] Zhengming Ding and Yun Fu. Robust multiview data analysis through collective low-rank subspace. *IEEE Transactions on Neural Networks and Learning Systems*, 29(5):1986–1997, 2018.
- [Guo and Ye, 2019] Jun Guo and Jiahui Ye. Anchors bring ease: An embarrassingly simple approach to partial multi-view clustering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 118–125, 2019.
- [Hardt *et al.*, 2016] Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning. In *NeurIPS*, page 3323–3331, 2016.
- [Hjelm *et al.*, 2019] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization, 2019.
- [Hu and Chen, 2019] Menglei Hu and Songcan Chen. One-pass incomplete multi-view clustering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 3838–3845, 2019.
- [Kenfack *et al.*, 2023] Patrik Joslin Kenfack, Adín Ramírez Rivera, Adil Mehmood Khan, and Manuel Mazzara. Learning fair representations through uniformly distributed sensitive attributes. In *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pages 58–67, 2023.
- [Kleindessner *et al.*, 2019a] Matthäus Kleindessner, Samira Samadi, Pranjal Awasthi, and Jamie Morgenstern. Guarantees for spectral clustering with fairness constraints. In *International conference on machine learning*, pages 3458–3467. PMLR, 2019.
- [Kleindessner *et al.*, 2019b] Matthäus Kleindessner, Samira Samadi, Pranjal Awasthi, and Jamie Morgenstern. Guarantees for spectral clustering with fairness constraints. In *ICML*, pages 3458–3467, 09–15 Jun 2019.
- [Lee *et al.*, 2021] Woojin Lee, Hyungjin Ko, Junyoung Byun, Taeho Yoon, and Jaewook Lee. Fair clustering with fair correspondence distribution. *Information Sciences*, 581:155–178, 2021.
- [Li *et al.*, 2014a] Shao-Yuan Li, Yuan Jiang, and Zhi-Hua Zhou. Partial multi-view clustering. In *AAAI*, AAAI’14, page 1968–1974. AAAI Press, 2014.
- [Li *et al.*, 2014b] Shao-Yuan Li, Yuan Jiang, and Zhi-Hua Zhou. Partial multi-view clustering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 28, 2014.
- [Li *et al.*, 2020a] Peizhao Li, Han Zhao, and Hongfu Liu. Deep fair clustering for visual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9070–9079, 2020.
- [Li *et al.*, 2020b] Yunfan Li, Peng Hu, Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng. Contrastive clustering, 2020.
- [Li *et al.*, 2022] Xingfeng Li, Quansen Sun, Zhenwen Ren, and Yinghui Sun. Dynamic incomplete multi-view imputing and clustering. In *Proceedings of the 30th ACM International Conference on Multimedia*, MM ’22, page 3412–3420, New York, NY, USA, 2022. Association for Computing Machinery.
- [Li *et al.*, 2023] Xingfeng Li, Yinghui Sun, Quansen Sun, Jia Dai, and Zhenwen Ren. Distribution consistency based fast anchor imputation for incomplete multi-view clustering. In *Proceedings of the 31st ACM International Conference on Multimedia*, MM ’23, page 368–376, New York, NY, USA, 2023. Association for Computing Machinery.
- [Lin *et al.*, 2021] Yijie Lin, Yuanbiao Gou, Zitao Liu, Boyun Li, Jiancheng Lv, and Xi Peng. Completer: Incomplete multi-view clustering via contrastive prediction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11174–11183, 2021.
- [Lin *et al.*, 2023] Yijie Lin, Yuanbiao Gou, Xiaotian Liu, Jinfeng Bai, Jiancheng Lv, and Xi Peng. Dual contrastive prediction for incomplete multi-view representation learning. *IEEE TPAMI*, 45(4):4447–4461, 2023.
- [Liu *et al.*, 2020] Xinwang Liu, Miaomiao Li, Chang Tang, Jingyuan Xia, Jian Xiong, Li Liu, Marius Kloft, and En Zhu. Efficient and effective regularized incomplete multi-view clustering. *IEEE transactions on pattern analysis and machine intelligence*, 43(8):2634–2646, 2020.
- [Macqueen, 1967] J. Macqueen. Some methods for classification and analysis of multivariate observations. *Proc. Symp. Math. Statist. and Probability*, 5th, 1, 1967.
- [McCulloch and Pitts, 1943] Warren S. McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *biol math biophys*, 1943.
- [Tao *et al.*, 2019] Zhiqiang Tao, Hongfu Liu, Huazhu Fu, and Yun Fu. Multi-view saliency-guided clustering for image cosegmentation. *IEEE Transactions on Image Processing*, 28(9):4634–4645, 2019.

- [Tao *et al.*, 2020] Zhiqiang Tao, Hongfu Liu, Sheng Li, Zhengming Ding, and Yun Fu. Marginalized multiview ensemble clustering. *IEEE Transactions on Neural Networks and Learning Systems*, 31(2):600–611, 2020.
- [Wang *et al.*, 2018] Qianqian Wang, Zhengming Ding, Zhiqiang Tao, Quanxue Gao, and Yun Fu. Partial multi-view clustering via consistent gan. In *IEEE ICDM*, pages 1290–1295, 2018.
- [Wang *et al.*, 2021] Qianqian Wang, Zhengming Ding, Zhiqiang Tao, Quanxue Gao, and Yun Fu. Generative partial multi-view clustering with adaptive fusion and cycle consistency. *IEEE Transactions on Image Processing*, 30:1771–1783, 2021.
- [Wang *et al.*, 2022] Qianqian Wang, Zhiqiang Tao, Quanxue Gao, and Licheng Jiao. Multi-view subspace clustering via structured multi-pathway network. *IEEE TNLS*, pages 1–7, 2022.
- [Wen *et al.*, 2020] Jie Wen, Zheng Zhang, Zhao Zhang, Zhihao Wu, Lunke Fei, Yong Xu, and Bob Zhang. Dimcnet: Deep incomplete multi-view clustering network. In *Proceedings of the 28th ACM international conference on multimedia*, pages 3753–3761, 2020.
- [Wen *et al.*, 2023] Jie Wen, Gehui Xu, Chengliang Liu, Lunke Fei, Chao Huang, Wei Wang, and Yong Xu. Localized and balanced efficient incomplete multi-view clustering. In *Proceedings of the 31st ACM International Conference on Multimedia*, MM '23, page 2927–2935, New York, NY, USA, 2023. Association for Computing Machinery.
- [Xia *et al.*, 2022] Wei Xia, Quanxue Gao, Qianqian Wang, and Xinbo Gao. Tensor completion-based incomplete multiview clustering. *IEEE Transactions on Cybernetics*, 52(12):13635–13644, 2022.
- [Xie *et al.*, 2016] Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis, 2016.
- [Xu *et al.*, 2013] Chang Xu, Dacheng Tao, and Chao Xu. A survey on multi-view learning, 2013.
- [Xu *et al.*, 2019] Cai Xu, Ziyu Guan, Wei Zhao, Hongchang Wu, Yunfei Niu, and Beilei Ling. Adversarial incomplete multi-view clustering. In *IJCAI*, volume 7, pages 3933–3939, 2019.
- [Xu *et al.*, 2023] Jie Xu, Chao Li, Liang Peng, Yazhou Ren, Xiaoshuang Shi, Heng Tao Shen, and Xiaofeng Zhu. Adaptive feature projection with distribution alignment for deep incomplete multi-view clustering. *IEEE TIP*, 32:1354–1366, 2023.
- [Yin and Sun, 2022] Jun Yin and Shiliang Sun. Incomplete multi-view clustering with cosine similarity. *Pattern Recognition*, 123:108371, 2022.
- [Yin *et al.*, 2015] Qiyue Yin, Shu Wu, and Liang Wang. Incomplete multi-view clustering via subspace learning. *CIKM*, 2015.
- [Zafar *et al.*, 2017] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P. Gummadi. Fairness constraints: Mechanisms for fair classification, 2017.
- [Zhan *et al.*, 2019] Kun Zhan, Chaoxi Niu, Changlu Chen, Feiping Nie, Changqing Zhang, and Yi Yang. Graph structure fusion for multiview clustering. *IEEE Transactions on Knowledge and Data Engineering*, 31(10):1984–1993, 2019.
- [Zhang *et al.*, 2023] Chao Zhang, Jingwen Wei, Bo Wang, Zechao Li, Chunlin Chen, and Huaxiong Li. Robust spectral embedding completion based incomplete multi-view clustering. In *Proceedings of the 31st ACM International Conference on Multimedia*, MM '23, page 300–308, New York, NY, USA, 2023. Association for Computing Machinery.
- [Zhao *et al.*, 2024] Bowen Zhao, Qianqian Wang, Zhiqiang Tao, Wei Feng, and Quanxue Gao. Dfmvc: Deep fair multi-view clustering. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 8090–8099, 2024.
- [Zheng *et al.*, 2023a] Lecheng Zheng, Yada Zhu, and Jingrui He. Fairness-aware multi-view clustering. In *Proceedings of the 2023 SIAM International Conference on Data Mining (SDM)*, pages 856–864. SIAM, 2023.
- [Zheng *et al.*, 2023b] Lecheng Zheng, Yada Zhu, and Jingrui He. Fairness-aware multi-view clustering, 2023.