

基于改进松弛嵌入空间的多视图聚类

张瑛

(台州科技职业学院汽车与信息工程学院 浙江 台州 318000)

摘要 针对传统聚类方法缺乏统一特征表示,存在保守性的缺陷,提出一种基于改进松弛嵌入空间的多视图聚类方法。在统一的框架下联合学习一个综合的潜在嵌入表示矩阵、全局相似矩阵和一个精确指标矩阵。进一步放松全局相似矩阵的约束,并在此基础上提出一种改进的松弛多视图聚类嵌入空间,使得该方法具有更低的计算复杂度和更多的数据点对之间的相关性。实验结果表明,该方法能够获得鲁棒性更强、准确度更高的聚类结果。

关键词 多视图聚类 嵌入空间 相似矩阵 松弛因子

中图分类号 TP393.06

文献标志码 A

DOI:10.3969/j.issn.1000-386x.2024.04.041

MULTIPLE VIEW CLUSTERING BASED ON IMPROVED SLACK EMBEDDING SPACE

Zhang Ying

(School of Automotive and Information Engineering, Taizhou Vocational College of Science & Technology, Taizhou 318000, Zhejiang, China)

Abstract In view of the lack of unified feature representation and defects of conservatism of traditional clustering methods, a multiple view clustering method based on improved slack embedding space is proposed. In a unified framework, a comprehensive potential embedding representation matrix, a global similarity matrix and an accurate index matrix were jointly learned. Furthermore, the constraint of global similarity matrix was slack, and an improved slack multiple view clustering embedding space was proposed, which made the proposed method have lower computational complexity and more correlation between data point pairs. The experimental results show that the proposed method can obtain more robust and more accurate clustering results.

Keywords Multiple view clustering Embedding space Similarity matrix Slack factor

0 引言

多视图聚类是近十年来机器学习领域的一个研究热点。在多视图中,同一实例可以由从多个资源或不同特征子集获得的多个视图来表示^[1]。例如,在一个网页中,不同类型的数据如文本、视频和图像,可以被考虑在内,因为它们是网页的不同方面。考虑到多视图的多样性,研究如何有效地集成此类数据是非常必要的。

处理多视图数据的一种简单方法是将所有特征串联成一个新的特征向量,然后将该特征向量输入到单视图聚类方法中,得到最终的聚类结果^[2]。然而,这种聚类策略忽略了各个视图之间的不同特性以及相关

性。Chen等^[3]提出了一种鲁棒的粗糙模糊C-均值多视图聚类算法,它综合了模糊集的概率和可能性成员关系,既能处理噪声环境中的重叠聚类,又能处理粗糙集在聚类定义中的不确定性和模糊性。洪敏等^[4]基于粗糙集理论提出了一种粗糙C-均值聚类方法,将每个簇分别用一个上下近似原型描述,其中边界区域被定义为上下近似的差分。但是,这些方法会受到原始视图质量差的影响。因为目前大多数方法直接从原始数据集获取聚类特征,其中可能存在大量的噪声和离群点,可能导致聚类性能下降。

为了处理潜在的噪声,刘良凤等^[5]通过将标准子空间聚类扩展到多视图情况,学习每个子空间表示,并采用希尔伯特-施密特独立性准则(HSIC)探索多个子空间表示的多样性。夏冬雪等^[6]提出同时恢复底层多

视图子空间和与不同视图相关的投影的方法。Chen 等^[7]采用典型相关分析(CCA)将多个视图投影到一个低维子空间,然后利用学习的表示进行数据聚类。尽管这些方法取得了很好的聚类效果,但是它们仍然缺乏发现多视图数据的统一特征表示的能力。另外,上述方法大多倾向于将相似度矩阵的构造和聚类指标矩阵的计算分别考虑,不能在统一的框架下同时进行,因此存在一定的保守性。

针对上述局限性,提出一种基于改进松弛嵌入空间的多视图聚类方法。在统一的框架下联合学习一个综合的潜在嵌入表示矩阵、全局相似矩阵和一个精确指标矩阵。进一步提出基于潜在嵌入空间的松弛多视图聚类算法(R-MCLES),避免了二次规划优化问题,提高了数据点对之间的相关性。

1 算法

1.1 总体概述

首先简要介绍一下本文算法的框架,整体结构如图 1 所示。

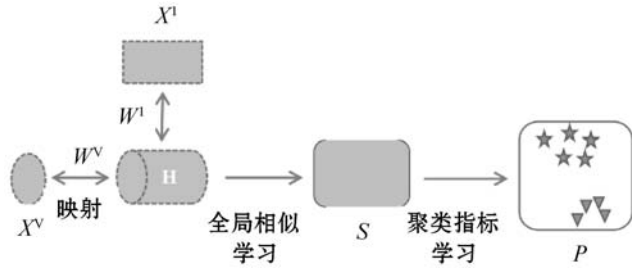


图 1 整体结构

在多视图聚类中,输入是一组包含 n 个样本的多视图观测值,用 V 个不同的视图表示,即 $\mathbf{X} = [\mathbf{X}^1; \mathbf{X}^2; \dots; \mathbf{X}^V] \in \mathbf{R}_{v=1}^{\sum_{v=1}^V d^v \times n}$, 其中 $\mathbf{X}^V \in \mathbf{R}^{d^v \times n}$ 是第 V 个视图的特征矩阵, d^v 是维数。本文方法目的是发现一个潜在的嵌入表示,它对来自多个视图的互补信息进行编码,同时获得全局相似性矩阵和聚类指标。根据多视点观测,每个数据点的潜在嵌入表示 $\mathbf{h}_i \in \mathbf{R}^{d \times 1}$ 应该以下述方式来推断,即所有的多视图都是从潜在嵌入表示矩阵 $\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n] \in \mathbf{R}^{d \times n}$ 中提取,通过自我表达学习可以得到全局相似矩阵。因此,一般目标函数可以表示为:

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{S}} L_H(\mathbf{X}, \mathbf{WH}) + \alpha L_S(\mathbf{H}, \mathbf{HS}) + \beta \Omega(\mathbf{S}) \quad (1)$$

式中: $\mathbf{W} = [\mathbf{W}^1; \mathbf{W}^2; \dots; \mathbf{W}^V] \in \mathbf{R}_{v=1}^{\sum_{v=1}^V d^v \times d}$ 是映射矩阵,有助于从潜在嵌入表示 \mathbf{H} 中恢复多视点观测值; $\mathbf{S} \in \mathbf{R}^{n \times n}$ 是全局相似矩阵;参数 α 和 β 用来平衡三项:第一

个是原始观测值和潜在嵌入表示的损失函数,第二项是基于潜在嵌入表示的自我表达学习,最后一项是关于全局结构的正则化。给定一个数据集,然后对该数据集进行映射得到潜在嵌入表示矩阵,再进行全局相似学习得到全局相似矩阵,最后通过聚类指标学习得到相应的聚类结果。

本文提出潜在嵌入空间中的松弛多视图聚类(R-MCLES),在同一个优化框架中同时学习潜在嵌入表示、全局相似性矩阵和聚类指标。

1.2 潜在嵌入空间中的多视图聚类

1.2.1 MCLES 模型

根据一个潜在表示可以提取多个视图的假设,本文考虑多个视图之间较低维的潜在嵌入空间进行数据聚类,而不是直接使用原始的单视图特征空间。首先,利用隐式嵌入表示矩阵 \mathbf{H} 之间的三角关系、多视图 \mathbf{X} 的观测值、对应的映射矩阵 $\mathbf{W} = \{\mathbf{W}^v\}_{v=1}^V$ 。

$$\mathbf{x}_i^v = \mathbf{W}^v \mathbf{h}_i \quad (2)$$

通过式(2)计算 \mathbf{H} ,则对应于求解以下优化问题:

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}} \quad & \|\mathbf{X} - \mathbf{WH}\|_F^2 \\ \text{s. t.} \quad & \|\mathbf{W}_{:,j}\|_2 \leq 1 \end{aligned} \quad (3)$$

式中: \mathbf{X} 和 \mathbf{W} 分别是多视图观测值和映射矩阵。引入 \mathbf{W} 是为了防止 \mathbf{H} 变得太小。一般而言,通过结合多个视图的互补信息,隐式嵌入表示矩阵 \mathbf{H} 会使信息更加全面。与大多数基于数据 \mathbf{X} 的原始特征空间构造相似矩阵 \mathbf{S} 的方法不同,该方法从潜在嵌入表示矩阵 \mathbf{H} 中学习相似矩阵 \mathbf{S} ,有效地提高了学习相似矩阵的鲁棒性和准确性。根据全局相似性学习,有

$$\begin{aligned} \min_{\mathbf{S}} \quad & \|\mathbf{H} - \mathbf{HS}\|_F^2 + \beta \|\mathbf{S}\|_F^2 \\ \text{s. t.} \quad & \mathbf{S}_{:,j}^T \mathbf{1} = 1, \mathbf{0} \leq \mathbf{S} \leq \mathbf{1} \end{aligned} \quad (4)$$

式中: β 是权衡参数。理想状态下, \mathbf{S} 中相关参数数量与数据集 \mathbf{X} 的簇数应该相同,即 c 个。换言之, \mathbf{S} 是具有一定排列方式的块对角矩阵。然而,式(4)中的解可能不满足期望的结果。为了解决这一问题,引入秩。

引理 1^[7] \mathbf{S} 的拉普拉斯矩阵 \mathbf{L}_S 的特征值 0 的数量 c 等于具有相似矩阵 \mathbf{S} 的图中连通分量的个数。

引理 1 表明,如果相似矩阵 \mathbf{S} 完全由 c 连通分量组成,则秩(\mathbf{L}_S) = $n - c$ 。因此,式(4)中的问题可以改写为

$$\begin{aligned} \min_{\mathbf{S}} \quad & \|\mathbf{H} - \mathbf{HS}\|_F^2 + \beta \|\mathbf{S}\|_F^2 \\ \text{s. t.} \quad & \mathbf{S}_{:,j}^T \mathbf{1} = 1, \mathbf{0} \leq \mathbf{S} \leq \mathbf{1}, \text{rank}(\mathbf{L}_S) = n - c \end{aligned} \quad (5)$$

通过将式(3)中的潜在嵌入学习和式(5)中的全局相似性学习整合到一个统一的框架中,可以得到:

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}, \mathbf{S}} \quad & \|\mathbf{X} - \mathbf{WH}\|_F^2 + \alpha \|\mathbf{H} - \mathbf{HS}\|_F^2 + \beta \|\mathbf{S}\|_F^2 \\ \text{s. t.} \quad & \|\mathbf{W}_{:,j}\|_2 \leq 1, \mathbf{S}_{:,j}^T \mathbf{1} = 1, \mathbf{0} \leq \mathbf{S} \leq \mathbf{1}, \text{rank}(\mathbf{L}_S) = n - c \end{aligned} \quad (6)$$

式中: $\alpha > 0$ 和 $\beta > 0$ 是权衡参数。事实上,解决式(6)中的问题有些困难,因为 $L_S = D - \frac{S^T + S}{2}$, 其中 $D \in$

$\mathbf{R}^{n \times n}$ 是第 i 个对角元素为 $\sum_j \frac{s_{ij} + s_{ji}}{2}$ 的对角矩阵。

由于 L_S 是半正定的,有 $\sigma_i(L_S) \geq 0$, 其中 $\sigma_i(L_S)$ 代表 L_S 的第 i 个最小特征值。众所周知,具有秩约束的优化问题具有组合复杂性。为了解决这个问题,可以在目标函数中加入秩约束作为正则项。根据文献[8], $\text{rank}(L_S) = n - c$ 相当于 $\sum_{i=1}^c \sigma_i(L_S) = 0$ 。因此,简化了问题,式(6)中的问题等价于:

$$\begin{aligned} \min_{W, H, S} \quad & \|X - WH\|_F^2 + \alpha \|H - HS\|_F^2 + \beta \|S\|_F^2 + \gamma \sum_{i=1}^c \sigma_i(L_S) \quad (7) \\ \text{s. t.} \quad & \|W_{:,j}\|_2 \leq 1, S_{:,j}^T \mathbf{1} = 1, \mathbf{0} \leq S \leq \mathbf{1} \end{aligned}$$

如果 γ 值非常大,上述最小化会使正则化项 $\sum_{i=1}^c \sigma_i(L_S) \rightarrow 0$, 满足约束 $\text{rank}(L_S) = n - c$ 。尽管如此,式(7)中的优化问题由于处理最后一个项有些困难,为了解决这个问题,根据文献[9]可知:

$$\sum_{i=1}^c \sigma_i(L_S) = \min_{P^T P} \text{Tr}(P^T L_S P) \quad (8)$$

式中: $P \in \mathbf{R}^{n \times c}$ 是聚类指标矩阵。因此,式(7)可以描述为:

$$\begin{aligned} \min_{W, H, S, P} \quad & \underbrace{\|X - WH\|_F^2}_{\text{潜在嵌入学习}} + \underbrace{\alpha \|H - HS\|_F^2 + \beta \|S\|_F^2}_{\text{全局优化学习}} + \underbrace{\gamma \text{Tr}(P^T L_S P)}_{\text{聚类指标学习}} \quad (9) \\ \text{s. t.} \quad & \|W_{:,j}\|_2 \leq 1, S_{:,j}^T \mathbf{1} = 1, \mathbf{0} \leq S \leq \mathbf{1}, P^T P = I \end{aligned}$$

式中:三个权衡参数 $\alpha > 0, \beta > 0$ 和 $\gamma > 0$ 来平衡以下四个项:第一项是建立潜在嵌入表示矩阵 H 和每个映射矩阵 W^r 的模型,以重新生成观测值;第二项是惩罚相似性学习中的结构错误;第三项用于避免没必要的解 $S = I$;最后一项保证相似矩阵满足秩约束,直接得到聚类指标矩阵 P 。潜在嵌入学习由多个视图之间的互补信息进行完善,并通过全局相似性学习和聚类指标学习加以改进,潜在嵌入学习有助于全局相似性学习且可以改进聚类指标学习,全局相似性学习有助于聚类指标学习且可以改进潜在嵌入学习。事实上,由于潜在嵌入表示、相似性矩阵和聚类指标矩阵之间的相互作用,在统一框架中存在着相互自学的特性。

1.2.2 MCLES 优化

引入了一种交替优化方案来解决式(9)中的问题。更新权值 W , 通过修正除 W 外的所有变量,可将式(9)中的问题简化为解决式(10)问题。

$$\begin{aligned} \min_W \quad & \|X - WH\|_F^2 \\ \text{s. t.} \quad & \|W_{:,j}\|_2 \leq 1 \end{aligned} \quad (10)$$

可以通过引入一个变量 $G \in \mathbf{R}^{\sum_{v=1}^V d^v \times d}$ 来优化式(10)问题:

$$\begin{aligned} \min_{W, G} \quad & \|X - WH\|_F^2 \\ \text{s. t.} \quad & W = G, \|G_{:,j}\|_2 \leq 1 \end{aligned} \quad (11)$$

式(11)的最优解可通过交替方向乘子算法(ADMM)获得:

$$\begin{cases} W^{r+1} = \arg \min_W \|X - WH\|_F^2 + \rho \|W - G^r + T^r\|_F^2 \\ G^{r+1} = \arg \min_G \rho \|W^{r+1} - G + T^r\|_F^2 \quad \text{s. t.} \quad \|G_{:,j}\|_2 \leq 1 \\ T^{r+1} = T^r + W^{r+1} - G^{r+1} \quad \text{更新 } \rho \end{cases} \quad (12)$$

式中: r 为迭代步数; $T \in \mathbf{R}^{\sum_{v=1}^V d^v \times d}$ 是中间变量。在优化过程中的每一步,由于 ADMM 算法具有良好的收敛性能,基于 ADMM 的 W 优化算法收敛速度快。更新 H , 通过修正除 H 外的所有变量,式(9)中的问题相当于求解:

$$\min_H \|X - WH\|_F^2 + \alpha \|H - HS\|_F^2 \quad (13)$$

通过将式(13)中的 H 设零,即可得到最优解 H^* , 且满足:

$$W^T W H^* + H^* \times \alpha (I - S)(I - S)^T = W^T X \quad (14)$$

更新 S : 通过改变除 S 外的所有变量,式(9)中的问题等价于:

$$\begin{aligned} \min_S \quad & \|H - HS\|_F^2 + \frac{\beta}{\alpha} \|S\|_F^2 + \frac{\gamma}{\alpha} \text{Tr}(P^T L_S P) \quad (15) \\ \text{s. t.} \quad & S_{:,j}^T \mathbf{1} = 1, \mathbf{0} \leq S \leq \mathbf{1} \end{aligned}$$

为了方便起见,引入一个变量 $K \in \mathbf{R}^{n \times n}$, 相当于 $H^T H$ 。因此,式(15)中的问题等价于:

$$\begin{aligned} \min_S \quad & \text{Tr}(K - 2KS + S^T K S P) + \frac{\beta}{\alpha} \|S\|_F^2 + \frac{\gamma}{\alpha} \text{Tr}(P^T L_S P) \\ \text{s. t.} \quad & S_{:,j}^T \mathbf{1} = 1, \mathbf{0} \leq S \leq \mathbf{1} \end{aligned} \quad (16)$$

式(16)中的问题可以按列优先的准则为:

$$\begin{aligned} \min_{S_{:,i}} \quad & K_{ii} - 2K_{:,i} + S_{:,i}^T K S_{:,i} + \frac{\beta}{\alpha} S_{:,i}^T S_{:,i} + \frac{\gamma}{2\alpha} \mathbf{b}_i^T S_{:,i} \quad (17) \\ \text{s. t.} \quad & S_{:,i}^T \mathbf{1} = 1, \mathbf{0} \leq S \leq \mathbf{1} \end{aligned}$$

式中: $\mathbf{b}_i \in \mathbf{R}^{n \times 1}$ 是第 j 个元素 \mathbf{b}_{ij} 为 $a = B$ 的列向量 $\mathbf{b}_j = \|\mathbf{P}_{i,:} - \mathbf{P}_{:,j}\|^2$ 。式(17)可以进一步简化为:

$$\begin{aligned} \min_{S_{:,i}} \quad & S_{:,i}^T \left(\frac{\beta}{\alpha} I + K \right) S_{:,i} + \left(\frac{\gamma \mathbf{b}_i^T}{2\alpha} - 2K_{i,:} \right) S_{:,i} \quad (18) \\ \text{s. t.} \quad & S_{:,i}^T \mathbf{1} = 1, \mathbf{0} \leq S \leq \mathbf{1} \end{aligned}$$

式(18)中的问题可以用许多现有的二次规划方法来求解。更新 P 。通过更改除 P 以外的所有变量,式(9)中的问题可以改写为:

$$\min_P \text{Tr}(P^T L_S P), P^T P = I \quad (19)$$

由对应于 c 个最小特征值的 L_S 的 c 个特征向量可

以得到最优解 \mathbf{P} 。通过交替优化方法,变量 \mathbf{W} 、 \mathbf{H} 、 \mathbf{S} 和 \mathbf{P} 经过迭代进行更新,直到收敛。提出的 MCLES 方法的步骤如算法 1 所示。

算法 1 MCLES

输入:多视图矩阵 $\mathbf{X} = [\mathbf{X}^1; \mathbf{X}^2; \dots; \mathbf{X}^v]$, 聚类数 c , 参数 α 、 β 和 γ , 潜在表达的嵌入维度 d 。

输出: \mathbf{W} 、 \mathbf{H} 、 \mathbf{S} 和 \mathbf{P} 。

初始化: $\mathbf{W} = \mathbf{0}$, $\mathbf{S} = \mathbf{0}$ 和 $r = 1$, 对于 \mathbf{H} 和 \mathbf{P} 随机赋值。

1. repeat
2. repeat
3. $r \leftarrow r + 1$
4. 根据式(12)更新 \mathbf{W}

5. until 收敛

6. 根据式(14)更新 \mathbf{H}

7. 对于每一个 i , 通过解式(18)更新 \mathbf{S} 的第 i 列值

8. 更新 \mathbf{P} , 由对应于 c 最小特征值的 $\mathbf{L}_s = \mathbf{D} - \frac{\mathbf{S}^T + \mathbf{S}}{2}$ 的 c 特征向量构成

9. until

1.3 潜在嵌入空间中的松弛多视图聚类

1.3.1 R-MCLES 模型

一般来说,一旦二次规划算法运用到其中一种方法的优化中,就会产生很高的计算复杂度,就像上面提出的 MCLES 方法一样。为了避免二次规划的优化问题,提出一种基于潜在嵌入空间的松弛多视图聚类算法(R-MCLES),减轻了对全局相似矩阵的约束。因此,R-MCLES 方法的目标函数可以重新表示为:

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{S}, \mathbf{P}} \underbrace{\|\mathbf{X} - \mathbf{WH}\|_F^2}_{\text{潜在嵌入学习}} + \alpha \underbrace{\|\mathbf{H} - \mathbf{HS}\|_F^2}_{\text{全局相似学习}} + \beta \underbrace{\|\mathbf{S}\|_F^2}_{\text{聚类指标学习}} + \gamma \text{Tr}(\mathbf{P}^T \mathbf{L}_s \mathbf{P}) \quad (20)$$

$$\text{s. t. } \|\mathbf{W}_{:,j}\|_2 \leq 1, \mathbf{S} \geq \mathbf{0}, \text{diag}(\mathbf{S}) = \mathbf{0}, \mathbf{P}^T \mathbf{P} = \mathbf{I}$$

1.3.2 R-MCLES 优化

具体而言,除了全局相似矩阵 \mathbf{S} 的变量外,相关变量的优化算法与 MCLES 算法相同。因此, \mathbf{S} 更新步骤如下。通过修改除了 \mathbf{S} 以外的全部变量,式(20)描述的问题等价于:

$$\min_{\mathbf{S}} \|\mathbf{H} - \mathbf{HS}\|_F^2 + \frac{\beta}{\alpha} \|\mathbf{S}\|_F^2 + \frac{\gamma}{\alpha} \text{Tr}(\mathbf{P}^T \mathbf{L}_s \mathbf{P}) \quad (21)$$

$$\text{s. t. } \mathbf{S}_{:,i}^T \geq \mathbf{0}, \text{diag}(\mathbf{S}) = \mathbf{0}$$

此外,通过引入变量 $\hat{\mathbf{K}} \in \mathbf{R}^{n \times n}$, 式(21)等价于:

$$\min_{\mathbf{S}} \text{Tr}(\hat{\mathbf{K}} - 2\hat{\mathbf{K}}\mathbf{S} + \mathbf{S}^T \hat{\mathbf{K}}\mathbf{S}) + \frac{\beta}{\alpha} \|\mathbf{S}\|_F^2 + \frac{\gamma}{\alpha} \text{Tr}(\mathbf{P}^T \mathbf{L}_s \mathbf{P}) \quad (22)$$

$$\text{s. t. } \mathbf{S} \geq \mathbf{0}, \text{diag}(\mathbf{S}) = \mathbf{0}$$

式中: $\hat{\mathbf{K}}$ 代表 $\mathbf{H}^T \mathbf{H}$ 的 tanh 函数, 即 $F(\mathbf{H}^T \mathbf{H}) = \tanh(a\mathbf{H}^T \mathbf{H} + b)$ 。实际上,在一定程度上,这个 tanh 函数可以看作是 $\mathbf{H}^T \mathbf{H}$ 的另一种表示方式。为了解决式(22)中的问题,引入辅助变量 $\mathbf{Z} \in \mathbf{R}^{n \times n}$ 使上述目标函数可

分离。因此,式(22)等价于:

$$\min_{\mathbf{S}} \text{Tr}(\hat{\mathbf{K}} - 2\hat{\mathbf{K}}\mathbf{S} + \mathbf{S}^T \hat{\mathbf{K}}\mathbf{S}) + \frac{\beta}{\alpha} \|\mathbf{Z}\|_F^2 + \frac{\gamma}{\alpha} \text{Tr}(\mathbf{P}^T \mathbf{L}_s \mathbf{P}) \quad (23)$$

$$\text{s. t. } \mathbf{S} \geq \mathbf{0}, \text{diag}(\mathbf{S}) = \mathbf{0}, \mathbf{Z} = \mathbf{S}$$

利用增广拉格朗日乘法(ALM)来解决上述问题,相应的增广拉格朗日函数为:

$$L(\mathbf{Z}, \mathbf{S}, \mathbf{Y}) = \text{Tr}(\hat{\mathbf{K}} - 2\hat{\mathbf{K}}\mathbf{S} + \mathbf{S}^T \hat{\mathbf{K}}\mathbf{S}) + \frac{\beta}{\alpha} \|\mathbf{Z}\|_F^2 + \frac{\gamma}{\alpha} \text{Tr}(\mathbf{P}^T \mathbf{L}_s \mathbf{P}) + \frac{\mu}{2} \left\| \mathbf{Z} - \mathbf{S} + \frac{\mathbf{Y}}{\mu} \right\| \quad (24)$$

式中: μ 为惩罚参数; $\mathbf{Y} \in \mathbf{R}^{n \times n}$ 为拉格朗日乘数。通过固定其他变量,上述问题可以交替最小化 \mathbf{Z} 、 \mathbf{S} 和 \mathbf{Y} 。

对于变量 \mathbf{Z} , 通过令 $\mathbf{E} = \mathbf{S} - \frac{\mathbf{Y}}{\mu}$, 可以更新式(25)

中的元素。

$$\mathbf{Z}_{ij} = \max\left(|\mathbf{E}_{ij}| - \frac{\beta}{\alpha\mu}, 0\right) \cdot \text{sign}(\mathbf{E}_{ij}) \quad (25)$$

对于变量 \mathbf{S} , 通过令 $\mathbf{F} = \mathbf{Z} - \frac{\mathbf{Y}}{\mu}$, 可以更新式(26)

中的元素。

$$\min_{\mathbf{S}_i} \mathbf{S}_i^T \left(\frac{\mu}{2} \mathbf{I} + \hat{\mathbf{K}} \right) \mathbf{S}_i + \left(\frac{\gamma}{2\alpha} \mathbf{b}_i^T - \mu \mathbf{F}_i^T - 2\hat{\mathbf{K}}_{i,:} \right) \mathbf{S}_i \quad (26)$$

可以通过令式(26)中的 \mathbf{S}_i 为 0, 得到 \mathbf{S} 。对于变量 \mathbf{Y} , 相当于更新式(27)。

$$\mathbf{Y} = \mathbf{Y} + \mu(\mathbf{Z} - \mathbf{S}) \quad (27)$$

算法 2 总结了 R-MCLES 方法中全局相似矩阵的优化过程。除了全局相似矩阵 \mathbf{S} 的优化部分外, R-MCLES 的整体优化算法几乎与 MCLES 算法相同。

算法 2 采用 ALM 算法更新整体结构

输入: 潜在嵌入表示 \mathbf{H} , 聚类指标矩阵 \mathbf{P} , 参数 α 、 β 和 γ 。

输出: \mathbf{S} 。

初始化: $\mathbf{S} = \mathbf{Z} = \mathbf{0}$, $\mathbf{Y} = \mathbf{0}$, $\mu = 1$, 用 $\mathbf{H}^T \mathbf{H}$ 的 tanh 函数初始化 \mathbf{K} 。

1. 根据式(25)更新 \mathbf{Z} 。

2. 对于每一个 i , 通过解决式(26)更新 \mathbf{S} 的第 i 列。

3. 根据式(27)更新 \mathbf{Y} 。

1.4 算法时间复杂度与收敛性

所提出的 MCLES 和 R-MCLES 方法的优化分别由四个子问题组成。对于 MCLES 方法, 在交替优化方案下,

利用 ADMM 算法更新 \mathbf{W} , 其复杂度为 $O\left(\left(\sum_{v=1}^V d^v\right)^2 d\right)$ 。

用 Bartels-Stewart 算法求 \mathbf{H} 的计算复杂度为 $O(d^3)$ 。为了更新 \mathbf{S} , 二次规划的计算复杂度为 $O(n^3)$ 。 \mathbf{P} 的复杂度是 $O(cn^2)$ 。因此, 对于每次迭代, MCLES 方法

的总体计算复杂度为 $O\left(\left(\sum_{v=1}^V d^v\right)^2 d + n^3 + d^3 + cn^2\right)$ 。

此外, 对于 R-MCLES 方法, 只需进一步讨论 \mathbf{S} 的计算即可, 其复杂度为 $O(n^2)$, 远小于 MCLES 方法中 \mathbf{S} 的

计算复杂度。因此,对于每次迭代,R-MCLES 方法的总体计算复杂度为 $O((\sum_{v=1}^V d^v)^2 d + d^3 + (c+1)n^2)$ 。

2 实验与结果分析

2.1 实验设置

2.1.1 测试数据集

六个测试数据集分别是 Yale、MSRCv1、ORL、BBC-Sport、WebKB Texas 和 3Sources。

(1) Yale:是一个广泛使用的面部图像数据集,由 165 幅灰度图像组成,共有 15 位人员,每位人有 11 幅图像。本实验使用了三类图像,它们的尺寸分别是 4 096、3 304 和 6 750。

(2) MSRCv1:是一个图像数据集,共有七个类别,210 个对象。包括树、建筑物、飞机、汽车和牛等七类。本实验采用 MSRCv1 数据集的四个类数据,分别是 CM 特征(视图 1)、GIST 特征(视图 2)、LBP 特征(视图 3)和 GENT 特征(视图 4)。

(3) ORL:是一个广泛使用的人脸图像数据集,共有 40 位志愿者,每位有 10 幅图片,共 400 幅图片。本实验使用了三种特征,即强度特征(视图 1)、LBP 特征(视图 2)和 Gabor 特征(视图 3)。

(4) BBCSport^[12]:是一个由英国广播公司体育新闻体育网站 2004 年—2005 年 5 个主题的 544 个文档组成的数据集,每个主题新闻对应一个类别。本实验采用 BBCSport 数据集的两个视图,其维度分别为 3 183 和 3 203。

(5) WebKB-Texas:这个数据集由大学计算机系的网页组成。共有 187 个网页,每个网页有 5 个指标。在本实验中,WebKB Texas 数据集包含两个视图,分别是内容特性(视图 1)和城市特性(视图 2)。

(6) 3Sources:是一个多源的新闻数据集,包含三个来源,即英国广播公司、卫报和路透社。该数据集有 169 个新闻对象、6 个类、3 种媒体报道。本实验使用了三个视图,分别是 BBC(视图 1)、卫报(视图 2)和路透社(视图 3)。

2.1.2 参与比较的方法

将本文方法与两种经典的单视图聚类方法和六种最新的多视图聚类方法进行了比较。

(1) 谱聚类(SC)^[6]:对每个视图分别进行单视图 SC 方法。

(2) ConPCA^[7]:是一种扩展的 SC 方法。首先将所有视图的特征结合起来,然后应用 PCA(主成分分

析)方法提取低维表示,最后将低维表示输入到谱聚类算法中,得到最终的聚类结果。

(3) Co-Reg^[10]:对聚类假设进行了共调控,从而可以在视图之间形成相同的聚类关系。

(4) Co-Tr^[11]:不管视图的约束,假设一个点被真正的底层聚类分配给同一个簇。

(5) Min-Dis^[12]:在谱聚类算法的基础上,采用“最小分歧”策略构造二部图。

(6) RMSC^[13]:是一种利用标准马尔可夫链进行聚类的多视图谱聚类方法,鲁棒性较好。

(7) LMSC^[14]:根据多个视图的共同潜在结构发现一个子空间表示,然后将其输入到谱聚类算法中。

(8) 多视图聚类的图学习(MVGL)^[15]:根据每个视图的优化图得到全局图。

对于上述八种方法,按原文献建议进行参数调整,从而生成最佳结果。在实验中,由于采用的随机初始化方法,因此每运行 20 次,报告结果的平均准确度和标准差。

2.1.3 评价指标

本文采用了四种常用的评价指标对模型进行评估,即准确度(ACC)、归一化互信息(NMI)、纯度(PUR)和 Rand 指数(RI)^[12]。每个指标的值越高表示性能越好。

2.2 参数分析

本节对所提出的 MCLES 和 R-MCLES 方法的四个参数 d 、 α 、 β 和 γ 进行分析,在 $[10, 100]$ 、 $[0.8, 8]$ 、 $[0.1, 10]$ 和 $[0.001, 0.01]$ 范围内对这四个参数进行调整。此调整策略也用于后面的收敛性分析、比较实验和可视化。

R-MCLES 算法基于六个测试数据集,有关维度 d 的 ACC 和 NMI 的结果如图 2 和图 3 所示。可以看出,R-MCLES 的性能总体上是比较稳定的,即在大范围内,ACC 和 NMI 的变化相对较小。但结果在不同的数据集上有显著差异,主要是因为不同原始多视图数据集的视图维度不同。

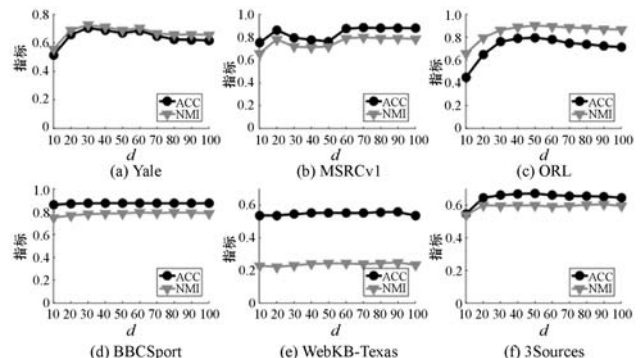


图2 MCLES 参数 d 分析

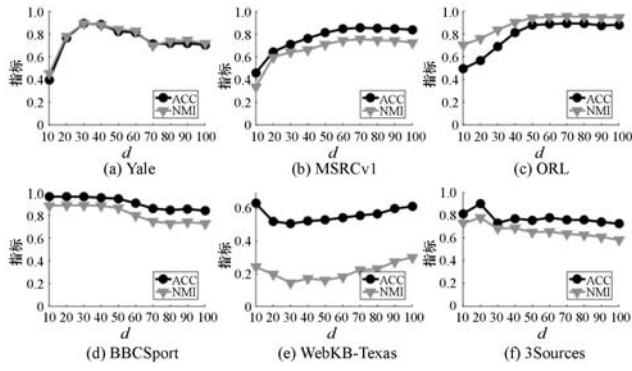


图3 R-MCLES 参数 d 分析

图4 - 图9 分别给出了三个权衡参数 α , β 和 γ 的分析结果。可以清楚地看出,三个权衡参数的取值范围都很大。然而,与维数 d 的情况类似,由于不同数据集具有不同的性质,实现相对最佳和稳定结果的最佳范围因数据集而异。因此, α , β 和 γ 的值因数据集而异。

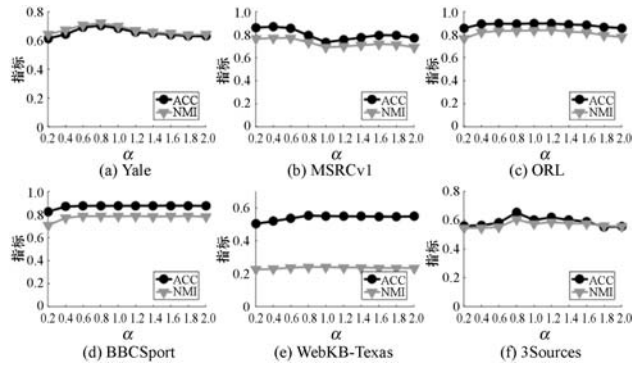


图4 MCLES 参数 α 分析

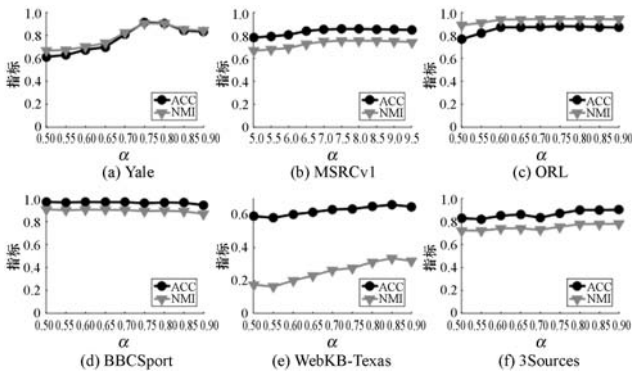


图5 R-MCLES 参数 α 分析

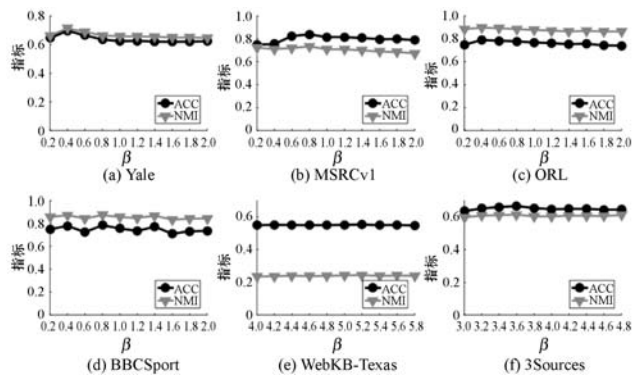


图6 MCLES 参数 β 分析

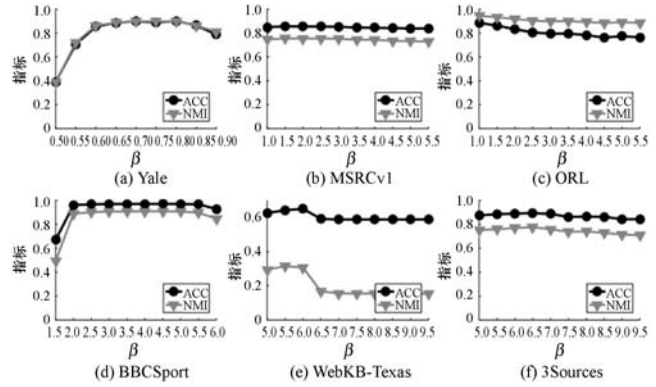


图7 R-MCLES 参数 β 分析

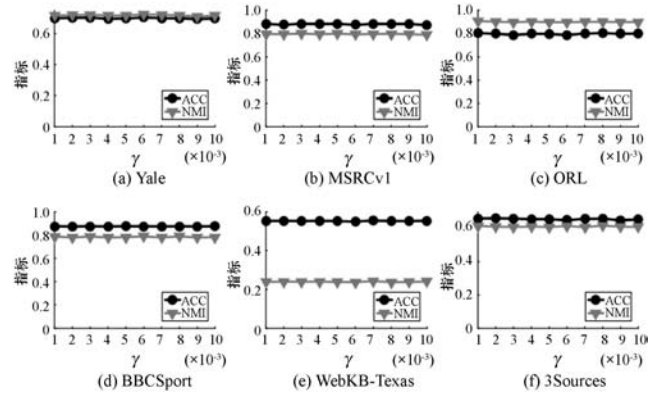


图8 MCLES 参数 γ 分析

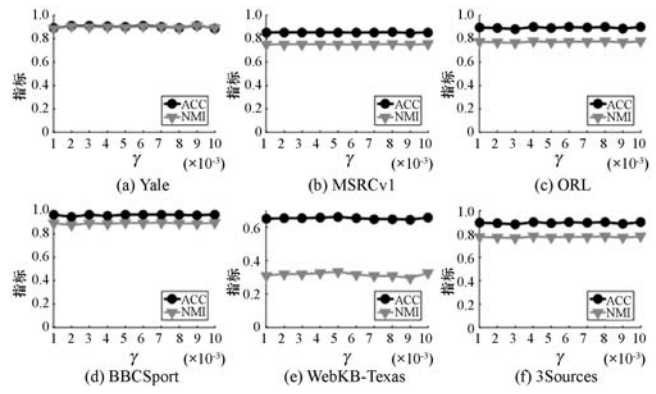
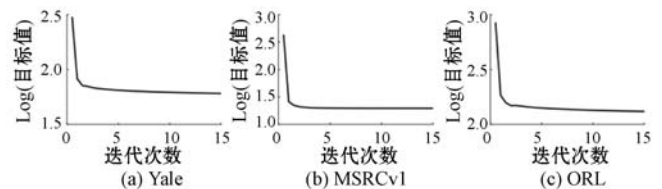


图9 R-MCLES 参数 γ 分析

2.3 收敛性分析

本节对本文模型进行收敛性分析。与参数分析使用相同的六个数据集,即 Yale、MSRCv1、ORL、BBCSport、WebKB Texas 和 3Sources。图 10 和图 11 分别绘制了 MCLES 和 R-MCLES 每一步迭代的结果。从图 10 中可以发现,在六个数据集的迭代过程中,目标值迅速减少,经过 30 步的迭代,可以收敛。

同样,虽然 R-MCLES 方法前几次迭代会有一些波动,但也会在大约 30 次迭代后收敛。



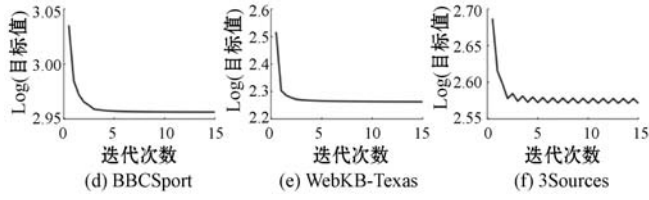


图 10 MCLES 收敛性分析

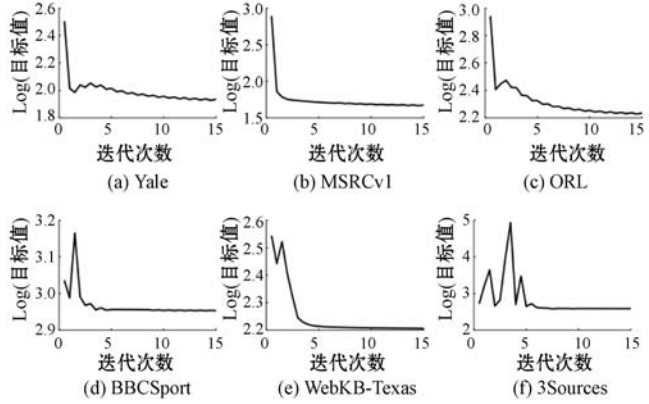


图 11 R-MCLES 收敛性分析

2.4 比较结果

本文比较了 10 种聚类方法基于 6 个基准数据集的 ACC、NMI、PUR 和 RI 的比较结果,并将结果分别统计在表 1 – 表 6 中,依次为图像数据集和文档数据集。表中均包含均值和标准偏差。此外,对于单视图光谱聚类(SC)方法,第一视图的结果表示为 SC1,第二视图的结果表示为 SC2。

表 1 Yale 数据集比较结果

方法	ACC	NMI	Purity	RI
SC 1	0.567(0.036)	0.602(0.023)	0.575(0.034)	0.924(0.004)
SC 2	0.573(0.035)	0.605(0.022)	0.581(0.033)	0.927(0.004)
SC 3	0.636(0.040)	0.657(0.030)	0.645(0.040)	0.935(0.005)
ConPCA	0.564(0.037)	0.609(0.028)	0.575(0.035)	0.925(0.004)
Co-Reg	0.604(0.005)	0.642(0.004)	0.614(0.005)	0.931(0.000)
Co-Tr	0.628(0.007)	0.659(0.006)	0.634(0.008)	0.934(0.001)
Min-Dis	0.597(0.010)	0.630(0.007)	0.603(0.009)	0.928(0.001)
RMSC	0.534(0.015)	0.586(0.010)	0.545(0.013)	0.920(0.002)
LMSC	0.673(0.010)	0.695(0.012)	0.677(0.010)	0.935(0.003)
MVGL	0.630(0.000)	0.638(0.000)	0.642(0.000)	0.924(0.000)
MCLES	0.708 (0.019)	0.729 (0.019)	0.709 (0.018)	0.940 (0.006)
R-MCLES	0.923 (0.014)	0.910 (0.016)	0.923 (0.014)	0.978 (0.005)

表 2 MSRCv1 数据集比较结果

方法	ACC	NMI	Purity	RI
SC 1	0.418(0.011)	0.332(0.020)	0.467(0.015)	0.813(0.004)
SC 2	0.689(0.054)	0.598(0.048)	0.726(0.040)	0.884(0.014)

续表 2

方法	ACC	NMI	Purity	RI
SC 3	0.615(0.006)	0.506(0.013)	0.646(0.007)	0.853(0.002)
SC 4	0.702(0.014)	0.547(0.011)	0.702(0.014)	0.872(0.002)
ConPCA	0.615(0.003)	0.506(0.009)	0.647(0.008)	0.853(0.002)
Co-Reg	0.600(0.005)	0.519(0.003)	0.628(0.003)	0.851(0.001)
Co-Tr	0.782(0.007)	0.700(0.008)	0.794(0.007)	0.911(0.002)
Min-Dis	0.590(0.008)	0.514(0.006)	0.605(0.007)	0.852(0.002)
RMSC	0.389(0.011)	0.282(0.009)	0.414(0.011)	0.794(0.002)
LMSC	0.655(0.100)	0.558(0.098)	0.671(0.092)	0.864(0.024)
MVGL	0.671(0.000)	0.577(0.000)	0.704(0.000)	0.863(0.000)
MCLES	0.882(0.004)	0.795(0.009)	0.882(0.004)	0.941(0.002)
R-MCLES	0.856(0.003)	0.755(0.005)	0.856(0.003)	0.929(0.001)

表 3 ORL 数据集比较结果

方法	ACC	NMI	Purity	RI
SC 1	0.651(0.031)	0.800(0.015)	0.685(0.025)	0.977(0.001)
SC 2	0.771(0.025)	0.888(0.012)	0.802(0.020)	0.985(0.001)
SC 3	0.708(0.027)	0.842(0.012)	0.738(0.019)	0.981(0.001)
ConPCA	0.653(0.025)	0.803(0.014)	0.690(0.023)	0.977(0.000)
Co-Reg	0.692(0.004)	0.837(0.001)	0.729(0.003)	0.980(0.000)
Co-Tr	0.751(0.005)	0.879(0.003)	0.785(0.005)	0.984(0.000)
Min-Dis	0.721(0.008)	0.858(0.003)	0.758(0.006)	0.982(0.000)
RMSC	0.748(0.011)	0.879(0.005)	0.785(0.009)	0.984(0.000)
LMSC	0.818 (0.034)	0.916 (0.020)	0.854 (0.029)	0.988 (0.002)
MVGL	0.735(0.000)	0.865(0.000)	0.795(0.000)	0.971(0.000)
MCLES	0.792(0.021)	0.897(0.011)	0.835(0.016)	0.985(0.002)
R-MCLES	0.907 (0.021)	0.963 (0.007)	0.925 (0.015)	0.995 (0.000)

表 4 BBCSport 数据集比较结果

方法	ACC	NMI	Purity	RI
SC 1	0.845(0.000)	0.672(0.000)	0.845(0.000)	0.888(0.000)
SC 2	0.511(0.000)	0.234(0.000)	0.571(0.000)	0.623(0.000)
ConPCA	—	—	—	—
Co-Reg	0.694(0.005)	0.536(0.002)	0.735(0.004)	0.788(0.001)
Co-Tr	0.875(0.003)	0.718(0.002)	0.877(0.002)	0.912(0.001)
Min-Dis	0.855(0.006)	0.785(0.006)	0.873(0.004)	0.926(0.003)
RMSC	0.848(0.022)	0.777(0.010)	0.868(0.012)	0.924(0.008)
LMSC	0.919 (0.002)	0.836 (0.004)	0.919 (0.002)	0.952 (0.001)
MVGL	0.419(0.000)	0.08(0.000)	0.422(0.000)	0.333(0.000)
MCLES	0.880(0.003)	0.807(0.011)	0.880(0.003)	0.942(0.004)
R-MCLES	0.972 (0.003)	0.910 (0.007)	0.972 (0.003)	0.972 (0.003)

表 5 WebKB-Texas 数据集比较结果

方法	ACC	NMI	Purity	RI
SC 1	0.550(0.000)	0.021(0.000)	0.561(0.000)	0.381(0.000)
SC 2	0.533(0.015)	0.272(0.011)	0.692(0.024)	0.659(0.009)
ConPCA	0.552(0.018)	0.285 (0.013)	0.709 (0.020)	0.676(0.007)
Co-Reg	0.589(0.002)	0.199(0.004)	0.635(0.005)	0.642(0.003)
Co-Tr	0.600 (0.006)	0.248(0.005)	0.677(0.005)	0.689 (0.003)
Min-Dis	0.531(0.007)	0.197(0.003)	0.653(0.003)	0.657(0.003)
RMSC	0.497(0.006)	0.235(0.003)	0.667(0.005)	0.653(0.001)
LMSC	0.484(0.045)	0.206(0.027)	0.650(0.021)	0.626(0.018)
MVGL	0.598(0.000)	0.280(0.000)	0.689(0.000)	0.562(0.000)
MCLES	0.550(0.012)	0.242(0.009)	0.683(0.008)	0.673(0.006)
R-MCLES	0.662 (0.018)	0.336 (0.044)	0.716 (0.028)	0.696 (0.022)

表 6 3Sources 数据集比较结果

方法	ACC	NMI	Purity	RI
SC 1	0.670(0.051)	0.559(0.026)	0.740(0.015)	0.829(0.013)
SC 2	0.697(0.021)	0.608(0.015)	0.797(0.007)	0.849(0.004)
SC 3	0.621(0.026)	0.604(0.018)	0.759(0.006)	0.831(0.007)
ConPCA	0.763 (0.013)	0.690 (0.022)	0.825 (0.018)	0.875 (0.005)
Co-Reg	0.551(0.003)	0.486(0.002)	0.677(0.002)	0.767(0.001)
Co-Tr	0.598(0.005)	0.553(0.005)	0.741(0.005)	0.807(0.002)
Min-Dis	0.527(0.007)	0.483(0.006)	0.688(0.005)	0.764(0.004)
RMSC	0.542(0.014)	0.505(0.012)	0.687(0.009)	0.762(0.006)
LMSC	0.710(0.010)	0.649(0.015)	0.810(0.011)	0.853(0.007)
MVGL	0.349(0.000)	0.121(0.000)	0.414(0.000)	0.360(0.000)
MCLES	0.660(0.021)	0.619(0.020)	0.770(0.015)	0.835(0.012)
R-MCLES	0.915 (0.020)	0.802 (0.025)	0.915 (0.020)	0.936 (0.014)

根据这六个表格,可以得出以下结论。

(1) 总体上,与其他方法相比,MCLES 方法的稳定性和聚类性都较出色,几乎排名前三位。例如,对于 Yale 和 MSRCv1 数据集,MCLES 方法比现有的所有方法的性能都好。然而,在某些文档如 3sources 数据集上,MCLES 的性能不如 R-MCLES,因为它不能在严格的约束下有效地考虑更多数据对的底层语义信息。

(2) R-MCLES 方法通过放宽全局相似性矩阵和潜在表示矩阵的内积相似性约束,进一步充分利用了数据的潜在语义信息,在 90% 的测试数据集上的测试结果能够取得较好的聚类精度。R-MCLES 的 NMI 在 Yale 和 ORL 数据集上的运行结果比 MCLES 分别高 18.18% 和 6.54%。虽然 R-MCLES 方法在 MRSCv1 数据集上的性能并不优于 MCLES,但其性能也接近

MCLES,并且优于其他传统方法。

(3) 本文算法在不同的数据集上具有很强的鲁棒性,尽管性能并不总是最好的。但是有些算法的性能不是很稳定,即在某些数据集上可以获得很好的结果,但在其他数据集上却没有很好的结果。注意,ConPCA 方法不能在 BBCSport 数据集上运行,因为 BBCSport 数据集上的特征太稀疏,无法运行 SVD。

(4) R-MCLES 方法在图像数据集和文档数据集上的性能都较好。例如,在 ORL 数据集(图像数据集)上,R-MCLES 的 ACC 和 NMI 分别比次优方法高出 8.88% 和 4.6%。对于 3Sources 数据集(文档数据集),R-MCLES 的 ACC 和 NMI 分别比次优方法高出 20.38% 和 15.22%。因此,改善后的 R-MCLES 方法在不同类型的数据集上表现得更好。

综上所述,本文方法通过直接学习基于潜在嵌入表示矩阵的相似度矩阵和聚类指标矩阵,获得了鲁棒性更强、准确度更高的聚类结果。

2.5 可视化

本节将模型进行可视化,t 分布随机邻域嵌入(t-SNE)用于将每个视图的原始特征以及 MCLES 和 R-MCLES 获得的潜在嵌入表示映射到二维空间中,然后将二维空间中的数据点进行可视化。由于篇幅限制,只使用 Yale 和 MSRCv1 数据集,它们的结果分别见图 12 和图 13,用不同的灰度表示不同的类别。可以发现,本文方法可以清楚地表现出潜在的簇结构。也就是说,与每个视图的原始特征相比,潜在嵌入表示包含更好的簇结构。因此,R-MCLES 比 MCLES 能够学习到更好的潜在嵌入表示,这进一步证实了降低全局相似矩阵约束能够完善模型的理论基础。

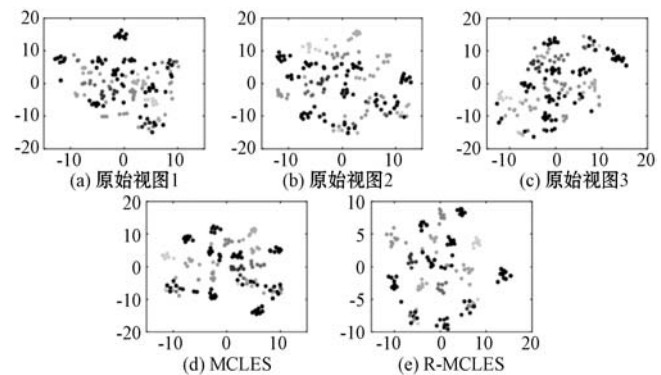
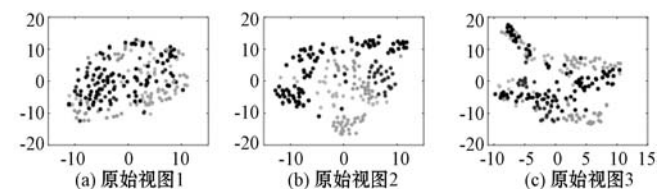


图 12 Yale 数据集可视化



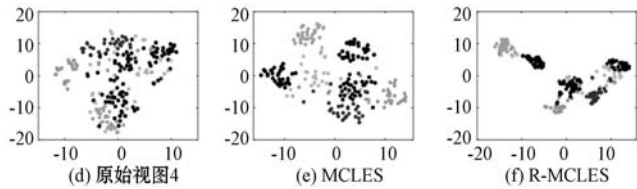


图13 MSRCv1数据集可视化

2.6 计算时间对比

通过对比各个方法在三个数据集上聚类时间来分析各个方法的计算时间代价,结果如表7所示。

表7 计算时间对比 单位:s

数据集	Yale	MSRCv1	ORL
SC 1	389.8	245.6	215.6
SC 2	456.9	311.5	255.4
SC 3	388.7	265.3	221.6
ConPCA	189.3	124.6	93.6
Co-Reg	284.6	189.3	168.5
Co-Tr	332.1	205.6	166.6
Min-Dis	597.6	363.3	289.6
RMSC	241.2	177.6	158.6
LMSC	298.6	196.3	155.9
MVGL	547.3	339.6	277.9
MCLES	302.2	198.9	144.9
R-MCLES	212.3	134.6	113.2

可以看出,基于改进松弛嵌入空间的多视图聚类相对域 MCLES 极大地降低了计算时间,有效缩短了计算时间。

3 结语

针对传统方法缺乏统一特征表示,存在保守性的缺陷,提出一种基于改进松弛嵌入空间的多视图聚类方法。最后分析多个数据集的实验结果可以得出如下结论:

(1) 基于改进松弛嵌入空间的多视图聚类能够充分利用视图之间的语义信息,进一步挖掘各个多视图数据集之间的隐藏信息,因此能够得到精度更高的聚类效果。

(2) 基于改进松弛嵌入空间的多视图聚类过放宽全局相似性矩阵和潜在表示矩阵的内积相似性约束,有效地降低了计算的复杂度,节省了计算成本,降低了计算时间。

(3) 本文算法在不同的数据集上均能够取得较好

的聚类结果,不论是图像数据集还是文档数据集,证明了本文方法具有很强的泛化能力,另外对于参数变化具有一定的鲁棒性。

参考文献

- [1] 洪敏,贾彩燕,王晓阳. K-means 型多视图聚类中的初始化问题研究[J]. 计算机科学与探索,2019,13(4):574-585.
- [2] Li R H, Zhang C Q, Hu Q H, et al. Flexible multi-view representation learning for subspace clustering[C]//28th International Joint Conference on Artificial Intelligence,2019:2916-2922.
- [3] Chen M S, Huang L, Wang C D, et al. Multi-view clustering in latent embedding space[J]. Proceedings of the AAAI Conference on Artificial Intelligence,2020,34(4):1125-1135.
- [4] 洪敏,贾彩燕,李亚芳,等. 样本加权的多视图聚类算法[J]. 计算机研究与发展,2019,56(8):1677-1685.
- [5] 刘良凤,刘三阳. 基于权重差异度的动态模糊聚类算法[J]. 吉林大学学报(理学版),2019,57(3):574-582.
- [6] 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(8):1627-1638.
- [7] Chen X J, Sun W Y, Wang B, et al. Spectral clustering of customer transaction data with a two-level subspace weighting method[J]. IEEE Transactions on Cybernetics,2019,49(9):3230-3241.
- [8] 李杏峰,黄玉清,任珍文. 联合低秩稀疏的多核子空间聚类算法[J]. 计算机应用,2020,40(6):1648-1653.
- [9] 夏菁,丁世飞. 基于低秩稀疏约束的自权重多视角子空间聚类[J]. 南京大学学报(自然科学),2020,56(6):862-869.
- [10] Zhan K, Zhang C Q, Guan J P, et al. Graph learning for multiple view clustering[J]. IEEE Transactions on Cybernetics,2018,48(10):2887-2895.
- [11] 李占芳,李慧云,刘新为. 分类稀疏低秩表示的子空间聚类方法[J]. 系统科学与数学,2018,38(8):852-865.
- [12] 解昊,赵志刚,吕慧显,等. 非负局部约束低秩子空间聚类算法[J]. 计算机工程与应用,2018,54(23):137-143,155.
- [13] 张越美. 基于子空间学习的多视图聚类方法研究[D]. 西安:西安电子科技大学,2019.
- [14] 王丽娟,丁世飞,丁玲. 基于迁移学习的软子空间聚类算法[J]. 南京大学学报(自然科学),2020,56(4):515-523.
- [15] 徐鲲鹏,陈黎飞,孙浩军,等. 类属型数据核子空间聚类算法[J]. 软件学报,2020,31(11):3492-3505.