

多示例嵌入学习的实例关联性挖掘与强化

杨梅^{1,3,4}, 邓雯¹, 张本文², 闵帆^{1,3,4*}

(1.西南石油大学计算机科学学院, 四川 成都 610500; 2.四川民族学院理工学院, 四川 康定 626001; 3.西南石油大学人工智能研究院, 四川 成都 610500; 4.西南石油大学机器学习研究中心, 四川 成都 610500)

摘要:提出了多示例嵌入学习(multi-instance learning, MIL)的实例关联性挖掘与强化算法(multi-instance embedding learning with instance affinity mining and reinforcement, MEMR),包括3个技术。关联性挖掘技术基于自定义的关联性指标,首先在负实例空间中选择初始负代表实例集,然后根据正、负实例间的差异性,选择初始正代表实例集。关联性强化技术分别评估初始正、负代表实例集与整个实例空间的正负关联性,获得整体关联性更强的代表实例集。包嵌入技术通过嵌入函数将包转换为单向量进行学习。实验在4类应用领域和7种对比算法上进行。结果表明,MEMR的准确性总体优于其他对比算法,特别是在图像检索和网页推荐数据集上具有显著优势。

关键词:关联性挖掘;关联性强化;嵌入方法;实例选择;多示例学习

中图分类号: TP181 **文献标志码:** A

引用格式: 杨梅,邓雯,张本文,等.多示例嵌入学习的实例关联性挖掘与强化[J].山东大学学报(理学版),2024,59(1):35-45.

Multi-instance embedding learning with instance affinity mining and reinforcement

YANG Mei^{1,3,4}, DENG Wen¹, ZHANG Benwen², MIN Fan^{1,3,4*}

(1. School of Computer Science, Southwest Petroleum University, Chengdu 610500, Sichuan, China; 2. School of Polytechnic, Sichuan Minzu College, Kangding 626001, Sichuan, China; 3. Institute for Artificial Intelligence, Southwest Petroleum University, Chengdu 610500, Sichuan, China; 4. Lab of Machin Learning, Southwest Petroleum University, Chengdu 610500, Sichuan, China)

Abstract: We propose the multi-instance embedding learning with instance affinity mining and reinforcement (MEMR) algorithm, including three techniques. The affinity mining technique is based on a custom affinity metric. First, the initial negative representative instance set (INRI) is selected in the negative instance space. Then, the initial positive representative instance set (IPRI) is chosen according to the difference between positive and negative instances. The affinity reinforcement technique evaluates the positive (negative) affinity between IPRI (INRI) and the entire instance space to obtain a representative instance set with stronger overall affinity. The bag embedding technique converts bags into single vectors for learning through the designed embedding function. Experiments are carried out across four application domains and seven comparison algorithms. The results show that MEMR generally outperforms other comparison algorithms in accuracy, especially in image retrieval and web recommendation datasets.

Key words: affinity mining; affinity reinforcement; embedding method; instance selection; multi-instance learning

0 引言

MIL 的处理对象称为包,每个包由多个实例组成。MIL 最早用于药物活性检测^[1],即将每个分子视作

收稿日期:2022-08-02; 网络出版时间:2023-04-24 14:23:15

网络出版地址: <http://kns.cnki.net/kcms/detail/37.1389.N.20230423.1348.004.html>

基金项目:国家自然科学基金资助项目(62006200);四川省自然科学基金资助项目(2019YJ0314);中央引导地方科技发展专项项目(2021ZYD0003);浙江省海洋大数据挖掘与应用重点实验室开放课题(OBDMA202102)

第一作者简介:杨梅(1982—),女,副教授,硕士,研究方向为多示例学习、深度学习。E-mail: yangmei@swpu.edu.cn

* 通信作者简介:闵帆(1973—),男,教授,博士,研究方向为粒计算、主动学习。E-mail: minfan@swpu.edu.cn

包,同一分子的不同形状视作实例。如果一个包至少包含一个可以用于制造药物的实例,则被标记为正,否则标记为负,以此来预测新分子是否可以用于制造药物。由于这种独特的性质与实际应用高度契合,因此 MIL 被广泛用于机器学习和计算机视觉等领域,例如图像分类^[2-4]、文本分类^[5-7]、医学诊断^[8-9]及生物方程标注^[10]。

根据 MIL 算法的实现原理,可以将现有算法大致分为两类^[11]:(1) 传统方法立足原始数据空间来预测包标签。如 MI-Boosting^[12]假设包中每个实例对包标签的贡献是相同且独立的,从而将单实例预测与包标签概率估计相联系;MIRSVM^[13]只关注包之间的关系,基于包代表训练 SVM,找到包代表之间的最佳分隔超平面;MIORDM^[14]从实例空间中选代表实例并研究它们的边际分布信息,以构建最佳分隔超平面;(2) 基于嵌入的方法通过将包转换为新特征空间中的单向量,从而将 MIL 问题简化为单实例学习(single-instance learning, SIL)问题。如 MIKI^[11]分析正实例的分布变化,基于加权多类模型选择高积极性的实例原型,从而将包转换为包含实例原型信息的向量;JointMIL^[15]考虑包的组成以及实例在分类中的相关性,提出一种中心损失,以共同学习实例级和包级嵌入;ELDB^[16]基于数据的空间分布与强化技术选择具有判别性的关键包集合,并通过集成策略得到加权模型,最终获得具有可区分性的单向量。流行的 MIL 嵌入方法通常忽略了正、负实例间的关联性,并缺乏增强能力,可能导致所选实例的代表性较弱,影响模型的最终性能。

本文提出多示例嵌入学习的实例关联性挖掘与强化算法(multi-instance embedding learning with instance affinity mining and reinforcement, MEMR),如图 1 所示。其中,关联性挖掘技术基于设计的关联性指标,在考虑正、负实例差异性的基础上,生成具有强代表性的初始正代表实例集和负代表实例集。关联性强化技术分别计算初始正、负代表实例集与整个实例空间的正、负关联性,并与初始正、负代表实例集内部的关联性进行比较,以获得整体关联性更强的代表实例集。最后,包嵌入技术基于包中实例与代表实例的相似关系,将包转换到新特征空间。在 34 个数据集上的实验验证了 MEMR 的性能。这些数据集源自不同的应用领域,包括图像检索、医学图像、文本分类及网页推荐。实验结果表明,MEMR 整体上优于其他对比算法,尤其在图像检索和网页推荐数据集上表现突出。

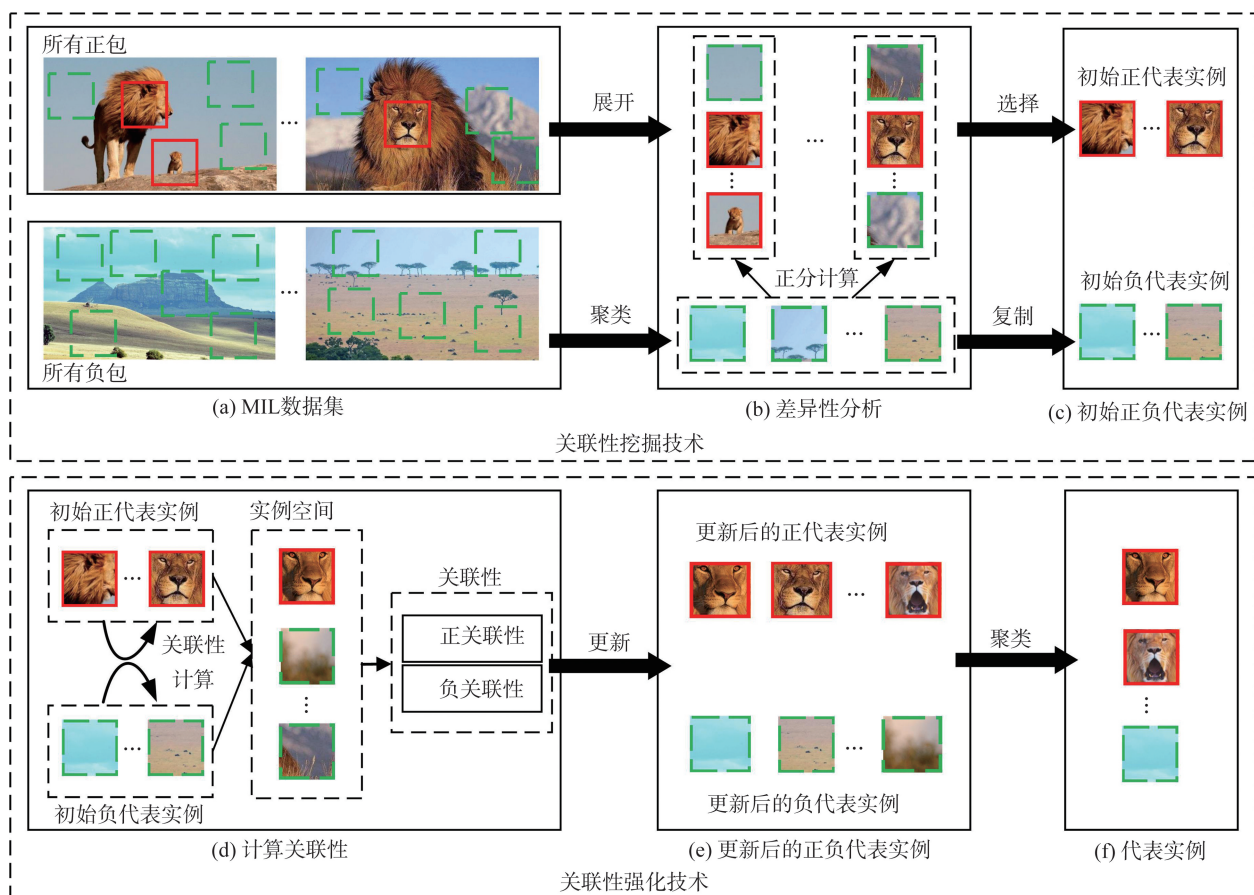


图 1 MEMR 算法的主体流程

Fig.1 Main flow of the MEMR algorithm

本文的主要贡献是:(1)提出了一个代表实例的关联性挖掘技术。基于关联性指标,挖掘正、负实例之间的关联性,提升所选正、负代表实例集的代表性。(2)提出了一个具有增强能力的代表实例关联性强化技术。利用整个实例空间的特征信息,评估并比较初始正、负代表实例集和实例空间的关联性,以强化代表实例集的整体关联性。

1 相关工作

MIL 嵌入方法的思想是将包嵌入为单向量,以利用 SIL 方法解决 MIL 问题。流行的 MIL 嵌入方法可以分为 4 类:(1) 基于统计的方法通过计算数据的统计量来表示一个包。如 Simple-MI^[17] 计算包中每个属性的平均值,并将其作为新向量的属性值;(2) 基于核的方法侧重于设计用于嵌入的内核。如 miFV^[18] 使用高斯混合模型对实例特征空间的密度建模,并基于 Fisher 核获得包的单向量表示;(3) 基于包的方法的主要任务是寻找关键包。如 Bamic^[19] 将包级 k -Medoids 聚类的簇心作为关键样本,并用 Hausdorff 距离度量包与每个簇心的距离,从而得到嵌入向量;(4) 基于实例的方法聚焦于关键实例的选择,然后根据包与实例的相似性,将每个包嵌入到由关键实例定义的特征空间。如 MILES^[20] 根据包和所有实例之间的距离度量,将包嵌入到新特征空间。

基于实例的 MIL 嵌入方法的关键实例选择结果会很大程度地影响模型的性能。MILD^[21] 提出基于实例标签的消歧方法,识别正包中的真正实例;MILFM^[22] 进行部分实例修剪,将正包中的所有实例和负包中的聚类中心视为实例原型;miVLAD^[23] 基于整个实例空间聚类,将簇心作为关键实例,并根据包中实例与簇心的差异实现嵌入;MILDM^[24] 构建判别性关键实例池,使包嵌入结果可以较容易地相互分离;StableMIL^[25] 确定多实例学习和因果推理之间的内在联系,选择可以改变负包标签的因果实例作为代表。

在以上这些算法中,代表实例的质量大多依赖于实例选择并且缺少正、负实例关联性的研究。本文基于关联性指标研究正、负实例关联性以进行实例选择,首先获得初始代表实例,然后增强代表实例的代表性,从而提升模型性能。

2 实例关联性挖掘与强化算法

首先介绍符号系统并定义关联性指标,然后详细介绍实例选择的 2 个关键技术。其中,关联性挖掘技术用于获得初始正代表实例集和负代表实例集;关联性强化技术用于评估整个实例空间的关联性,并强化代表实例的整体关联性。最后介绍包嵌入技术。

2.1 符号系统

令 $\mathcal{X} = \mathbf{R}^d$ 表示 d 维实例空间, $\mathcal{F} = \{\mathbf{B}_i\}_{i=1}^N$ 表示给定数据集, $\mathbf{Y} = [y_1, y_2, \dots, y_N]$ 表示标签向量,其中 $\mathbf{B}_i = \{\mathbf{x}_{ij}\}_{j=1}^{n_i}$ 表示第 i 个包, $y_i \in \{-1, +1\}$ 表示 \mathbf{B}_i 的标签, \mathbf{x}_{ij} 表示 \mathbf{B}_i 的第 j 个实例, N 表示 \mathcal{F} 中包的个数, n_i 则表示 \mathbf{B}_i 中实例的个数。表 1 列出了本文中一些重要符号。

表 1 符号说明
Table 1 Notation description

符号	含义	符号	含义
\mathcal{X}	实例空间	\mathbf{V}_i	\mathbf{B}_i 的嵌入向量
\mathcal{F}	数据集	N	\mathcal{F} 中包的个数
\mathbf{Y}	标签向量	n_i	\mathbf{B}_i 中实例的个数
\mathbf{B}_i	第 i 个包	d	实例的维度
\mathbf{x}_{ij}	\mathbf{B}_i 的第 j 个实例	m	正包的个数
y_i	\mathbf{B}_i 的标签	ψ	\mathbf{C} 中代表实例的个数
\mathbf{C}	代表实例集		

2.2 关联性指标

一个实例与一个集合的平均相似性越高,表示它们之间的关联性越强。本文将实例与集合之间的关联性指标定义如下:

$$A(\mathbf{x}^*, \mathbf{R}) = \frac{1}{h} \sum_{k=1}^h \exp(-\|\mathbf{x}^* - \mathbf{r}_k\|_2), \quad (1)$$

其中: \mathbf{x}^* 表示待评估的实例; $\mathbf{R} = \{\mathbf{r}_k\}_{k=1}^h$ 表示实例集合; h 为 \mathbf{R} 中实例的数量。 $A(\mathbf{x}^*, \mathbf{R})$ 越大, \mathbf{x}^* 与集合 \mathbf{R} 的关联性越强, 则 \mathbf{x}^* 包含的有用信息越多, 更适合被选做代表实例。指标中采用的欧几里得距离是常用于度量实例之间距离^[16,19,26]的方法, 且通过平均化, 在充分利用实例 \mathbf{x}^* 与集合内每个实例的关系的同时, 降低边缘点的影响, 从而提高应对噪声的鲁棒性。MEMR 算法将关联性指标与实例的代表性度量联系起来, 通过挖掘并增强代表实例的关联性, 选取信息丰富的实例作为代表实例, 以提升 MIL 学习模型的性能。

2.3 关联性挖掘技术

基于 MIL 标准假设^[1], 正实例的数量远小于负实例的数量, 而已有的实例选择方法通常忽略正、负实例之间的差异性^[20-23], 可能导致所选代表实例集中未包含正实例。关联性挖掘技术利用正、负实例之间的差异性, 分别选择正、负代表实例集, 是解决该问题的途径之一。

如算法 1 所示, 第 1—2 行分别获得正包的集合 \mathcal{F}^+ 和负实例空间 \mathcal{F}_n 。第 3 行利用负包中只包含负实例的数据特点, 使用 k -means 算法从 \mathcal{F}_n 中选择 m 个聚类中心作为初始负代表实例集 $\mathbf{R}^- = \{\mathbf{r}_k^-\}_{k=1}^m$ 。第 4—11 行利用正包中至少包含一个正实例的数据特点, 从每个正包 \mathbf{B}_i 中选择一个最正的实例, 从而获得初始正代表实例集 \mathbf{R}^+ 。具体做法是: 首先, 利用式(1)计算 $\mathbf{x}_{ij} \in \mathbf{B}_i$ 与 \mathbf{R}^- 的关联性 $A(\mathbf{x}_{ij}, \mathbf{R}^-)$ 。根据正、负实例间差异最大化原则, 即关联性越小, \mathbf{x}_{ij} 与 \mathbf{R}^- 之间差异性越大, \mathbf{x}_{ij} 更可能为正实例, 因此, \mathbf{x}_{ij} 为正的得分为:

$$s_{ij} = 1/A(\mathbf{x}_{ij}, \mathbf{R}^-). \quad (2)$$

然后, 比较包中每个实例的得分, 将得分最高的实例 \mathbf{x}_{ij^*} 记为该包中最正的实例, 且作为 \mathbf{B}_i 的正代表实例, 并将 \mathbf{x}_{ij^*} 并入 \mathbf{R}^+ , 其中:

$$j^* = \arg \max_j s_{ij}. \quad (3)$$

为了便于描述, 将最后得到的初始正代表实例集表示为 $\mathbf{R}^+ = \{\mathbf{r}_k^+\}_{k=1}^m$ 。

算法 1 insInit (\mathcal{F})

输入:

数据集 $\mathcal{F} = \{\mathbf{B}_i\}_{i=1}^N$;

标签向量 $\mathbf{Y} = [y_1, y_2, \dots, y_N]$;

输出:

初始正(负)代表实例集 $\mathbf{R}^+(\mathbf{R}^-)$;

1. $\mathcal{F}^+ = \{\mathbf{B}_i \in \mathcal{F} \mid y_i = +1\}$;
2. $\mathcal{F}_n = \cup \mathbf{B}_i \in \mathcal{F} \setminus \mathcal{F}^+$;
3. \mathbf{R}^- = 根据 k -Means 算法从 \mathcal{F}_n 中获得的 m 个聚类中心;
4. $\mathbf{R}^+ = \emptyset$;
5. for ($\mathbf{B}_i \in \mathcal{F}^+$) do;
6. for ($\mathbf{x}_{ij} \in \mathbf{B}_i$) do;
7. 根据 Eq. (2) 计算 s_{ij} ;
8. end for;
9. 根据 Eq. (3) 更新 j^* ;
10. $\mathbf{R}^+ \leftarrow \mathbf{R}^+ \cup \{\mathbf{x}_{ij^*}\}$;
11. end for.
12. 输出 \mathbf{R}^+ 和 \mathbf{R}^- .

2.4 关联性强化技术

关联性挖掘技术是利用正、负实例之间的差异性, 分别选择初始正、负代表实例集。为了进一步利用整个实例空间 \mathcal{R} 的特征信息, 关联性强化技术基于关联性指标设计集合内互评和集合外它评标准, 旨在强化初始正、负代表实例集, 以提升代表实例集的整体关联性, 如图 1(d)—(f) 所示。

首先, 使用集合内互评标准度量初始正、负代表实例集的内部关联性。根据式(1), 计算初始正代表实

例 $r_k^+ \in R^+$ 的关联性 $A(r_k^+, R^+)$ 和负代表实例 $r_k^- \in R^-$ 的关联性 $A(r_k^-, R^-)$ 。然后,基于集合外它评标准度量 X 中每个实例 x_j^* 相对于 R^+ 和 R^- 的外部关联性,即正关联性 $A(x_j^*, R^+)$ 和负关联性 $A(x_j^*, R^-)$ 。通过大小比较,若 x_j^* 的正关联性更强,则将 x_j^* 并入中间集 R^{+*} ,反之并入 R^{-*} ,其中 R^{+*} 和 R^{-*} 已分别初始化为 R^+ 和 R^- 。其次,分别比较 R^{+*} 和 R^{-*} 中实例的正关联性和负关联性,选择其中正关联性和负关联性最强的 m 个实例构成新的正代表实例集 R^+ 和负代表实例集 R^- 。最后,考虑到代表实例的平衡分布,将 $R^* = R^+ \cup R^-$ 聚类为 ψ 簇,获得最终的代表实例集 $C = \{c_k\}_{k=1}^{\psi}$,其中 c_k 为聚类中心^[27]。

2.5 包嵌入技术

包嵌入技术基于选定的 C ,将每个包转换为一个单向量,从而将 MIL 问题简化为 SIL 问题。如图 2 所示,首先,计算包 B_i 中每个实例 x_{ij} 与所有代表实例之间的距离,并将 x_{ij} 分配给离它距离最近的代表实例 c_k , x_{ij} 称为 c_k 的近邻实例。令 B_i 中分配给 c_k 的所有近邻实例记作 Ω_{ik} ,则 B_i 对于 c_k 的分量可以计算为

$$v_{ik} = \sum_{x_{ij} \in \Omega_{ik}} x_{ij} - c_k \quad (4)$$

得到 ψ 个分量之后,嵌入函数将这些分量串联成一个 $D = \psi \times d$ 维的向量 V_i ,

$$V_i = \parallel_{k=1}^{\psi} v_{ik} \quad (5)$$

最后,对嵌入向量进行归一化操作可以提升模型性能^[23]。使用 $v_{il} \leftarrow \text{sign}(v_{il}) \sqrt{|v_{il}|}$ 对 V_i 的每个元素进行处理,随后通过 $V_i \leftarrow V_i / \|V_i\|_2$ 将 V_i 归一化。在得到所有包的嵌入向量后,基于新特征空间训练一个 SIL 分类器 $\mathcal{F}(\cdot)$ 来预测新包的标签。

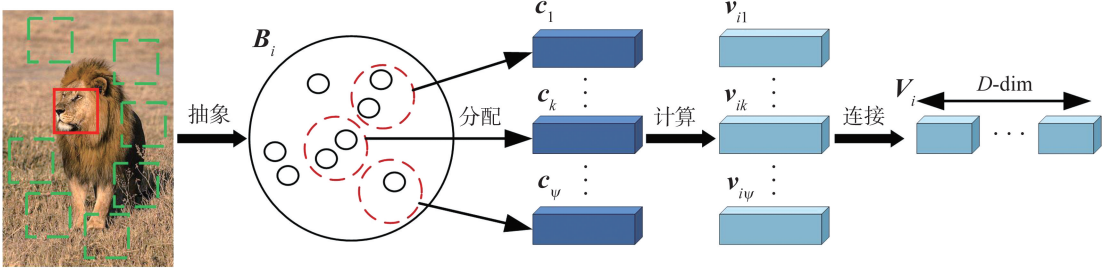


图 2 包嵌入技术示意图
Fig.2 Bag embedding technique

2.6 算法描述

算法 2 展示了 MEMR 的实现过程。第 1—2 行为关联性挖掘技术,根据算法 1 得到原始数据集 \mathcal{S} 的初始正、负代表实例集 R^+ 和 R^- 。第 3—14 行为关联性强化技术,其中:第 4 行获得实例空间 \mathcal{B} ;第 5 行通过关联性指标评估 R^+ 和 R^- 中实例的内部关联性;第 6 行初始化 R^{+*} 和 R^{-*} ;第 7—10 行度量 \mathcal{B} 中每个实例相对于 R^+ 和 R^- 的正关联性和负关联性,将正关联性更强的实例并入 R^{+*} ,反之并入 R^{-*} ;第 11—12 行从 R^{+*} 和 R^{-*} 中分别选取整体关联性最强的 m 个实例构成新的正、负代表实例集 R^+ 和 R^- ;第 13—14 行基于强化后的结果,通过聚类得到最终的代表实例集 C 。第 15—20 行为包嵌入技术,即计算整个数据集的嵌入向量并进行归一化操作。第 21 行基于新特征空间训练分类器 $\mathcal{F}(\cdot)$ 。在测试阶段,通过该分类器 $\mathcal{F}(V'_i)$,可以预测包 B'_i 的标签。

算法 2 MEMR

输入:

数据集 $\mathcal{S} = \{B_i\}_{i=1}^N$;

标签向量 $\mathbf{Y} = [y_1, y_2, \dots, y_N]$;

代表实例数量 ψ ;

输出:

分类器 $\mathcal{F}(\cdot)$;

代表实例集 C ;

训练:

1. / * Step 1. 初始化正负代表实例集 */;

2. $\mathbf{R}^+, \mathbf{R}^- = \text{insInit}(\mathcal{F})$;
3. /* Step 2. 生成代表实例集 */;
4. $\mathcal{B} = \bigcup_{i=1}^N \mathbf{B}_i$;
5. 根据 Eq. (1) 计算 $A(\mathbf{r}_k^+, \mathbf{R}^+)$ 和 $A(\mathbf{r}_k^-, \mathbf{R}^-)$;
6. $\mathbf{R}^{+*} = \mathbf{R}^+, \mathbf{R}^{-*} = \mathbf{R}^-$;
7. for ($\mathbf{x}_j^* \in \mathcal{B}^*$) do;
8. 计算 \mathbf{x}_j^* 的 $A(\mathbf{x}_j^*, \mathbf{R}^+)$ 和 $A(\mathbf{x}_j^*, \mathbf{R}^-)$;
9. If $A(\mathbf{x}_j^*, \mathbf{R}^+) > A(\mathbf{x}_j^*, \mathbf{R}^-)$ then $\mathbf{R}^{+*} = \cup \{\mathbf{x}_j^*\}$ else $\mathbf{R}^{-*} = \cup \{\mathbf{x}_j^*\}$;
10. end for;
11. \mathbf{R}^+ 是 \mathbf{R}^{+*} 中前 m 个强关联性的实例;
12. \mathbf{R}^- 是 \mathbf{R}^{-*} 中前 m 个强关联性的实例;
13. $\mathbf{R}^* = \mathbf{R}^+ \cup \mathbf{R}^-$;
14. 对 \mathbf{R}^* 使用 k -means 算法得到 ψ 个聚类中心 $\mathbf{C} = \{c_1, \dots, c_\psi\}$;
15. /* Step 3. 包嵌入 */;
16. for ($\mathbf{B}_i \in \mathcal{F}$) do;
17. 根据 Eq. (4) 和 Eq. (5) 计算 \mathbf{B}_i 的嵌入向量 \mathbf{V}_i ;
18. $v_{il} \leftarrow \text{sign}(v_{il}) \sqrt{|v_{il}|}$, 其中 v_{il} 是 \mathbf{V}_i 的第 l 个特征值;
19. $\mathbf{V}_i \leftarrow \mathbf{V}_i / \|\mathbf{V}_i\|_2$;
20. end for.
21. 使用 $\{(\mathbf{V}_i, y_i)\}_{i=1}^N$ 训练分类器 $\mathcal{F}(\cdot)$.
- 测试:
22. 对包 \mathbf{B}'_i 执行 17~19 步得到嵌入向量 \mathbf{V}'_i ;
23. 根据分类器得到预测结果 $\mathcal{F}(\mathbf{V}'_i)$.

3 实验与结果

为了验证 MEMR 算法的有效性,选取 4 个应用领域中 34 个 MIL 数据集进行实验,并与 7 个 MIL 算法进行对比。实验验证手段为 10 次十折交叉验证,性能评估指标采用准确率 (accuracy),SIL 分类器采用 kNN、J48、SVM,所有算法的最终性能展示为最优分类器的结果。主要从参数分析、消融实验、性能对比及效率分析展示实验结果。

3.1 数据集简介

实验使用的 4 类 MIL 数据集包括图像检索、医学图像、文本分类,及网页推荐。表 2 列出了这些数据集的关键属性。下面介绍每类数据集的领域知识。

表 2 数据集的详细属性
Table 2 Detailed properties of the used datasets

数据集	包数量			实例数	维度
	正包	负包	包		
Elephant	100	100	200	1 391	230
Fox	100	100	200	1 320	230
Tiger	100	100	200	1 220	230
Messidor	654	546	1 200	12 352	687
Ucsb_breast	26	32	58	2 002	708
Newsgroups	993	1 007	2 000	80 137	200
Web	490	527	1 017	30 807	6 211

(1) 图像检索:其任务是识别图像中是否包含用户感兴趣的目标对象。Elephant、Fox、Tiger^[28]是该领域常用的数据集,其中每一个图像视作一个包,图像中每一个区域视作一个实例。

(2) 医学图像:其任务是利用 MIL 辅助医疗诊断。Messidor^[29]是糖尿病视网膜病变筛查数据集,包括 654 位糖尿病患者和 546 位健康患者的眼底图像。Ucsb_breast^[30]包括 58 个弱标记的染色图像,数据来自 32 位良性和 26 位恶性乳腺癌患者。

(3) 文本分类:其任务是预测文本中是否包含目标主题的文章。Newsgroups^[31]是文本分类领域最常用的数据集之一,包括 20 个子数据集,每个数据集由 100 个包组成,每个包包包含来自 20 个不同主题的文章实例。

(4) 网页推荐:其任务是根据网页中链接的内容为用户推荐感兴趣的网页。每个网页表示为一个包,网页中的链接视为实例,用户感兴趣的链接代表正实例。Web^[32]包括 9 个子数据集,每个数据集中的实例均具有高维性。

3.2 对比算法

实验对比 4 类前沿的基于嵌入的 MIL 方法。(1) 基于统计的方法:Simple-MI^[17]无需参数设置;(2) 基于核的方法:miFV^[18]中高斯模型的分量数设置为 1;(3) 基于包的方法:ELDB^[16]的包选择模式设置为“a”,代表包比例设置为 0.9;(4) 基于实例的方法:MILFM^[22]聚类中心数量设置为 40,miVLAD^[23]聚类中心数量设置为 1,MILDLM^[24]中判别实例数量设置为包的数量,StableMIL^[25]的实例阈值设置为 0.25。

3.3 参数分析

根据算法 2 可知,MEMR 的分类性能受代表实例数量 ψ 的影响。图 3 展示了 ψ 在 5 个代表数据集上的参数分析结果。代表数据集选取于图像检索、医学图像、文本分类,及网页推荐 4 个领域,包括 Elephant、Ucsb_breast、News.aa、News.sm,及 Web4。图 3 中横坐标表示 ψ ,纵坐标表示准确率。结果显示,当 ψ 为 1 或 2 时,MEMR 在大多数情况下都能达到最佳性能,而随着数量增加,5 个数据集的准确率呈总体下降趋势且差异缩小,因此,MEMR 中 ψ 的取值范围为 {1,2}。

3.4 消融实验

图 4 展示了 MEMR 的消融实验结果。图中 MEM 表示直接对图 1(c) 中的初始正、负代表实例集聚类,选取代表实例集 C 。为了确保公平比较,对于同一数据集使用相同的参数设置,即代表实例数量和分类器相同。图中轴线表示准确率,每个数据集分别对应一根轴线。结果显示,MEMR 在 5 个数据集上的准确率均比 MEM 更高,证明了关联性强化技术在关联性挖掘技术的基础上提升了 C 的整体关联性,验证了本算法的有效性。

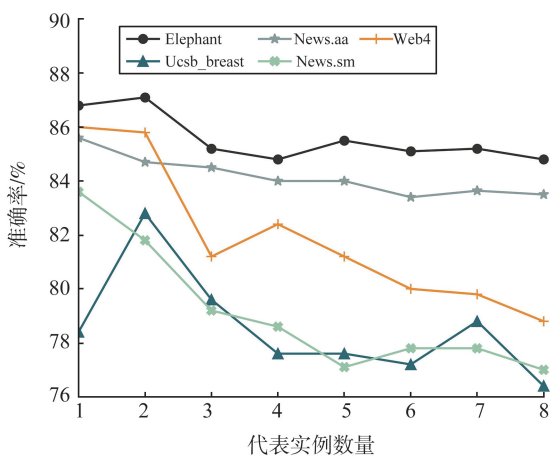


图 3 MEMR 代表实例数量的参数分析
Fig.3 Parameter analysis of the number of representative instances of MEMR

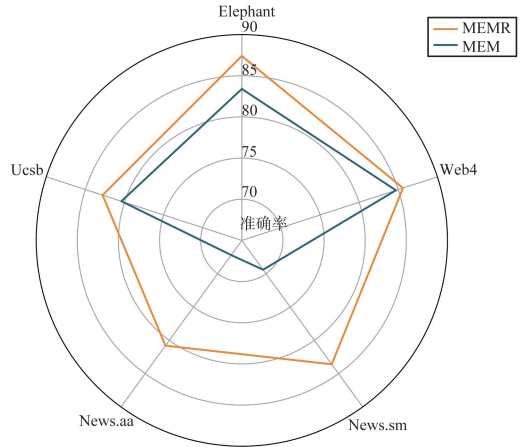


图 4 MEMR 与 MEM 的消融实验 (%)
Fig.4 Ablation experiment of MEMR and MEM (%)

3.5 性能对比

表 3—5 列出了所有对比算法在 4 类数据集上的最优平均准确率及相应标准差,每个数据集的最优结果使用粗体标注。表 3 展示了在图像检索和医学图像数据集上的对比结果。实验结果表明,MEMR 算法在图像检索数据集上有显著优势,尤其在 Tiger 数据集上,MEMR 的准确率高出第二名 6%,高出最后一名近

20%,这可能是由于关联性挖掘和强化技术可以更有效地选择代表实例。同时,MEMR在Messidor和Ucsb_breast数据集上也取得较好表现,仅次于miFV算法,原因可能是miFV的高斯混合模型可以更好地获取这类数据集的信息。表4给出所有算法在20个文本分类数据集上的实验结果,MEMR在超过一半的文本数据集上取得最优结果,并且明显优于MILFM、MILDLM、StableMIL,原因可能是文本数据集非常稀疏,实例的特征信息很相似,而MEMR的关联性强化技术利用整个实例特征空间对代表实例进行二次强化,在一定程度上提升了代表实例的代表性。表5展示了所有对比算法在网页推荐数据集上的结果。Web数据集的特点是高维和稀疏性,MEMR在大多数网页数据集上显示出最高的分类性能,在一定程度上证明了其应对高维和稀疏数据的能力。

表3 图像检索和医学图像数据集的平均准确率
Table 3 Average accuracy on image retrieval and medical image datasets

Dataset	Simple-MI	MILFM	miFV	miVLAD	MILDLM	Stable-MIL	ELDB	MEMR
Elephant	82.5±0.84	81.5±1.22	86.0±1.76	84.7±0.98	76.5±1.64	63.2±2.58	75.4±1.99	87.1±1.28
Fox	61.9±0.86	60.8±2.60	61.2±0.75	63.3±1.75	54.2±3.47	59.7±4.33	58.8±2.66	64.6±0.86
Tiger	81.1±1.16	76.3±1.29	79.1±0.58	84.9±7.63	69.0±1.41	65.7±2.06	67.4±2.73	85.1±0.37
Messidor	61.8±0.84	62.1±0.53	70.5±0.53	67.5±0.28	64.0±0.24	62.2±0.47	56.8±1.53	69.3±0.30
Ucsb_breast	81.2±2.71	55.6±2.33	85.6±0.80	80.0±1.79	56.0±2.19	54.4±0.20	63.0±7.62	82.8±4.12

表4 文本分类数据集的平均准确率
Table 4 Average accuracy on text categorization datasets

Dataset	Simple-MI	MILFM	miFV	miVLAD	MILDLM	Stable-MI	ELDB	MEMR
News.aa	83.6±0.80	52.6±7.86	83.8±1.60	84.0±2.28	54.6±7.06	52.6±4.03	84.6±1.78	86.0±2.28
News.cg	78.0±0.63	54.6±1.20	80.2±0.98	79.6±0.80	53.2±5.31	50.2±5.11	79.5±2.42	80.8±1.33
News.co	57.4±3.56	49.6±2.24	72.2±1.33	69.2±1.60	52.2±4.31	47.4±4.03	63.1±3.07	71.6±2.06
News.csi	75.4±0.80	57.6±3.20	79.8±1.17	80.0±1.55	56.6±6.62	50.2±5.49	78.1±2.02	81.2±2.40
News.csm	77.8±0.75	52.8±6.79	77.2±0.75	78.0±1.10	43.4±3.38	51.0±5.02	76.4±3.89	80.6±0.80
News.cw	71.0±3.16	57.8±2.79	86.6±0.80	82.6±1.02	56.8±4.31	54.2±4.49	79.6±1.43	81.0±1.10
News.mf	58.8±0.98	51.2±2.32	71.0±1.26	72.2±1.94	46.8±2.71	52.6±6.65	64.4±2.37	74.0±2.68
News.ra	75.4±0.49	52.4±1.62	78.4±1.36	81.6±1.02	51.8±6.68	52.0±5.02	71.5±2.27	82.4±1.74
News.rm	77.2±2.56	54.8±3.25	85.6±2.58	82.8±0.75	57.0±5.10	54.0±1.90	81.7±1.83	83.2±1.17
News.rsb	74.6±1.02	54.6±3.44	84.8±0.40	83.2±0.75	48.2±3.43	54.2±3.06	79.2±3.33	83.6±1.36
News.rsh	80.8±0.98	50.4±0.49	87.8±1.33	89.6±1.02	47.4±5.85	51.0±4.77	77.2±3.19	90.0±1.55
News.sc	73.8±0.40	58.6±1.85	75.2±1.60	83.0±1.10	49.4±3.83	50.2±4.53	70.0±2.87	82.8±1.60
News.se	92.0±0.00	53.0±0.00	92.6±0.80	92.4±0.49	55.6±1.96	51.0±3.58	88.2±1.32	92.4±1.36
News.sm	72.4±1.36	57.2±0.75	83.2±1.72	81.0±1.10	52.6±4.96	51.0±7.16	80.9±2.81	83.6±1.36
News.sr	77.4±0.80	50.2±1.17	79.8±2.56	79.8±2.32	51.4±2.80	57.6±2.65	80.7±1.16	80.6±1.62
News.ss	82.2±0.40	54.2±1.72	87.2±1.17	85.6±1.20	50.2±2.79	50.0±1.41	78.8±1.75	88.6±1.20
News.tpg	77.2±1.17	52.6±1.36	77.8±1.17	81.8±0.75	44.0±4.77	51.0±2.61	75.4±2.91	81.0±1.67
News.tpmd	83.0±1.79	60.0±3.58	79.0±0.63	83.6±1.02	55.6±2.94	55.6±4.59	76.7±2.26	85.2±1.17
News.tpmc	66.2±3.82	62.4±1.02	75.8±1.47	76.0±1.67	53.6±4.03	56.8±2.99	65.5±1.78	77.2±0.75
News.trm	61.6±1.02	52.6±1.02	75.0±1.10	78.0±2.28	47.2±3.87	51.0±3.63	66.4±2.27	76.2±1.47

表5 网页推荐数据集的平均准确率
Table 5 Average accuracy on Web recommendation datasets

Dataset	Simple-MI	MILFM	miFV	miVLAD	MILDLM	Stable-MI	ELDB	MEMR
Web1	80.7±2.53	81.6±0.86	83.4±1.06	79.6±1.09	83.6±1.15	83.0±1.23	81.1±1.72	84.0±1.36
Web2	82.3±2.34	81.0±0.36	83.0±0.93	80.0±1.00	82.7±0.81	82.7±1.00	74.2±2.31	81.4±1.56
Web3	80.7±2.90	81.9±1.68	82.1±0.93	81.2±1.87	81.6±0.73	81.4±1.23	79.0±1.99	82.0±1.45
Web4	80.9±1.00	79.8±2.26	80.3±1.09	83.8±0.36	78.7±1.36	77.6±0.45	79.1±2.59	85.6±1.56
Web5	78.1±1.52	77.8±2.41	78.1±1.15	82.5±0.89	79.0±1.41	78.1±0.61	74.1±2.82	83.2±1.09
Web6	79.8±2.90	82.0±1.34	77.8±0.45	84.9±0.93	82.7±1.52	76.5±0.68	79.7±1.77	86.7±1.87
Web7	64.1±0.73	60.0±1.99	68.0±2.02	72.9±2.18	61.8±4.53	61.0±3.06	52.6±3.49	74.5±1.41
Web8	64.7±1.45	61.6±2.66	71.2±1.59	76.3±2.70	56.1±2.02	59.0±3.19	48.0±3.16	78.9±2.10
Web9	68.0±2.53	57.4±3.01	74.0±4.05	77.2±1.15	56.7±2.34	54.9±3.73	46.5±2.93	81.2±2.97

为了比较不同算法的整体性能,使用 Bonferroni-Dunn 检验^[33] 比较各算法的平均排名。平均排名表示算法在每个数据集上准确率结果排名的平均值,如表 6 所示。MEMR 在图像检索、文本分类,及网页推荐数据集上排名第一,在医学图像数据集上排名第二,并且具有最高整体平均排名。这个结果表明,MEMR 可以应用于大多数 MIL 分类任务,并且具有较好的整体性能。图 5 报告了在 0.05 显著性水平下的临界差异 (critical difference, CD) 图。图中每个算法的平均排名沿轴标记,处于右侧的算法排名更高。将与 MEMR 平均排名之差小于一个 CD 域 (CD=1.60) 的算法用粗线连接起来。结果显示,MEMR 优于连接的 miFV 和 miVLAD 算法,并显著优于在 CD 域外的其他算法。

表 6 8 个算法在 4 类数据集上的平均排名
Table 6 Mean rank of eight algorithms on four classes of datasets

Datasets	Simple-MI	MILFM	miFV	miVLAD	MILDLM	Stable-MIL	ELDB	MEMR
图像检索	3.33	5.00	3.33	2.33	6.67	7.33	7.00	1.00
医学图像	5.00	6.50	1.00	3.50	5.00	6.50	6.50	2.00
文本分类	4.45	6.45	2.55	2.25	7.25	7.30	4.20	1.55
网页推荐	4.89	5.11	3.33	3.67	4.33	5.78	7.33	1.56
平均排名	4.50	5.97	2.74	2.71	6.29	6.85	5.41	1.53

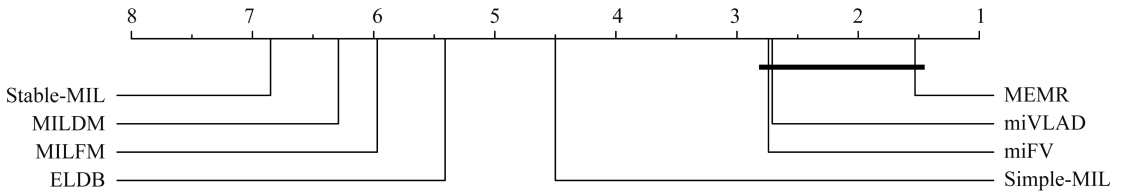


图 5 MEMR 与 7 种算法的 Bonferroni-Dunn 检验
Fig.5 Bonferroni-Dunn Test of MEMR and 7 algorithms

3.6 效率分析

表 7 展示了 MEMR 与 7 种对比算法在 5 个数据集上的时间复杂度以及 CPU 运行时间。MEMR 的时间成本包括构建代表实例集与特征转换的时间消耗。其中构建代表实例集包括选取初始代表实例集及强化代表实例集关联性两个步骤,时间成本均为 $O(dn)$ 。此外,特征转换的成本为 $O(dN)$,因此,MEMR 的时间复杂度为 $O(dn)$,其中 d 是实例维度, n 是实例空间的大小, N 是数据集的大小。CPU 运行时间的平均排名显示,MEMR 的速度仅次于 Simple-MI 和 miVLAD,这可能是因为 Simple-MI 不需要计算实例间的距离。miVLAD 通过时间复杂度较低的 k -means 算法选取代表实例。MEMR 的分类性能优于这 2 种算法。

表 7 8 个算法在 5 个数据集上一次 10CV 的 CPU 运行时间
Table 7 The CPU runtime (s) of one time 10CV of the eight algorithms on the five datasets

Datasets ($d/n/N$)	Simple-MI	MILFM	miFV	miVLAD	MILDLM	Stable-MIL	ELDB	MEMR
Time Complexity	$O(dN)$	$O(dn^2)$	$O(dn)$	$O(dn)$	$O(dn^2)$	$O(dn^2)$	$O(dn^2)$	$O(dn)$
Elephant (230/1 320/200)	0.141	39.974	5.561	1.062	34.179	16.962	24.854	5.234
Ucsb_breast (708/2 002/58)	0.125	83.145	8.576	1.390	65.969	152.230	63.860	4.624
News.aa (200/5 443/100)	0.078	446.780	5.717	1.250	382.254	66.806	349.120	3.905
News.sm (200/3 094/100)	0.063	130.253	5.030	1.094	121.533	40.662	114.382	3.406
Web4 (6 059/3 423/113)	1.109	533.881	711.693	9.748	453.940	1 169.954	785.484	44.241
Mean rank	1.00	7.40	4.00	2.00	6.40	6.20	6.00	3.00

4 总结与展望

本文提出多示例嵌入学习的实例关联性挖掘与强化算法。通过利用正实例与负实例之间的差异性,得到初始正、负代表实例集;进一步地,利用整个实例空间的特征信息,度量代表实例的整体关联性,对代表实例集进行强化。实验结果显示,MEMR 算法在近 62% 的数据集上取得最好的实验结果。同时,MEMR 的平均排名为 1.53,优于其余 7 种对比算法,尤其是在图像检索和网页推荐数据集上的排名分别达到 1.00 和 1.56。除此之

外, MEMR 的效率优于大部分 MIL 前沿算法。以上实验结果证明 MEMR 具有优秀的分类性能。

MEMR 进一步的研究与改进方向包括:(1) 更全面的度量指标。目前 MEMR 对初始代表实例集和最终代表实例的选择仅基于关联性指标,未来还可以设计多角度的度量方式;(2) 更有效的强化技术。目前基于关联性排序实现强化过程,未来可以研究更灵活的增强方式。

参考文献:

- [1] DIETTERICH T G, LATHROP R H, LOZANO-PEREZ T. Solving the multiple instance problem with axis-parallel rectangles [J]. *Artificial Intelligence*, 1997, 89(1/2):31-71.
- [2] LI Daxiang, ZHANG Yue. Multi-instance learning algorithm based on LSTM for Chinese painting image classification [J]. *IEEE Access*, 2020, 8:179336-179345.
- [3] SHARMA Y, SHRIVASTAVA A, EHSAN L, et al. Cluster-to-conquer: a framework for end-to-end multi-instance learning for whole slide image classification [C]//*Medical Imaging with Deep Learning*. Lübeck: PMLR, 2021:682-698.
- [4] 王刚,许信顺.一种新的基于多示例学习的场景分类方法[J]. *山东大学学报(理学版)*, 2010, 45(7):108-113.
WANG Gang, XU Xinchun. A new multi-instance learning method for scene classification [J]. *Journal of Shandong University (Natural Science)*, 2010, 45(7):108-113.
- [5] YI Lin, ZHANG Honggang. Regularized instance embedding for deep multi-instance learning [J]. *Applied Sciences*, 2020, 10(1):64-78.
- [6] KUMAR V, CHEMMENGATH S, GUPTA Y, et al. Multi-instance training for question answering across table and linked text [J/OL]. *arXiv*, 2021. <https://arxiv.org/abs/2112.07337>.
- [7] YANG Mei, ZENG Wenxi, MIN Fan. Multi-instance embedding learning through high-level instance selection [C]//*Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Chengdu: Springer, 2022:122-133.
- [8] TIAN Yuchi, WANG Jiawei, YANG Wenjie, et al. Deep multi-instance transfer learning for pneumothorax classification in Chest X-Ray images [J]. *Medical Physics*, 2022, 49(1):231-243.
- [9] MANIVANNAN S, COBB C, BURGESS S, et al. Subcategory classifiers for multiple-instance learning and its application to retinal nerve fiber layer visibility classification [J]. *IEEE Transactions on Medical Imaging*, 2017, 36(5):1140-1150.
- [10] WEI Xiushen, YE Hanjia, MU Xin, et al. Multi-instance learning with emerging novel class [J]. *IEEE Transactions on Knowledge and Data Engineering*, 2019, 33(5):2109-2120.
- [11] ZHANG Yalin, ZHOU Zhihua. Multi-instance learning with key instance shift [C]//*International Joint Conference on Artificial Intelligence*. Melbourne: IJCAI, 2017:3441-3447.
- [12] XU Xin, FRANK E. Logistic regression and boosting for labeled bags of instances [C]//*Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Berlin: Springer, 2004:272-281.
- [13] MELKI G, CANO A, VENTURA S. MIRSVM: multi-instance support vector machine with bag representatives [J]. *Pattern Recognition*, 2018, 79:228-241.
- [14] LUAN Tianxiang, LUO Tingjin, ZHUGE Wenzhang, et al. Optimal representative distribution margin machine for multi-instance learning [J]. *IEEE Access*, 2020, 8:74864-74874.
- [15] CHIKONTWE P, KIM M, NAM S J, et al. Multiple instance learning with center embeddings for histopathology classification [C]//*International Conference on Medical Image Computing and Computer-Assisted Intervention*. Lima: Springer, 2020: 519-528.
- [16] YANG Mei, ZHANG Yuxuan, WANG Xizhao, et al. Multi-instance ensemble learning with discriminative bags [J]. *IEEE Transactions on Systems, Man, and Cybernetics; Systems*, 2021, 52(9):5456-5467.
- [17] AMORES J. Multiple instance classification: review, taxonomy and comparative study [J]. *Artificial Intelligence*, 2013, 201: 81-105.
- [18] WEI Xiushen, WU Jianxin, ZHOU Zhihua. Scalable multi-instance learning [C]//*International Conference on Data Mining*. Shenzhen: IEEE, 2014:1037-1042.
- [19] ZHANG Minling, ZHOU Zhihua. Multi-instance clustering with applications to multi-instance prediction [J]. *Applied Intelligence*, 2009, 31(1):47-68.
- [20] CHEN Yixin, BI Jinbo, WANG J Z. MILES: multiple-instance learning via embedded instance selection [J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2006, 28(12):1931-1947.
- [21] LI Wujun. MILD: multiple-instance learning via disambiguation [J]. *IEEE Transactions on Knowledge and Data Engineering*, 2009, 22(1):76-89.
- [22] HONG Richang, WANG Meng, GAO Yue, et al. Image annotation by multiple-instance learning with discriminative feature

- mapping and selection[J]. IEEE Transactions on Cybernetics, 2013, 44(5):669-680.
- [23] WEI Xiushen, WU Jianxin, ZHOU Zhihua. Scalable algorithms for multi-instance learning[J]. IEEE Transactions on Neural Networks and Learning Systems, 2017, 28(4):975-987.
- [24] WU Jia, PAN Shirui, ZHU Xingquan, et al. Multi-instance learning with discriminative bag mapping[J]. IEEE Transactions on Knowledge and Data Engineering, 2018, 30(6):1065-1080.
- [25] ZHANG Weijia, LI Jiuyong, LIU Lin. Robust multi-instance learning with stable instances[J/OL]. arXiv, 2019. <https://arxiv.org/pdf/1902.05066.pdf>.
- [26] FU Zhouyu, ROBLES-KELLY A, ZHOU Jun. MILIS: multiple instance learning with instance selection[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2010, 33(5):958-977.
- [27] MIN Fan, ZHANG Shiming, CIUCCI D, et al. Three-way active learning through clustering selection[J]. International Journal of Machine Learning and Cybernetics, 2020, 11(5):1033-1046.
- [28] ANDREWS S, TSOCHANTARIDIS I, HOFMANN T. Support vector machines for multiple-instance learning[J]. Neural Information Processing Systems, 2002, 14:561-568.
- [29] DECENCIERE E, ZHANG Xiwei, CAZUGUEL G, et al. Feedback on a publicly distributed image database: the Messidor database[J]. Image Analysis and Stereology, 2014, 33(3):231-234.
- [30] KANDEMIR M, HAMPRECHT F A. Computer-aided diagnosis from weak supervision: a benchmarking study[J]. Computerized Medical Imaging and Graphics, 2015, 42:44-50.
- [31] ZHOU Zhihua, SUN Yuyin, LI Yufeng. Multi-instance learning by treating instances as non-iid samples[C]//Proceedings of the 26th Annual International Conference on Machine Learning. Montreal: ACM, 2009:1249-1256.
- [32] XU Bicun, TING Kaiming, ZHOU Zhihua. Isolation set-kernel and its application to multi-instance learning[C]//Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. Anchorage: ACM, 2019:941-949.
- [33] DEMSAR J. Statistical comparisons of classifiers over multiple data sets[J]. The Journal of Machine Learning Research, 2006, 7:1-30.

(编辑:于善清)

(上接第 10 页)

- [23] LAI T L, ROBBINS H. Asymptotically efficient adaptive allocation rules[J]. Advances in Applied Mathematics, 1985, 6(1):4-22.
- [24] DWORK C. Differential privacy: a survey of results[C]//International Conference on Theory and Applications of Models of Computation. Heidelberg: Springer, 2008:1-19.
- [25] 方滨兴. 释放数据使用权将成为未来技术发展取向[N/OL]. 中国新闻网, 2022-05-19[2023-12-05], <https://news.sciencenet.cn/htmlnews/2022/5/479297.shtm>.
- [26] WASSERMAN L, ZHOU S. A statistical framework for differential privacy[J]. Journal of the American Statistical Association, 2010, 105(489):375-389.
- [27] DUCHI J C, JORDAN M I, WAINWRIGHT M J. Privacy aware learning[J]. Journal of the ACM, 2014, 61(6):1-57.
- [28] LI T, SAHU A K, TALWALKAR A, et al. Federated learning: challenges, methods, and future directions[J]. IEEE Signal Processing Magazine, 2020, 37(3):50-60.
- [29] TAN C, SUN F, KONG T, et al. A survey on deep transfer learning[C]// International Conference on Artificial Neural Networks. Cham: Springer, 2018:270-279.
- [30] FINN C, XU, K, LEVINE S. Probabilistic model-agnostic meta-learning[C]// NIPS'18: Proceedings of the 32nd International Conference on Neural Information Processing Systems. Montréal: Curran Associates Inc., 2018:9537-9548.
- [31] VILALTA R, DRISSI Y. A perspective view and survey of meta-learning[J]. Artificial Intelligence Review, 2002, 18(2):77-95.
- [32] 张钹. 人工智能进入后深度学习时代[J]. 智能科学与技术学报, 2019, 1(1):4-6.
ZHANG Ba. Artificial intelligence is entering the post deep-learning era[J]. Chinese Journal of Intelligent Science and Technology, 2019, 1(1):4-6.
- [33] 张钹, 朱军, 苏航. 迈向第三代人工智能[J]. 中国科学: 信息科学, 2020, 50:1281-1302.
ZHANG Ba, ZHU Jun, SU Hang. Toward the third generation of artificial intelligence[J]. Science China: Information Sciences, 2020, 50:1281-1302.

(编辑:李艺)