

## 面向子空间聚类的多视图统一表示学习网络

林毓秀<sup>1,2</sup> 刘慧<sup>1,2</sup> 于晓<sup>1,2</sup> 张彩明<sup>2,3</sup>

<sup>1</sup> (山东财经大学计算机科学与技术学院 济南 250014)

<sup>2</sup> (山东省数字媒体技术重点实验室 (山东财经大学) 济南 250014)

<sup>3</sup> (山东大学软件学院 济南 250010)

(linyx0109@mail.sdufe.edu.cn)

## A Multi-view Unified Representation Learning Network for Subspace Clustering

Lin Yuxiu<sup>1,2</sup>, Liu Hui<sup>1,2</sup>, Yu Xiao<sup>1,2</sup>, and Zhang Caiming<sup>2,3</sup>

<sup>1</sup> (School of Computer Science and Technology, Shandong University of Finance and Economics, Jinan 250014)

<sup>2</sup> (Digital Media Technology Key Lab of Shandong Province (Shandong University of Finance and Economics), Jinan 250014)

<sup>3</sup> (Software college, Shandong University, Jinan 250010)

**Abstract** Multi-view subspace clustering aims to explore the rich information across views to guide the clustering process. The key to it lies in how to effectively learn the unified representation and subspace representation between views. Recently, deep clustering methods have achieved promising effects through the powerful representation capability of neural networks. However, the inherent multi-source heterogeneity of multi-view data makes existing methods commonly achieve independent encoding of each view with a unimodal encoder, which not only increases the number of model parameters but also has poor model generalization performance. Besides, the low-rank subspace representation has been shown to promote clustering performance, and the traditional nuclear norm regularization does not consider the difference of different singular values, leading to a biased estimation. To this end, this paper proposes a multi-view subspace clustering method with Transformer (MURLN). MURLN first introduces the Transformer into our multi-view representation learning network, which projects different views into the low-dimensional feature space with the same mapping rule by sharing parameters, aligning the latent representation of each view. Secondly, a weighted fusion strategy for intra-view samples is conducted to construct a unified representation in a rational way. In addition, we use the weighted Schatten- $p$  norm as low-rank regularization constraints to generate lower-rank subspace representations, which further improves the clustering performance. Extensive experiments on multi-view datasets verify the effectiveness and superiority of our proposed method.

**Key words** multi-view subspace clustering; Transformer; weighted fusion; low-rank representation; weighted Schatten  $p$ -norm

**摘要** 多视图子空间聚类旨在挖掘多视图丰富的信息来指导高维数据聚类,其研究关键在于如何有效地学习多个视图的统一表示和子空间表示。近年来,深度聚类方法利用神经网络强大的表征能力取得了优异的性能。然而,多视图数据固有的多源异构性使得大多数现有方法以单模态编码器实现对各个视图的独立编码,不仅增加

**收稿日期:** 2024-00-00; **修回日期:** 2020-00-00; (学生第一作者)

**基金项目:** 国家自然科学基金项目 (62072274, U22A2033); 中央引导地方科技发展项目(YDZX2022009); 山东省泰山学者特聘专家基金项目 (tstp20221137)

This work is supported by the National Natural Science Foundation of China (62072274, U22A2033), the Central Guidance on Local Science and Technology Development Project (YDZX2022009), and the Special Funds of Taishan Scholars Project of Shandong Province (tstp20221137).

**通信作者:** 刘慧 (liuh\_lh@sdufe.edu.cn)

了模型参数量而且降低了模型泛化性. 另一方面, 低秩子空间表示被证明能够促进聚类性能的提升, 传统的核范数正则化优化没有考虑不同奇异值隐含的信息量差异, 是矩阵秩的一个有偏估计. 为此, 本文提出了一种面向子空间聚类的多视图统一表示学习网络 (MURLN). 首先, MURLN 使用 Transformer 作为编码器, 通过共享参数将异构视图以相同的映射规则投影到低维特征空间, 从而对齐每个视图的潜在表示. 其次, 采用视图内样本加权融合的方法, 合理地构建统一表示. 此外, 利用加权 Schatten- $p$  范数构建低秩正则化约束, 生成更低秩的子空间表示, 进一步提高了聚类性能. 在多视图数据集上的广泛实验验证了所提方法的有效性和优越性.

**关键词** 多视图子空间聚类; Transformer; 加权融合; 低秩表示; 加权 Schatten- $p$  范数

**中图法分类号** TP391

信息技术的普及和迅速发展, 极大地丰富了数据在现实应用场景中的表现形式, 使其呈现出高维、多源和异构属性<sup>[1]</sup>. 例如, 在机器人领域, 人机交互依赖视觉、听觉、触觉等方面的感官信息; 在内容理解领域, 多媒体片段包含图片、文本以及音频信号; 在医疗诊断领域, 脑磁共振成像根据扫描参数的不同分为 T1WI、T2WI、T2-FLAIR 和增强扫描; 在人脸识别领域, 提取图像的 LBP, HOG、Gabor 等特征以降低阴影及光照变化的影响. 广义地讲, 从多角度、多渠道对同一观测对象进行的多样性描述被定义为多视图数据, 每一种模态或特征都是描述对象的一个视图<sup>[2]</sup>. 对比传统的单个视图, 多视图包含两个重要的信息分量, 即一致性信息和互补性信息<sup>[3]</sup>, 能够实现对数据样本的全面理解.

聚类作为一种有效的数据预处理技术, 旨在将无标签数据按内在相似性划分到不同类别, 已被广泛应用于智能医学<sup>[4]</sup>、推荐系统<sup>[5]</sup>、模式识别<sup>[6]</sup>等任务场景. 因此, 协同利用多视图数据蕴含的丰富信息来提升聚类效果具有重要的研究意义. 当前的多视图聚类方法可以粗略地划分为 4 类: 基于协同训练的方法<sup>[7]</sup>、基于多核学习的方法<sup>[8]</sup>、多视图图聚类<sup>[9]</sup>和多视图子空间聚类 (Multi-view Subspace Clustering, MvSC)<sup>[10-11]</sup>. MvSC 通过将高维多视图数据映射到低维特征子空间, 并从中获得统一的特征表示用于子空间自表达学习, 是实现高维数据聚类的一种有效方法. 得益于深度神经网络在特征提取和表示方面的显著优势, 深度多视图子空间聚类能够探索复杂场景中样本更深层的特征表达, 已成为近年来的研究热点<sup>[12-14]</sup>.

自编码器 (Autoencoders, AEs) 是一种典型的深度无监督学习模型, 能够将无标签样本的语义信息表示为特征向量<sup>[15]</sup>. 国内外学者将其与多视图子空间聚类结合, 提出了一系列方法. Abavisani 等人<sup>[16]</sup>为每个视图分别训练单独的卷积自编码器, 然后强制学习到的潜在表示具有相同的子空间结构. 在此基础上, Zhu 等人<sup>[17]</sup>使用两个自编码器对每个视图进行表征, 用于构建一组特定于视图的子空间表示和一个共享的子空间表示. Zhang 等人<sup>[18]</sup>提出一种深度嵌套的自

编码器模型 (Autoencoder in Autoencoder Networks, AE<sup>2</sup>-Nets) 来退化出视图间共享的紧致表示. 为了提高表示学习的质量, Li 等人<sup>[19]</sup>利用多尺度特征学习策略, 为每个视图提取不同尺度的低维嵌入. 总的来说, 由于多视图数据存在于不一致的数据空间, 这些方法通常采用“一对一”的编码网络对各个视图进行独立编码. 特别是在处理跨模态的多视图数据时, 需要根据具体数据类型选择不同的自编码器架构<sup>[20]</sup>, 如图 1 所示.

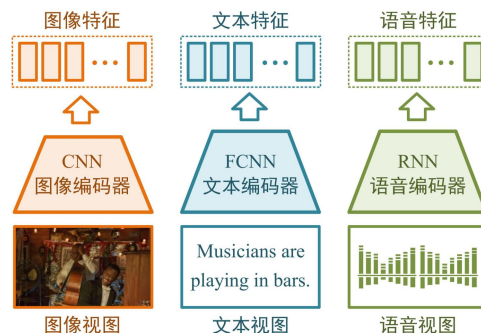


Fig. 1 Schematic diagram of multi-view encoding

图 1 多视图编码架构示意图

近些年, Transformer<sup>[21]</sup>在自然语言处理和计算机视觉中展现出强大的建模能力以及良好的并行计算能力, 逐渐成为处理文本和图像的通用编码架构<sup>[22-23]</sup>. 为此, 本文提出了一种新的面向子空间聚类的多视图统一表示学习网络 (Multi-view Unified Representation Learning Network, MURLN). 首先, 我们结合自注意力机制以及多视图学习的特点, 提出基于 Transformer 的多视图共享编码框架, 将多视图数据输入到同一个编码器中, 通过共享参数的形式进行训练, 从而提高了模型的泛化能力. 编码过程同时施加多样性约束以挖掘视图间的互补信息. 考虑到同一个观测样本在不同视图中的不同聚类表现, 本文将多视图融合的目标粒度从整个视图细化到视图内样本, 提出一种样本加权的融合方法从低维潜在空间中构建统一表示. 通过迭代优化模型参数, 自适应地学习视图中每个样本的合适权重. 其次, 为追求子空间表示的低秩性, 设计基于加权 Schatten- $p$  范数<sup>[24]</sup>的正则化约束. 网络结构如图 2 所示.

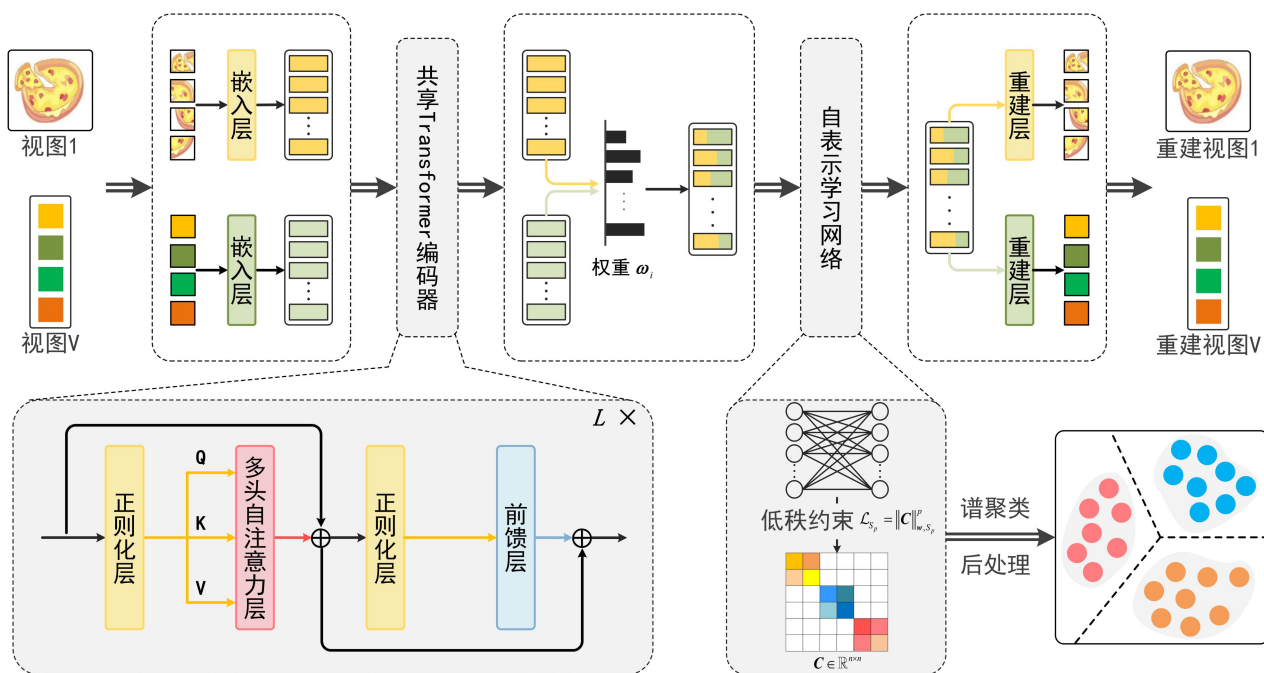


Fig. 2 Framework of the MURLN

图2 MURLN 的架构图

综上, 本文主要贡献包括 4 个方面:

1) 提出基于 Transformer 的多视图共享编码框架, 通过共享网络参数对齐不同视图的潜在表示, 消除多视图数据的异构性差异.

2) 设计了一种样本加权的多视图特征融合方法, 降低视图内低质量样本对融合结果的影响, 加强可靠样本的权重.

3) 提出基于加权 Schatten- $p$  范数的低秩正则化约束, 与现有方法中常用的核范数相比, 由于考虑了奇异值之间的显著差异, 生成的子空间表示能够更好地逼近原始秩函数.

4) 在多个类型的多视图数据集上进行了大量实验. 结果表明, 与其他方法相比, MURLN 实现了更优越的聚类性能.

本文其余部分组织如下: 第 1 节对一些相关工作进行回顾. 第 2 节详细介绍本文提出的方法. 第 3 节通过在多个数据集上的大量实验证明方法的有效性. 第 4 节对本工作进行总结和展望.

## 1 相关工作

### 1.1 符号约定和问题定义

为了方便描述, 首先对符号约定进行说明. 大写粗体字母表示矩阵 (例如,  $\mathbf{A}$ ), 小写粗体字母表示向量 (例如,  $\mathbf{a}$ ), 标量用非粗体字母表示 (例如,  $a$ ). 表 1 列出了文中常用的符号及定义.

给定具有  $V$  个视图的数据集  $\{\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(V)}\}$ , 其中,  $\mathbf{X}^{(v)} = [\mathbf{x}_1^{(v)}, \mathbf{x}_2^{(v)}, \dots, \mathbf{x}_N^{(v)}] \in \mathbb{R}^{d_v \times N}$  表示第  $v$  个视图维度为  $d_v$  的数据矩阵. 本文假设数据完整, 即所有视

图的样本数量均为  $N$ . 多视图子空间聚类旨在整合多视图的一致和互补性信息学习视图间共享的自表示系数矩阵  $\mathbf{C} \in \mathbb{R}^{N \times N}$ , 并在此基础上构建相似度矩阵用于谱聚类.

Table 1 Some Important Notations

表 1 主要符号

符号	定义
$\mathbf{x}_i^{(v)} \in \mathbb{R}^{d_v}$	视图 $v$ 中的第 $i$ 个样本
$\mathbf{I}$	单位矩阵
$ A $	绝对值操作
$\sigma_i(\mathbf{A})$	矩阵 $\mathbf{A}$ 的第 $i$ 个奇异值
$\text{diag}(\mathbf{A})$	以向量的形式返回矩阵 $\mathbf{A}$ 的对角线值
$\text{tr}(\cdot)$	迹函数

### 1.2 子空间聚类

基于谱聚类的子空间聚类依赖于数据的自表达特性<sup>[25-26]</sup>, 它假设空间中每个样本都可以由同一子空间中其他样本的线性组合表示, 非同一子空间中的样本对应的表示系数几乎为零. 经典的单视图子空间聚类可以形式化为如下优化问题:

$$\min_{\mathbf{C}} \|\mathbf{X} - \mathbf{XC}\|_F + \lambda \|\mathbf{C}\|_c, \quad \text{s.t. } \text{diag}(\mathbf{C}) = \mathbf{0} \quad (1)$$

其中,  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ ,  $D$  表示视图维度.  $\|\cdot\|_F$  表示重建损失,  $\|\cdot\|_c$  表示正则化约束项. 可选约束项  $\text{diag}(\mathbf{C}) = \mathbf{0}$  用来避免得到平凡解  $\mathbf{C} = \mathbf{I}$ . 通过优化目标函数式 (1), 得到表示系数矩阵  $\mathbf{C}$ , 并计算亲和矩阵  $\mathbf{W} = (\|\mathbf{C}\| + \|\mathbf{C}^T\|)/2$ . 最后, 执行谱聚类算法, 如 Ncut<sup>[27]</sup>, 得到聚类结果.

然而, 利用单视图的有限信息进行聚类存在较多

局限,一旦数据受损可能产生较差的聚类效果.另外,由于多视图数据的异构性,直接使用单视图聚类方法处理多视图数据也是不合理的.

### 1.3 多视图子空间聚类

多视图子空间通过探索并整合多视图的一致和互补性信息,显著提高了聚类准确性<sup>[28-29]</sup>.例如,DiMSC<sup>[30]</sup>设计基于希尔伯特-施密特独立准则的正则化约束来隐式地增强多视图表示的多样性.直接从原始数据空间进行自表达重构可能会受到噪声影响,LMSC<sup>[10]</sup>提出将输入数据映射到一个隐空间,然后在隐空间中学习统一的低秩子空间表示.与浅层模型不同,深度 MvSC 方法利用神经网络构建有效的数据表征,取得了不错的效果<sup>[31-32]</sup>.在 LMSC 的基础上,Zhang 等人<sup>[33]</sup>基于神经网络提出一个更泛化的模型(gLMSC).DMSC-UDL<sup>[12]</sup>通过卷积自编码器整合局部和全局结构信息,并利用哈达玛积判别约束强化视图内样本的簇间差异.Wang 等人<sup>[34]</sup>设计一个深度结构化多路径网络(SMpNet)用于多视图子空间聚类,利用编码网络中间层产生的隐藏特征信息,提高了自表示系数矩阵的质量.

现有深度方法中使用的自编码器可分为卷积自编码器和全连接自编码器.卷积自编码器<sup>[35]</sup>侧重于通过卷积层的级联提取图像数据的潜在表示,全连接自编码器则更适用于处理图像以外的数据,例如文本、1 维特征等.由于视图间的异构性,现有深度多视图子空间方法对不同视图采用特定的编码器单独编码,有较高的训练成本.尽管可以先利用预训练模型(例如,Faster-RCNN<sup>[36]</sup>和 BERT<sup>[37]</sup>)提取不同模态的特征<sup>[38]</sup>,但这种做法会带来较高的模型参数量,也很难实现多个潜在表示的对齐.Transformer 模型的出现很大程度上为该问题提供了一种有前途的解决方法.根据文献[23],依赖于自注意力机制的 Transformer 允许网络根据输入内容动态收集相关特征,被证明能够作为通用的基础编码模型用于多模态数据特征的学习.

为了区分不同视图对提升聚类性能的贡献程度,iCmSC<sup>[39]</sup>使用一个视图级融合算子实现多个潜在表示的加权融合.AMVDSN<sup>[40]</sup>引入注意力机制获得每个视图的动态权重,使融合过程更加关注重要的视图.实际上,质量较差的视图也可能存在有价值的观测样本,质量较好的视图中同样可能存在不重要的或者部分受损的样本.因此,将每个视图视为一个整体来分配权重的全局融合方法,忽视了同一样本在不同视图对聚类结果的贡献不同,甚至会对统一表示的学习造成不利影响.

### 1.4 低秩表示学习

基于谱聚类的子空间聚类算法性能高度依赖于自表示系数矩阵的质量.大多数现有方法通过构建适当的正则化约束来保证系数矩阵满足某些特定属性,常见的包括稀疏子空间聚类(Sparse Subspace Clustering, SSC)<sup>[25]</sup>和低秩表示(Low-Rank Representation, LRR)<sup>[26]</sup>.SSC 使用  $\ell_0$  或  $\ell_1$  范数正则化来约束子空间表示的稀疏性.LRR 将系数矩阵的秩函数凸松弛为核范数,通过求解核范数最小化加强子空间表示的低秩性.在多视图子空间聚类问题中,低秩表示由于其良好的聚类性能和噪声鲁棒性而被广泛研究<sup>[41-43]</sup>.根据标准核范数的定义, $\|C\|_* = \sum_i \sigma_i(C)$ ,其在最小化过程中将所有非零奇异值平均化,是矩阵秩的一个有偏估计.为此,构建低秩正则化约束需要充分考虑奇异值之间的差异性,这有利于挖掘数据中有价值的信息.

## 2 本文方法

本文设计了一个面向子空间聚类的多视图统一表示学习网络,网络结构如图 1 所示,主要包含五个模块.

### 2.1 模型结构

#### 2.1.1 特征嵌入模块

MURLN 首先提出一个特征嵌入模块,将不同视图的输入转换为相同维度的特征嵌入.更具体地说,对于包含  $V$  个视图的数据集,特征嵌入模块由  $V$  个特定于视图的嵌入层组成.以第  $v$  个视图输入  $\mathbf{X}^{(v)} \in \mathbb{R}^{d_v \times N}$  为例,嵌入层通常包含三步操作:1) 将输入数据划分为  $M$  个不重叠的块(patch);2) 将块线性映射(通常由可学习的线性层实现)到  $k$  维的嵌入空间,得到视图嵌入  $\mathbf{X}_{\text{emb}}^{(v)} \in \mathbb{R}^{M \times k \times N}$ ;3) 构造位置编码  $\mathbf{X}_{\text{pos}}^{(v)} \in \mathbb{R}^{M \times k \times N}$  来保留数据内的相对关系.嵌入层的输出  $\mathbf{X}_E^{(v)} \in \mathbb{R}^{M \times k \times N}$  最终表示为带有位置编码的视图嵌入:

$$\mathbf{X}_E^{(v)} = \mathbf{X}_{\text{emb}}^{(v)} + \mathbf{X}_{\text{pos}}^{(v)}. \quad (2)$$

由于视图间的异构性,嵌入层中的块划分和位置编码策略取决于输入视图的模态类型.以长和宽分别为  $H$  和  $W$  的  $C$  通道 2-D 图像视图  $\mathbf{X}^{(v)} \in \mathbb{R}^{H \times W \times C \times N}$  为例,参考文献[44],将其划分为  $M = H \times W / P^2$  个边长为  $P$  的图像块并拉伸成 1-D 特征向量,并按照图像块的位次索引编码位置.由于特征嵌入模块仅用于统一输入维度,且不同模态的块划分和位置编码策略已经得到了广泛的研究,因此其具体实现方式不作为本文研究的重点.

#### 2.1.2 共享编码器模块

在统一视图嵌入维度的基础上,进一步设计了一

个基于 Transformer 结构的共享编码器模块, 来挖掘不同视图的高级语义信息. 该编码器通过参数共享的形式将各个视图嵌入以相同的映射规则投影到低维潜在空间, 有利于实现视图间的对齐, 同时避免维度灾难和可能存在的噪声干扰. 具体来说, 编码器由  $L$  个相同的网络层组堆叠而成, 每个层组包括一个多头自注意力层和一个带有 GELU 激活函数的全连接前馈层. 在每层之前执行归一化操作, 每层之间由残差结构连接. 给定第  $v$  个视图的嵌入  $\mathbf{X}_E^{(v)}$ , 编码器中第  $l$  层组的计算过程为:

$$\mathbf{X}_0^{(v)} = \mathbf{X}_E^{(v)}, \quad (3)$$

$$\overline{\mathbf{X}}_l^{(v)} = \text{MSA}(\text{LN}(\mathbf{X}_{l-1}^{(v)})) + \mathbf{X}_{l-1}^{(v)}, \quad (4)$$

$$\mathbf{X}_l^{(v)} = \text{FFN}(\text{LN}(\overline{\mathbf{X}}_l^{(v)})) + \overline{\mathbf{X}}_l^{(v)}, \quad (5)$$

其中,  $\mathbf{x}_l^{(v)}$ ,  $l=1,2,\dots,L$  是编码模块中第  $l$  层的输出,  $\text{MSA}(\cdot)$ 、 $\text{LN}(\cdot)$  和  $\text{FFN}(\cdot)$  分别表示多头自注意力层、归一化层和全连接前馈层. 经过嵌入模块和共享编码器模块, 原始高维空间中的数据  $\mathbf{X}^{(v)}$  被逐步映射到低维特征空间, 并最终得到潜在表示  $\mathbf{F}^{(v)} = \mathbf{X}_L^{(v)} \in \mathbb{R}^{d \times N}$ , 其中  $d = M \times k$  为潜在表示的维度.

### 2.1.3 样本加权融合模块

利用多视图数据的互补性信息来构建一个统一表示, 有利于全面认识数据本质. 为此, MURLN 提出一种样本加权的多视图特征融合方法, 整合多个视图的潜在表示为一个统一特征表示  $\mathbf{F} \in \mathbb{R}^{d \times N}$ . 进一步地, 考虑到观测样本在不同视图中的质量差异, 简单地拼接特征会忽略视图间以及视图内的复杂关系和不平衡性. 另一方面, 在没有先验知识的情况下, 人工设计权重又具有盲目性. 本文将样本的权重关系融入目标函数中, 通过网络迭代优化动态地进行更新. 加权融合的过程形式化为:

$$\mathbf{F} = \sum_{v=1}^V \sum_{i=1}^N \omega_i^{(v)} \mathbf{f}_i^{(v)}, \quad (6)$$

其中,  $\boldsymbol{\omega}^{(v)} = [\omega_1^{(v)}, \omega_2^{(v)}, \dots, \omega_N^{(v)}] \in \mathbb{R}^N$  表示第  $v$  个视图对应的权重向量. 多视图特征表示的有效融合能够强化特征的表达能力, 有助于提升聚类性能.

### 2.1.4 自表示学习模块

遵循一致性假设, 所有视图应具有相似甚至相同的聚类结构. 自表示学习模块包含一个不带偏置和非线性激活函数的全连接层, 它能够自然地模拟数据的自表达特性, 从而描述成对样本之间的关系, 称为自表达层. 将统一表示  $\mathbf{F}$  和特定视图的潜在表示  $\mathbf{F}^{(v)}$  共同作为自表达层的输入, 学习所有视图共享的一致子

空间表示  $\mathbf{C} \in \mathbb{R}^{N \times N}$ ,  $\mathbf{c}_{i,j}$  反映了样本  $i$  和样本  $j$  的相似性关系.

现有子空间聚类方法表明<sup>[33]</sup>, 低秩表示可以保持数据的结构特征, 并且减少噪声等冗余信息的干扰, 有利于学习更理想的系数矩阵. 为此, 本文在自表示学习模块中定义一个基于加权 Schatten- $p$  范数的低秩约束强化子空间表示的低秩性. 与传统的基于核范数的低秩约束相比, 加权 Schatten- $p$  范数不仅在形式上更加灵活, 同时还考虑了不同秩分量之间的显著差异 (较大奇异值通常表示数据矩阵中嵌入的重要结构信息), 已在图像去噪、矩阵补全等任务中验证了有效性<sup>[45-46]</sup>. 定义为:

$$\begin{aligned} \|\mathbf{C}\|_{w, S_p} &= (\sum_{i=1}^N w_i \sigma_i(\mathbf{C})^p)^{1/p}, \\ \text{s.t. } w_i &= 1 / (\sigma_i(\mathbf{C}) + \varepsilon) \end{aligned} \quad (7)$$

其中,  $0 < p \leq 1$ .  $\mathbf{w} = (w_1, w_2, \dots, w_i, \dots, w_N)$  是非负权重向量, 通过不同的权重因子实现对非零奇异值不同程度的收缩.  $\varepsilon$  为一个极小值以避免分母为零, 通常取  $10^{-16}$ . 为避免开根过程, 公式 (7) 改写如下:

$$\|\mathbf{C}\|_{w, S_p}^p = \sum_{i=1}^N w_i \sigma_i(\mathbf{C})^p = \text{tr}(\mathbf{W} \boldsymbol{\Delta}^p), \quad (8)$$

其中,  $\mathbf{W}$  和  $\boldsymbol{\Delta}$  为对角矩阵, 主对角线上的元素分别由  $w_i$  和  $\sigma_i$  组成. 当  $p=1$  时, 公式 (8) 退化为加权核范数最小化问题, 当  $p \rightarrow 0$  时, Schatten-0 范数等价于求解系数矩阵的秩最小化问题. 因此, Schatten- $p$  范数在形式上比核范数更接近于目标矩阵的秩函数, 对低秩问题的求解也更为灵活.

### 2.1.5 视图重建模块

对于无监督的多视图聚类, 验证编码器能否学习到紧致的特征表示主要通过解码器将特征重建回原始数据  $\widehat{\mathbf{X}}^{(v)} \in \mathbb{R}^{d_v \times N}$ , 并最小化重建误差来实现. MURLN 使用单层全连接层作为解码器的具体实现.

## 2.2 模型优化

网络模型的训练过程可以形式化为对应目标函数的优化问题. 首先, 通过最小化输入样本与重建样本之间的重建损失, 端到端地训练编码器-解码器. 重建损失  $\mathcal{L}_{re}$  定义为:

$$\mathcal{L}_{re} = \sum_v \left\| \mathbf{X}^{(v)} - \widehat{\mathbf{X}}^{(v)} \right\|_F^2. \quad (9)$$

将得到的多视图潜在表示输入样本加权融合模块, 构造统一的特征表示, 提高表达能力. 另外, 为了有效挖掘视图间的互补性信息, MURLN 利用哈达玛积设计正则化约束, 增强多视图表示之间的多样性. 融合损失定义为:



$$\mathcal{L}_f = \sum_v \underbrace{\|F - F^{(v)}\|_F^2}_{\text{fusion error}} + \sum_{v \neq u} \underbrace{\|F^{(v)} \odot F^{(u)}\|_1}_{\text{diversity constraint}}, \quad (10)$$

其中,  $\odot$  表示哈达玛积, 隐式地鼓励统一表示充分编码互补信息.

自表示系数矩阵的学习受到统一表示和特定视图潜在表示两方面的共同约束, 定义为:

$$\mathcal{L}_{self} = \frac{1}{V} \sum_v \|(F^{(v)} - F^{(v)}C)\|_F^2 + \|F - FC\|_F^2. \quad (11)$$

根据对各部分损失函数的分析, 模型的整体目标函数形式如下:

$$\min \mathcal{L} = \min \mathcal{L}_{re} + \lambda_1 \mathcal{L}_f + \lambda_2 \mathcal{L}_{self} + \lambda_3 \mathcal{L}_{S_p}, \quad (12)$$

其中,  $\lambda_1 > 0, \lambda_2 > 0, \lambda_3 > 0$  为权衡参数.  $\mathcal{L}_{S_p} = \|C\|_{w, S_p}^p$  是作用在子空间表示上的正则化约束项.

### 2.3 训练细节

为了加速模型收敛, 本文使用“预训练+微调”的两阶段优化策略训练网络.

第一阶段预训练编码器和解码器. 首先从均值为零的高斯分布中随机采样初始化编解码器参数, 并随机采样  $b$  个样本  $x_i^{(v)}, i=1, 2, \dots, b$  作为一个训练批次. 然后利用随机梯度下降算法最小化公式 (9) 中的重建损失, 迭代更新编解码器的参数.

基于上述初始值, 第二阶段对完整网络进行微调, 其中未预训练的参数被随机初始化. 由于自表达层学习所有样本的关联关系, 因此将全部数据作为微调阶段的输入. 利用梯度下降算法对公式 (12) 迭代求解, 直至收敛. 一旦网络收敛, 学习到的自表示系数矩阵被用来计算亲和矩阵, 最后通过谱聚类后处理获得聚类结果. 具体的参数设置和实验结果在第 3 章中详细讨论. 算法 1 提供了本文模型的训练流程.

#### 算法 1. MURLN 算法.

输入: 多视图数据  $X^{(v)}$ , 样本数量  $N$ , 批大小  $b$ , 特征维度  $d$ , 参数  $p$ , 权衡参数  $\lambda_1, \lambda_2, \lambda_3$ ;

输出: 聚类结果.

- ① 初始化模型参数、超参数及自表示矩阵  $C$ ;
- ② /\* 步骤 1: 预训练自编码器 \*/
- ③ WHILE 未收敛 do
- ④ 采样一批数据:  $x_i^{(v)} \in X^{(v)}, i=1, 2, \dots, b$ ;
- ⑤ 最小化公式 (9), 更新特征嵌入模块、共享编码器模块和视图重建模块参数;
- ⑥ END WHILE
- ⑦ /\* 步骤 2: 微调完整网络 \*/
- ⑧ WHILE 未收敛 do
- ⑨ 在数据集  $X^{(v)}$  上最小化公式 (12), 更新所

有可训练参数及自表示矩阵  $C$ ;

⑩ END WHILE

⑪ 计算亲和矩阵, 执行谱聚类;

⑫ 返回聚类结果

## 3 实验结果与分析

### 3.1 数据集和实验设置

**实验环境.** 提出的方法在 Windows11 操作系统利用 PyTorch 深度学习库由 Python 编程语言实现. 实验环境的主要硬件参数为 NVIDIA GeForce RTX 3080 Ti GPU, 12G 显存, Intel(R) Core(TM) i9-10900F CPU, 32G 内存.

**数据集.** 实验在 7 个公开多视图数据集上进行, 涉及新闻文档 (BBCSport)、手写数字 (MNIST-USPS)、人脸目标 (Yale、Extended YaleB) 和通用对象 (RGB-D、Caltech101-20、MSRCV1) 等多种场景. 表 2 列出了数据集的详细信息.

**评价指标.** 本文使用 6 种通用的评价指标<sup>[41]</sup>综合评估方法的聚类性能, 包括准确性 (Accuracy, ACC)、标准化互信息 (Normalized Mutual Information, NMI)、调整兰德指数 (Adjusted Rand Index, AR)、F 分数 (F-Measure, FM)、精度 (Precision, P) 和召回率 (Recall, R). 所有指标的值越高, 聚类质量越好.

**对比方法.** 方法 MURLN 将与 10 种先进的多视图聚类方法进行比较, 包括基于多样性约束的多视图子空间聚类 (DiMSC<sup>[30]</sup>), 基于低秩表示的多视图聚类方法 (LMSC<sup>[10]</sup>、RMSL<sup>[28]</sup>、DSS-MSC<sup>[29]</sup>、FCMSC<sup>[43]</sup> 和 FLMSC<sup>[11]</sup>), 深度多视图子空间聚类 (MSCNLG<sup>[31]</sup>、DMSC-UDL<sup>[12]</sup>、SDMSC<sup>[19]</sup>、D<sup>2</sup>MVSC<sup>[32]</sup>). 所有对比方法遵循原文建议仔细地调整最优参数. 在处理多特征数据集时, DMSC-UDL 方法中使用的卷积自编码器被替换为工作[32]中的全连接自编码器. 为了减少随机性的影响, 所有方法独立运行 20 次, 报告各指标在每个数据集上的平均值.

**实验参数.** 采用 AdamW 优化器对网络进行优化, 初始学习率设置为  $5e-4$ , 动量因子和权值衰减分别设为 0.9 和 0.05. 优化器根据任务损失产生的梯度来更新模型参数. 在预训练阶段, 训练迭代次数设置为 500, 批大小为 256 (数据量不足时, 批大小为样本数量). 在微调阶段, 训练迭代次数设置为 200, 批大小设置为输入数据的样本量, 以学习完整的表示. 用四元组 [分块数  $M$ , 潜在表示的特征维度  $d$ , 网络层数  $L$ , 注意力头个数] 表示编码器的超参数. 对于

Yale 数据集, 编码器结构设置为[4, 256, 4, 4]. 对于 MSRCV1 数据集, 编码器结构设置为[4, 512, 4, 8]. 在 MNIST-USPS, BBCSport, Caltech101-20 和 Extended

YaleB 数据集中, 编码器结构为[8, 512, 8, 16]. 通过参数敏感性分析, 本文将模型在不同数据集上的超参数统一设置为  $\lambda_1 = 0.01$ ,  $\lambda_2 = 0.5$ ,  $\lambda_3 = 1$ .

Table 2 Dataset Details

表 2 数据集详细信息

场景类型	数据集	类别数	样本数	视图数	视图维度
新闻文档	BBCSport[33]	5	544	2	[3183, 3203]
手写字符	MNIST-USPS[20]	10	5000	2	[28×28, 28×28]
人脸目标	Yale[30]	15	165	3	[4096, 3304, 6750]
	Extended YaleB[19]	10	640	3	[2500, 3304, 6750]
通用对象	MSRCV1[10]	7	210	6	[1302, 48, 512, 100, 256, 210]
	RGB-D[17]	50	500	2	[64×64, 64×64]
	Caltech101-20[31]	20	2386	6	[48, 40, 254, 1984, 512, 928]

Table 3 Clustering results in Yale and Extended YaleB datasets (mean %)

表 3 数据集 Yale、Extended YaleB 上的聚类结果 (均值%)

方法	Yale						Extended YaleB					
	ACC	NMI	AR	FM	P	R	ACC	NMI	AR	FM	P	R
DiMSC	70.90	72.70	53.50	56.40	54.30	58.60	61.50	63.50	45.30	50.40	48.10	53.40
LMSC	75.20	73.50	55.10	56.40	54.30	57.10	52.42	52.14	26.82	35.10	30.60	41.16
RMSL	74.70	75.50	55.64	58.55	54.51	63.30	53.92	52.75	28.93	36.81	32.72	42.08
DSS-MVC	78.20	77.90	60.10	61.30	59.20	62.20	78.65	76.86	64.81	68.47	65.48	71.85
FCMSC	76.36	78.90	63.67	65.97	63.52	68.61	53.20	53.24	27.51	35.72	31.09	41.98
FLMSC	72.72	74.02	55.67	58.47	56.58	60.49	57.85	55.86	39.49	47.90	46.35	49.64
MSCNLG	91.67	90.30	82.10	83.21	81.76	84.70	51.23	48.88	25.45	33.75	29.91	38.72
DMSC-UDL	77.09	75.86	56.01	58.92	54.36	64.32	92.55	88.69	82.12	83.94	81.51	86.54
SDMSC	84.73	84.85	66.35	69.15	86.62	82.73	<u>98.17</u>	<u>95.90</u>	<u>94.03</u>	<u>95.82</u>	94.10	94.88
D <sup>2</sup> MVSC	<u>94.50</u>	<b>93.90</b>	<u>88.90</u>	<u>90.40</u>	<u>90.50</u>	<b>90.40</b>	<b>99.32</b>	<b>98.71</b>	<b>98.24</b>	<b>99.01</b>	<b>98.94</b>	<b>98.93</b>
MURLN	<b>94.95</b>	<u>93.88</u>	<b>89.37</b>	<b>90.59</b>	<b>90.68</b>	<u>89.72</u>	97.40	93.22	92.36	93.09	<u>94.31</u>	<u>95.27</u>

Table 4 Clustering results in Caltech101-20 and MSRCV1 datasets (mean %)

表 4 数据集 Caltech101-20、MSRCV1 中的聚类结果(均值%)

方法	Caltech101-20						MSRCV1					
	ACC	NMI	AR	FM	P	R	ACC	NMI	AR	FM	P	R
DiMSC	39.12	51.96	26.29	32.23	65.03	21.43	80.36	68.61	62.07	67.41	66.35	68.51
LMSC	55.78	65.12	42.95	48.88	76.61	35.89	80.60	65.30	59.90	65.20	61.20	66.30
RMSL	40.49	47.17	30.82	37.46	63.35	26.61	71.74	63.95	56.10	62.39	60.41	64.52
DSS-MVC	44.11	63.20	31.71	37.43	72.67	25.21	86.57	78.64	72.94	76.74	75.83	77.68
FCMSC	45.98	63.15	34.44	40.44	71.95	28.13	80.90	70.10	63.58	68.70	67.76	69.73
FLMSC	48.16	64.15	35.39	40.98	<u>77.09</u>	28.91	97.57	<b>96.12</b>	95.70	<b>96.16</b>	96.11	96.21
MSCNLG	50.23	58.24	36.65	46.70	47.98	45.51	93.08	80.69	83.16	87.09	89.00	85.26
DMSC-UDL	56.76	64.22	48.20	55.86	62.55	51.12	78.36	66.97	59.85	65.46	65.46	65.92
SDMSC	62.00	59.60	48.56	54.04	61.90	62.00	97.06	90.37	92.06	93.99	<b>97.76</b>	96.59
D <sup>2</sup> MVSC	<u>72.73</u>	<u>73.85</u>	<u>56.35</u>	<u>59.15</u>	76.62	<u>72.73</u>	<u>97.78</u>	93.00	<u>95.74</u>	95.36	96.45	<b>97.03</b>
MURLN	<b>80.19</b>	<b>78.37</b>	<b>85.21</b>	<b>87.09</b>	<b>86.54</b>	<b>89.30</b>	<b>98.34</b>	<u>95.99</u>	<b>96.41</b>	<u>96.07</u>	<u>96.82</u>	<u>96.87</u>

Table 5 Clustering results in BBCSport and MNIST-USPS datasets (mean %)

表 5 数据集 BBCSport、MNIST-USPS 中的聚类结果(均值%)

方法	BBCSport						MNIST-USPS					
	ACC	NMI	AR	FM	P	R	ACC	NMI	AR	FM	P	R
DiMSC	96.00	88.50	92.00	92.90	91.40	93.00	63.72	59.94	50.08	55.14	54.32	55.99
LMSC	91.20	82.60	84.20	88.70	87.30	87.70	73.33	75.48	65.35	68.99	65.57	72.81
RMSL	97.61	91.81	92.90	95.34	95.23	94.47	44.64	33.19	32.88	24.71	30.25	36.02
DSS-MVC	96.60	88.40	89.80	92.30	92.70	91.80	82.93	87.54	80.99	83.00	78.40	88.17
FCMSC	96.51	89.04	91.06	93.17	94.21	92.15	76.42	77.88	72.54	68.97	69.32	76.51
FLMSC	92.72	85.56	84.57	89.65	89.05	88.53	93.27	86.29	86.93	88.23	88.15	88.32
MSCNLG	88.79	72.50	73.39	79.91	77.73	82.23	67.60	77.95	62.66	67.04	57.12	81.12
DMSC-UDL	95.96	87.94	89.00	91.64	91.30	91.97	78.36	66.97	59.85	65.46	65.46	65.92
SDMSC	94.86	86.82	89.36	89.93	86.00	87.92	81.42	88.83	86.99	83.06	78.10	80.04
D <sup>2</sup> MVSC	<u>98.10</u>	<u>93.80</u>	<u>94.90</u>	<u>96.40</u>	<u>96.40</u>	<u>96.41</u>	<u>93.86</u>	<u>90.72</u>	<u>90.60</u>	<u>90.36</u>	<u>89.01</u>	<u>89.32</u>
MURLN	<b>98.70</b>	<b>95.07</b>	<b>94.96</b>	<b>96.41</b>	<b>96.93</b>	<b>96.82</b>	<b>95.63</b>	<b>93.11</b>	<b>92.87</b>	<b>93.49</b>	<b>92.55</b>	<b>93.94</b>

Table 6 Comparison of clustering performance of three coding networks (mean %)

表 6 三种编码网络的聚类性能比较(均值%)

数据集	视图(维数)	ACC			NMI		
		MURLN_CNN	MURLN_FC	MURLN	MURLN_CNN	MURLN_FC	MURLN
MNIST-USPS	视图-1 (28*28)	76.71	72.25	86.74	72.89	74.40	82.64
	视图-2 (28*28)	61.20	75.33	71.16	59.10	69.93	74.15
	MNIST-USPS	88.94	85.40	<b>95.63</b>	86.14	84.82	<b>93.11</b>
MNIST-USPS-4V	视图-1 (28*28)	-	41.31	85.67	-	19.56	78.40
	视图-2 (28*28)	-	50.90	72.78	-	42.23	67.44
	视图-3 (944)	-	52.56	62.81	-	44.08	59.66
	视图-4 (944)	-	35.04	67.34	-	19.75	62.19
	MNIST-USPS-4V	-	79.28	<b>93.27</b>	-	77.38	<b>91.78</b>

### 3.2 性能比较与分析

#### 3.2.1 聚类性能比较

实验结果如表 3-5 所示. 其中, 最优和次优结果分别用粗体和下划线标记. 从实验结果中可发现, 所提出的方法在所有数据集上都取得了有竞争力的聚类性能. 对结果的进一步分析表明, MURLN 的聚类有效性来自于两个方面.

一方面, 所提出的多视图统一表示学习网络能够利用多视图中的互补性信息表示数据复杂的非线性相关性, 这对提升聚类结果至关重要. 观察表 3 在 Yale 人脸数据集上的聚类结果, 获得前三名的方法分别是: MURLN、D<sup>2</sup>MVSC 和 MSCNLG, 反映了使用深度神经网络非线性地提取特征可以实现更好的聚类效果. Extended YaleB 数据集由于包含不同光线强度下的人脸图像, 样本关系与 Yale 相比更加复杂. 可以发现, 在 ACC 和 NMI 指标上, MURLN 比传统浅层方法中最优的 DSS-MVC 分别提高了

18.75%和 16.36%. 结合其他数据集中的实验结果, MURLN 无论是处理图像数据还是 1-D 特征数据, 性能都较对比方法有明显提升, 验证了多视图统一表示学习网络在特征表示方面的出色能力.

另一方面, 模型使用 Schatten- $p$  范数正则化约束学到了更低秩的自表示系数矩阵, 以及更清晰的簇结构. 在 Caltech101-20 数据集的 ACC 和 NMI 指标上, MURLN 比次优方法 D<sup>2</sup>MVSC 分别提高了 7.46%和 4.52%. 这是因为加权 Schatten- $p$  范数赋予不同奇异值不同的权重, 重视较大奇异值的作用, 保证了子空间表示的低秩性. 与 D<sup>2</sup>MVSC 方法使用的核范数相比, 加权 Schatten- $p$  范数在求解低秩问题上具有更好的灵活性与有效性.

#### 3.2.2 共享编码器效果分析

为深入探究基于 Transformer 的共享编码器在复杂多视图场景中的优势, 本文合成一个具有 2 种模式 4 个视图的人造数据集 (MNIST-USPS-4V). 具体地,



在 MNIST-USPS 原始数据上增加两个分别由 MNIST 和 USPS 数据集提取的 LBP 特征. 设计三种使用不同编码器的 MURLN 模型用于实验比较, 分别是基于全连接编码器的 MURLN\_FC, 卷积自编码器的 MURLN\_CNN 和本文所提方法 MURLN. 表 6 提供三种不同模型在 MNIST-USPS 和 MNIST-USPS-4V 上的 ACC 和 NMI 结果. 通过观察发现: 1) MURLN\_CNN 通过卷积操作能够保留图像数据的空间信息, 不适用于处理一维特征. 2) MURLN\_FC 可以通过将图像拉伸成一维向量来编码图像数据, 但向量化会导致图像空间信息的丢失, 降低了聚类性能. 3) 相比之下, 本文所提的共享编码器不仅具备处理不同模态数据的能力, 还能提高统一表示的学习质量, 从而获得较好的聚类结果.

进一步, 使用峰值信噪比 (Peak Signal-to-noise Ratio, PSNR) 衡量本文方法与卷积自编码器对图像数据的表示能力. PSNR 值越高, 表明重建结果越好. 选取 RGB-D 和 MNIST-USPS 数据集上的部分图像进行实验, 结果为: 卷积自编码器的平均 PSNR 值为 35.25dB, MURLN 则达到了 54.36dB.

### 3.3 可视化分析

图 3-4 可视化了 MURLN 在 Extended YaleB (第一行) 和 BBCSport (第二行) 数据集上学习得到的亲和矩阵和 t-sne 图. 显然, 学习到的亲和矩阵具有理想的块对角结构和较少的噪声, 能清楚地揭示样本之间的关系. 随着迭代次数的增加, 聚类结构变得更加清晰, 进一步验证了本文方法的有效性.

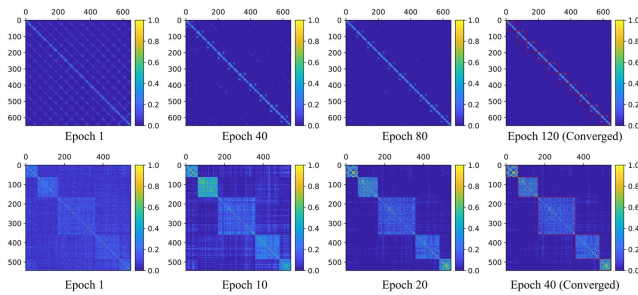


Fig. 3 Visualization of the affinity matrix

图 3 亲和矩阵可视化

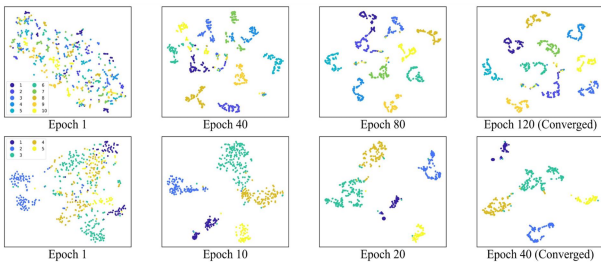


Fig. 4 Visualization of t-sne clustering

图 4 t-sne 聚类可视化

### 3.4 参数敏感性分析

MURLN 方法引入 4 个超参数 ( $\lambda_1, \lambda_2, \lambda_3, p$ ), 并在各数据集上以网格搜索策略来确定最优参数组合. 首先, 对于权衡超参数  $\lambda_1, \lambda_2, \lambda_3$ , 结合控制变量策略, 图 5 (a)-(i) 展示了在 Extend YaleB、Caltech101-20、BBCSport 数据集中, 使用不同超参数组合得到的聚类结果 (ACC 和 NMI). 以 BBCSport 为例, 当固定  $\lambda_3 = 0.001$ ,  $\lambda_1$  和  $\lambda_2$  在 ( $\{0.001, 0.005, 0.01\}$  和  $\{0.005, 0.01, 0.05, 0.1, 0.5\}$ ) 范围内取值, 可以实现较好的聚类效果. 综合其它数据集上的结果, 设置  $\lambda_1 = 0.01$ ,  $\lambda_2 = 0.5$ . 固定  $\lambda_1$  和  $\lambda_2$ , 调整  $\lambda_3$ , 如图 5 (g)-(i) 所示. 当  $\lambda_3$  的值在  $\{0.1, 0.5, 1.0\}$  范围内选择时, 聚类性能相对稳定. 经过实验分析, 本文设置  $\lambda_3 = 1.0$ .

图 6 提供了不同参数  $p$  值在 4 个数据集上的 ACC 和 NMI 结果.  $p$  以 0.1 的步长在  $\{0.1 \sim 1\}$  范围内取值. 通过观察可以发现, 参数  $p$  的选择受数据特性的影响. 例如, 对于 BBCSport 和 Caltech101-20 数据集,  $p$  设置为 0.7 和 0.8 可以实现最好的聚类性能. 对于 Extended YaleB 和 MSRCV1 数据集, 当  $p = 1$  时可以获得最优结果. 显然, 加权核范数作为秩函数的凸松弛并不总是近似求解秩最小化问题的最优选择, 说明本文利用 Schatten-p 范数设计低秩约束的合理性.

### 3.5 收敛性分析

图 5 (j)-(l) 展示了 MURLN 在 3 个数据集上的收敛曲线. 从图中可以看出, 随着迭代次数增加, 目标函数值逐渐下降并趋于稳定. 尽管在开始的几次迭代中发生振荡, 但最终在 80 次迭代左右收敛. 类似的, 聚类性能 (ACC 和 NMI) 在前几次迭代过程迅速升高并逐渐收敛. 表明 MURLN 具有良好的收敛性, 能够实现快速收敛. 在其他数据集上得到了类似的收敛分析.

### 3.6 消融研究

为了进一步分析不同模块对 MURLN 的贡献, 全面说明本文方法设计的合理性, 表 7 列出了对各模块进行消融实验的结果. “√” 表示 MURLN 模型包含该模块. 5 组实验表明, 整合所有模块的 MURLN 在所有数据集上均取得了最优性能, 每个模块都各自发挥作用, 去掉任何一个都会影响聚类效果.

## 4 结论

本文提出了一种新的多视图子空间聚类方法 (MURLN). 与现有方法为不同视图使用单独的编码器不同, MURLN 基于 Transformer 结构设计一个多视图统一表示学习网络, 所有视图各自经过特定的嵌入层来统一输入形式和维度, 然后共享同一个编码模

型. 此外, 提出基于加权 Schatten- $p$  范数的低秩正则化约束, 提高自表示系数矩阵学习的质量. 在多个公开数据集上的大量实验证明了所提方法的有效性和优越性.

利用全连接层实现数据自表达特性的方法普遍存在 full-batch training 的局限. 具体体现在随着样本

数量的增多, 全连接层的参数量迅速提升, 需要较多的计算和存储资源. 后续工作将关注于模型参数与样本数量的解耦研究, 提高大规模多视图数据集聚类的效率和精度. 此外, 由于数据收集和传输的复杂性, 现实场景中可能存在部分视图缺失现象<sup>[547]</sup>, 这也是未来的研究重点.

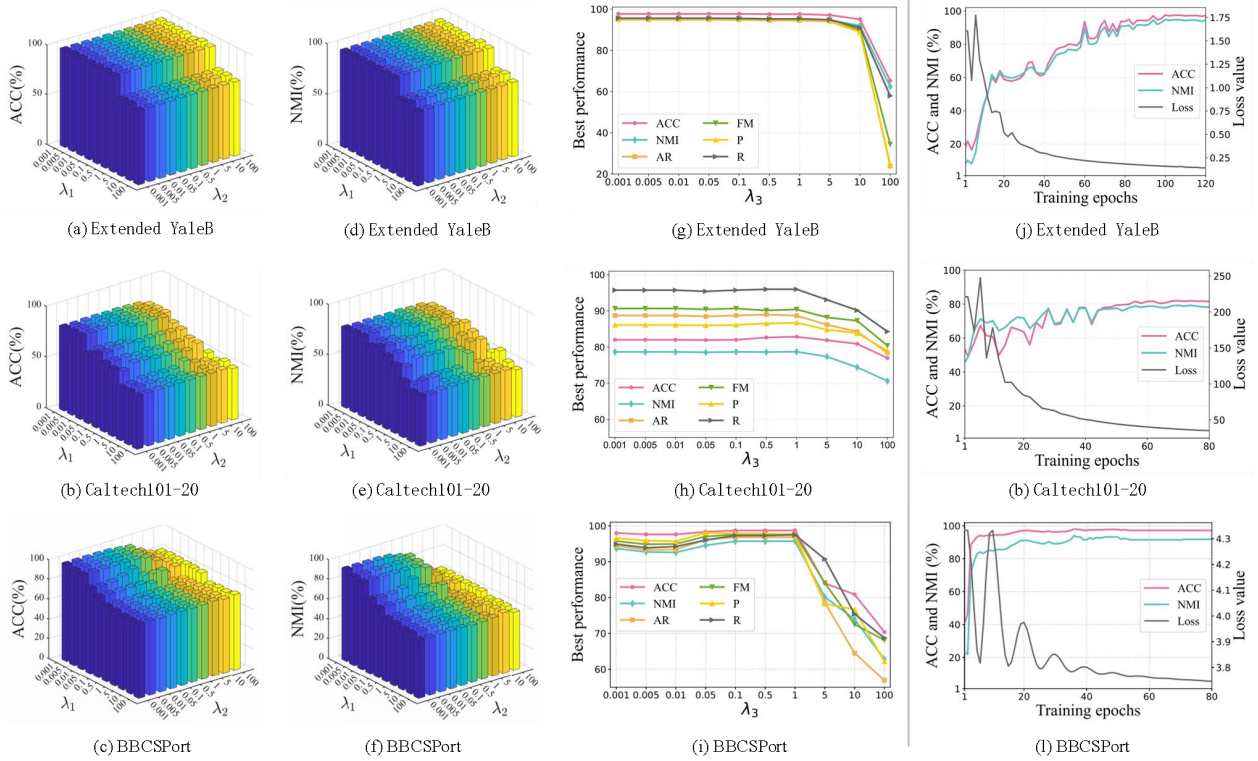


Fig. 5 Parameter sensitivity and convergence experiments on Extend YaleB, Caltech101-20, and BBCSport datasets

图 5 Extend YaleB、Caltech101-20、BBCSport 数据集上参数敏感性和收敛性实验

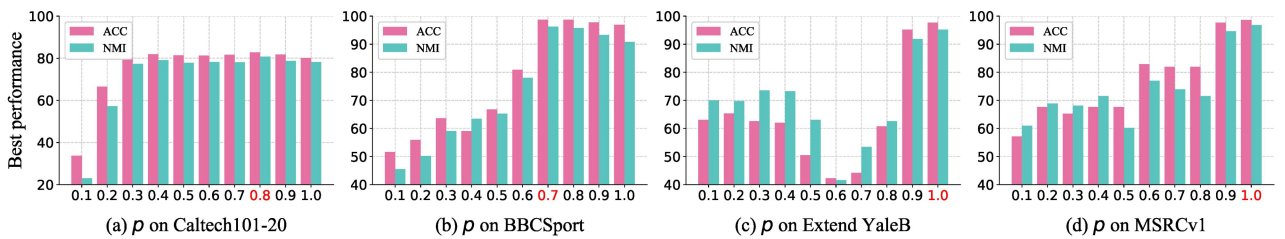


Fig. 6 Effect of parameter  $p$  on clustering performance

图 6 参数  $p$  对聚类性能的影响

Table 7 Ablation study of loss functions on MSRCV1 and Yale datasets

表 7 MSRCV1 和 Yale 数据集上损失函数的消融研究

	Components				MSRCV1				Yale			
	$\mathcal{L}_{re}$	$\mathcal{L}_f$	$\mathcal{L}_{self}$	$\mathcal{L}_{S_p}$	ACC	NMI	AR	FM	ACC	NMI	AR	FM
(a)	✓				87.13	78.59	76.43	75.86	86.78	83.61	80.29	80.05
(b)	✓	✓			92.86	91.37	87.21	88.03	93.16	91.21	85.79	84.73
(c)	✓		✓		95.51	89.97	91.76	90.44	92.37	92.01	86.19	88.32
(d)	✓			✓	94.77	90.40	88.84	89.16	91.53	89.86	84.15	88.35
(e)	✓	✓	✓	✓	<b>98.34</b>	<b>95.99</b>	<b>96.41</b>	<b>96.07</b>	<b>94.95</b>	<b>93.88</b>	<b>89.37</b>	<b>90.59</b>

**作者贡献声明:** 林毓秀负责网络模型的设计和实现, 论文的撰写和修改; 刘慧提出研究方向, 把握论文的创新性, 并指导论文修改; 于晓负责实验的整理和可视化分析; 张彩明提出了指导意见并修改论文。

## 参 考 文 献

- [1] Liu Xinwang. Research on the fundamental theory and methods of multi-view learning[J]. China Basic Science, 2022, 24(03): 27-34 (in Chinese)  
(刘新旺. 多视图学习的基础理论和方法研究[J]. 中国基础科学, 2022, 24(03): 27-34)
- [2] Liang Jiye, Liu Xiaolin. Research progress and prospects of multi-view clustering[J]. Journal of Shanxi University(Nat. Sci. Ed.), 2022, 45(03): 612-621 (in Chinese)  
(梁吉业, 刘晓琳. 多视图聚类研究进展与展望[J]. 山西大学学报(自然科学版), 2022, 45(03): 612-621)
- [3] Li Jinxing, Zhou Chuhao, Ji Xiaoqiang, et al. Multi-view instance attention fusion network for classification[J]. Information Fusion, 2024, 101: 101974
- [4] Zhang Dongxu, Yan Yang, Huang Yulin, et al. Unsupervised Cryo-EM images denoising and clustering based on deep convolutional autoencoder and K-Means++[J]. IEEE Transactions on Medical Imaging, 2023, 42(5): 1509-1521
- [5] Gong Weihua, Jin Rong, Pei Xiaobing, et al. Collaborative recommendation method based on community co-clustering in location based socail networks[J]. Journal of Computer Research and Development, 2019, 56(11): 2506-2517 (in Chinese)  
(龚卫华, 金蓉, 裴小兵, 等. LBSN 中基于社区联合聚类的协同推荐方法[J]. 计算机研究与发展, 2019, 56(11): 2506-2517)
- [6] Cao Congqi, Zhang Xin, Zhang Shizhou, et al. Weakly supervised video anomaly detection based on cross-batch clustering guidance[C]//Proc of the International Conference on Multimedia and Expo. Piscataway, NJ: IEEE, 2023: 2723-2728
- [7] Kumar A, Rai P, Daumé H. Co-regularized Multi-view Spectral Clustering[C]//Proc of the Annual Conference on Neural Information Processing Systems. Cambridge, MA: MIT, 2011: 1413-1421
- [8] Xia Dongxue, Yang Yan, Wang Hao, et al. Late fusion multi-view clustering based on local multi-kernel learning[J]. Journal of Computer Research and Development, 2020, 57(08): 1627-1638 (in Chinese)  
(夏冬雪, 杨燕, 王浩, 等. 基于邻域多核学习的后融合多视图聚类算法[J]. 计算机研究与发展, 2020, 57(08): 1627-1638)
- [9] Tang Chang, Liu Xinwang, Zhu Xinzong, et al. CGD: multi-view clustering via cross-view graph diffusion[C]//Proc of the AAAI Conference on Artificial Intelligence. Palo Alto, CA: AAAI, 2020: 5924-5931
- [10] Zhang Changqing, Hu Qinghua, Fu Huazhu, et al. Latent multi-view subspace clustering[C]//Proc of the Conference on Computer Vision and Pattern Recognition. Piscataway, NJ: IEEE, 2017: 4333-4341
- [11] Zhang Guangyu, Huang Dong, Wang Changdong. Facilitated low-rank multi-view subspace clustering[J]. Knowledge-Based Systems, 2022, 260: 110141
- [12] Wang Qianqian, Cheng Jiafeng, Gao Quanxue, et al. Deep multi-view subspace clustering with unified and discriminative learning[J]. IEEE Transactions on Multimedia, 2021, 23: 3483-3493
- [13] Liu Zhaohu, Song Peng. Deep low-rank tensor embedding for multi-view subspace clustering[J]. Expert Systems with Applications, 2023, 237: 121518
- [14] Wang Jing, Feng Songhe, Lyu Gengyu, et al. Triple-granularity contrastive learning for deep multi-view subspace clustering[C]//Proc of the ACM International Conference on Multimedia. New York: ACM, 2023: 2994-3002
- [15] Rumelhart D E, Hinton G E, Williams R J. Learning representations by back-propagating errors[J]. Nature, 1986, 323(6088): 533-536
- [16] Abavisani M, Patel V M. Deep multimodal subspace clustering networks[J]. IEEE Journal of Selected Topics in Signal Processing, 2018, 12(6): 1601-1614
- [17] Zhu Pengfei, Hui Binyuan, Zhang Changqing, et al. Multi-view deep subspace clustering networks[J]. arXiv preprint, arXiv: 1908.01978, 2019
- [18] Zhang Changqing, Liu Yeqing, Fu Huazhu. AE2-Nets: autoencoder in Autoencoder networks[C]//Proc of the Conference on Computer Vision and Pattern Recognition. Piscataway, NJ: IEEE, 2019: 2572-2580
- [19] Li Kai, Liu Hongfu, Zhang Yulun, et al. Self-guided deep multiview subspace clustering via consensus affinity regularization[J]. IEEE Transactions on Cybernetics, 2022, 52(12): 12734-12744
- [20] Cui Chenhong, Ren Yazhou, Pu Jingyu, et al. Deep multi-view subspace clustering with anchor graph[C]//Proc of the International Joint Conference on Artificial Intelligence. San Francisco, CA: Morgan Kaufmann, 2023: 3577-3585
- [21] Vaswani A, Shazeer N, Parmar N, et al. Attention is All you Need[C]//Proc of the Annual Conference on Neural Information Processing Systems. Cambridge, MA: MIT, 2017: 5998-6008
- [22] Xu Peng, Zhu Xiatian, Clifton D A. Multimodal learning with Transformers: a survey[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2023, 45(10): 12113-12132
- [23] Han Xue, Wang Yitong, Feng Junlan, et al. A survey of transformer-based multimodal pre-trained modals[J]. Neurocomputing, 2023, 515: 89-106
- [24] Xie Yuan, Gu Shuhang, Liu Yan, et al. Weighted Schatten p-norm minimization for image denoising and background subtraction[J]. IEEE Transactions on Image Processing, 2016, 25(10): 4842-4857
- [25] Elhamifar E, Vidal R. Sparse subspace clustering: algorithm, theory, and applications[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2013, 35(11): 2765-2781
- [26] Liu Guangcan, Lin Zhouchen, Yan Shuicheng, et al. Robust recovery of

- subspace structures by low-rank representation[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2013, 35(1): 171-184
- [27] Shi Jianbo, Malik J. Normalized cuts and image segmentation[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2000, 22(8): 888-905
- [28] Li Ruihuang, Zhang Changqing, Fu Huazhu, et al. Reciprocal multi-layer subspace learning for multi-view clustering[C]//Proc of the IEEE International Conference on Computer Vision. Piscataway, NJ: IEEE, 2019: 8171-8179
- [29] Zhou Tao, Zhang Changqing, Peng Xi, et al. Dual shared-specific multiview subspace clustering[J]. IEEE Transactions on Cybernetics, 2020, 50(8): 3517-3530
- [30] Cao Xiaochun, Zhang Changqing, Fu Huazhu, et al. Diversity-induced multi-view subspace clustering[C]//Proc of the Conference on Computer Vision and Pattern Recognition. Piscataway, NJ: IEEE, 2015: 586-594
- [31] Zheng Qinghai, Zhu Jihua, Ma Yuanyuan, et al. Multi-view subspace clustering networks with local and global graph information[J]. Neurocomputing, 2021, 449: 15-23
- [32] Wang Jiao, Wu Bin, Ren Zhenwen, et al. Decomposed deep multi-view subspace clustering with self-labeling supervision[J]. Information Sciences, 2024, 653: 119798
- [33] Zhang Changqing, Fu Huazhu, Hu Qinghua, et al. Generalized latent multi-view subspace clustering[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2020, 42(1): 86-99
- [34] Wang Qianqian, Tao Zhiqiang, Gao Quanxue, et al. Multi-view subspace clustering via structured multi-pathway network[J]. IEEE Transactions on Neural Networks and Learning Systems, 2022
- [35] Masci J, Meier U, Ciresan D C, et al. Stacked convolutional auto-encoders for hierarchical feature extraction[C]//Proc of the International Conference on Artificial Neural Networks. Berlin: Springer, 2011: 52-59
- [36] Ren Shaoqing, He Kaiming, Girshick R B, et al. Faster R-CNN: towards real-time object detection with region proposal networks[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2017, 39(6): 1137-1149
- [37] Devlin J, Chang Mingwei, Lee K, et al. BERT: pre-training of deep bidirectional transformers for language understanding[J]. arXiv preprint, arXiv: 1810.04805, 2018
- [38] Yan Xiaoqiang, Mao Yiqiao, Ye Yangdong, et al. Cross-modal clustering with deep correlated information bottleneck method[J]. IEEE Transactions on Neural Networks and Learning Systems, 2023
- [39] Wang Qianqian, Lian Huanhuan, Sun Gan, et al. iCmSC: incomplete cross-modal subspace clustering[J]. IEEE Transactions on Image Processing, 2020, 30: 305-317
- [40] Lu Runkun, Liu Jianwei, Zuo Xin. Attentive multi-view deep subspace clustering net[J]. Neurocomputing, 2021, 435: 186-196
- [41] Xie Yuan, Tao Dacheng, Zhang Wensheng, et al. On unifying multi-view self-representations for clustering by tensor multi-rank minimization[J]. International Journal of Computer Vision, 2018, 126(11): 1157-1179
- [42] Yu Xiao, Liu Hui, Lin Yuxiu, et al. Consensus guided auto-weighted multi-view clustering[J]. Journal of Computer Research and Development, 2022, 59(7): 1496-1508 (in Chinese).
- (于晓, 刘慧, 林毓秀, 等. 一致性引导的自适应加权多视图聚类[J]. 计算机研究与发展, 2022, 59(7): 1496-1508)
- [43] Zheng Qinghai, Zhu Jihua, Li Zhongyu, et al. Feature concatenation multi-view subspace clustering[J]. Neurocomputing, 2019, 379: 89-102
- [44] Dosovitskiy A, Beyer L, Kolesnikov A, et al. An image is worth 16x16 words: transformers for image recognition at scale[C]//Proc of the International Conference on Learning Representations. Washington DC: ICLR, 2021
- [45] Xie Yuan, Qu Yanyun, Tao Dacheng, et al. Hyperspectral image restoration via iteratively regularized weighted Schatten p-norm minimization[J]. IEEE Transactions on Geoscience and Remote Sensing, 2016, 54(8): 4642-4659
- [46] Nie Feiping, Huang Heng, Ding C H Q. Low-rank matrix recovery via efficient Schatten p-norm minimization[C]//Proc of the AAAI Conference on Artificial Intelligence. Palo Alto, CA: AAAI, 2012. 655-661
- [47] Zhao Boyu, Zhang Changqing, Chen Lei, et al. Generative model for partial multi-view clustering[J]. Acta Automatica Sinica, 2021, 47(8): 1867-1875 (in Chinese)
- (赵博宇, 张长青, 陈蕾, 等. 生成式不完整多视图数据聚类[J]. 自动化学报, 2021, 47(8): 1867-1875)



**Lin Yuxiu**, born in 1996. Ph.D., candidate. Her main research interests include data mining and visualization, multi-view learning and machine learning.

林毓秀, 1996 年生, 博士研究生. 主要研究方向为数据挖掘与可视化、多视图学习和机器学习.



**Liu Hui**, born in 1978, Ph.D., professor, PhD supervisor. Senior member of CCF. Her main research interests include machine learning, data mining and image processing.

(liuh\_lh@sdufe.edu.cn)

刘慧, 1978 年生, 博士, 教授, 博士生导师. CCF 会员. 主要研究方向为机器学习、数据挖掘和图像处理.



**Yu Xiao**, born in 1977. PhD, associated professor. Member of China Computer Federation. Her research interests include statistical machine learning and data mining.

于晓, 1977 年生, 博士, 副教授. CCF 会员. 主要研究方向为统计机器学习和数据挖掘.



**Zhang Caiming**, born in 1977. PhD, professor, PhD supervisor. Member of China Computer Federation. His main research interests include CAGD&CG, information visualization and

image processing.

张彩明, 1977 年生, 博士, 教授, 博士生导师. CCF 会员. 主要研究方向为 CAGD&CG、信息可视化和图像处理.