



智能系统学报

CAAI TRANSACTIONS ON INTELLIGENT SYSTEMS

对不平衡目标域的多源在线迁移学习

周晶雨, 王士同

引用本文:

周晶雨,王士同. 对不平衡目标域的多源在线迁移学习[J]. 智能系统学报, 2022, 17(2): 248–256.

ZHOU Jingyu,WANG Shitong. Multi-source online transfer learning for imbalanced target domains[J]. *CAAI Transactions on Intelligent Systems*, 2022, 17(2): 248–256.

在线阅读 View online: <https://dx.doi.org/10.11992/tis.202012019>

您可能感兴趣的其他文章

面对类别不平衡的增量在线序列极限学习机

Incremental online sequential extreme learning machine for imbalanced data

智能系统学报. 2020, 15(3): 520–527 <https://dx.doi.org/10.11992/tis.201904040>

基于最小最大概率机的迁移学习分类算法

Transfer learning classification algorithms based on minimax probability machine

智能系统学报. 2016, 11(1): 84–92 <https://dx.doi.org/10.11992/tis.201505024>

SMOTE过采样及其改进算法研究综述

Summary of research on SMOTE oversampling and its improved algorithms

智能系统学报. 2019, 14(6): 1073–1083 <https://dx.doi.org/10.11992/tis.201906052>

动态平衡采样的不平衡数据集分类方法

Imbalanced data ensemble classification using dynamic balance sampling

智能系统学报. 2016, 11(2): 257–263 <https://dx.doi.org/10.11992/tis.201507015>

应用于不平衡多分类问题的损失平衡函数

Application of the loss balance function to the imbalanced multi-classification problems

智能系统学报. 2019, 14(5): 953–958 <https://dx.doi.org/10.11992/tis.201808004>

一种基于密度的SMOTE方法研究

Research on the SMOTE method based on density

智能系统学报. 2017, 12(6): 865–872 <https://dx.doi.org/10.11992/tis.201706049>

微信公众平台



关注微信公众号，获取更多资讯信息

DOI: 10.11992/tis.202012019

网络出版地址: <https://kns.cnki.net/kcms/detail/23.1538.TP.20210621.1427.002.html>

对不平衡目标域的多源在线迁移学习

周晶雨, 王士同

(江南大学人工智能与计算机学院, 江苏 无锡 214122)

摘要: 多源在线迁移学习已经广泛地应用于相关源域中含有大量的标记数据且目标域中数据以数据流的形式达到的应用中。然而, 目标域的分类分布有时是不平衡的, 针对目标域每次以在线方式到达多个数据的不平衡二分类问题, 本文提出了一种可以对目标域样本过采样的多源在线迁移学习算法。该算法从前面批次的样本中寻找当前批次的样本的 k 近邻, 先少量生成多数类样本, 再生成少数类使得当前批次样本的分类分布平衡。每个批次合成样本和真实样本一同训练目标域函数, 从而提升目标域函数的分类性能。同时, 分别设计了在目标域的输入空间和特征空间过采样的方法, 并且在多个真实世界数据集上进行了综合实验, 证明了所提出算法的有效性。

关键词: 多源迁移学习; 在线学习; 目标域; 不平衡数据; 过采样; k 近邻; 输入空间; 特征空间

中图分类号: TP181 **文献标志码:** A **文章编号:** 1673-4785(2022)02-0248-09

中文引用格式: 周晶雨, 王士同. 对不平衡目标域的多源在线迁移学习 [J]. 智能系统学报, 2022, 17(2): 248-256.

英文引用格式: ZHOU Jingyu, WANG Shitong. Multi-source online transfer learning for imbalanced target domains[J]. CAAI transactions on intelligent systems, 2022, 17(2): 248-256.

Multi-source online transfer learning for imbalanced target domains

ZHOU Jingyu, WANG Shitong

(School of Artificial Intelligence and Computer Science, Jiangnan University, Wuxi 214122, China)

Abstract: Multi-source online transfer learning has been widely used in applications where the relevant source domain contains a large amount of labeled data and the data in the target domain is achieved in the form of data flow. However, the class distribution of the target domain is sometimes imbalanced. Aiming at the unbalanced binary classification problem wherein the target domain reaches multiple data online at a time, this paper proposes a multi-source online transfer learning algorithm by means of oversampling the target domain samples. First, the algorithm finds the k -nearest neighbors of the current batch of samples from the previous batch, then generates a small number of majority class samples, finally generating a minority class to balance the class distribution of the current batch of samples. Each batch of synthetic and real samples train the target domain function together, thereby improving the classification performance of the target domain function. At the same time, methods for oversampling in the input space and feature space of the target domain are designed respectively, and comprehensive experiments are conducted on multiple real-world data sets to prove the effectiveness of the proposed algorithm.

Keywords: multi-source transfer learning; online learning; target domain; imbalanced data; oversampling; k -nearest neighbor; input space; feature space

迁移学习的主要目的是利用源域的知识来提高目标域的学习性能, 多年来进行了广泛的研究^[1]。使用一些分布相似的现有数据来提取有用的信

息, 可以解决目标域的训练数据有限或标记成本太高的问题。在许多实际应用中, 与目标域分布相似的离线源域有多个, 所以可以轻松地从这些源域中收集辅助信息。为了应对不同来源对与目标域的贡献不同的问题, 许多复杂的基于提升方法的多源迁移学习算法^[2-3]被设计。基于提升方

收稿日期: 2020-12-16. 网络出版日期: 2021-06-21.

基金项目: 国家自然科学基金项目 (61572236).

通信作者: 王士同. E-mail: wxwangst@aliyun.com.

法的算法根据贡献高低对多个源域附加权重来生成集成分类器,合理利用每个源域的知识。

多源迁移学习通过多个源域中提取的知识来改善目标域上的学习任务的性能,近年来得到了越来越多的关注。Qian等^[4]提出了一个多域鲁棒优化的框架,用于学习多个域的单一模型。Huffman等^[5]提出了一种确定交叉熵损失和其他损失分布加权组合解的多源自适应算法。Peng等^[6]提出了多源域自适应矩匹配方法,利用多源域特征分布的矩进行动态对齐,将知识从多标记源域转移到未标记目标域。Kang等^[7]提出了一种在线多源多分类转移学习算法。这些现有的算法可以从多个源域迁移知识到目标域,而本文的目标是解决源域和目标域数据类别不平衡的多源在线迁移学习问题。

现有的大多数迁移学习工作都假设事先提供了源域和目标域的训练数据^[8]。但是,在某些实际应用中,目标域的数据可能以在线的方式到达。近十年,在线学习^[9-10]得到了广泛的研究。在线学习中,分类器在每个回合中接收一个实例及其标签,然后预测该实例,并根据预测结果和真实标签的损失信息更新分类器。Wang等^[11]提出一种基于最大最小概率机的迁移学习分类算法。Zhao等^[12]提出一种可以立即响应的且高效的在线学习算法来解决在线迁移学习任务。Wu等^[13]提出了一种具有多个源域的在线迁移学习算法,当目标数据到达时,多个源域分类器和目标域分类器同时做出预测,根据各分类器的权重组合最终预测结果,并更新各分类器的权重。

目前,大多数在线迁移学习都默认目标域的分类分布是平衡的,然而现实中存在很多不平衡的数据。例如,机器的故障诊断,医疗诊断以及军事应用等。在大多数现实世界的问题中,少数类实例的错误分类代价往往很大,减少少数类错误分类是至关重要的。处理不平衡数据集的方法可以分为对数据的采样方法^[14]、成本敏感方法和算法级方法^[15]。采样方法对数据集进行预处理,将类别修改至相对平衡。成本敏感方法对错误分类少数类实例的决策函数施加更大的惩罚。算法级的方法直接修改分类器来处理不平衡问题。

因此,本文提出一种针对目标域不平衡的多源在线迁移学习算法。其中,目标域每次到达一批数据。在算法中,从前面已经到达的批次中寻找当前批次样本的 k 近邻,形成种子和邻居对。然后在样本对之间的线段上适量生成合成的多数类样本,再合成少数类样本使当前批次的类别分布

相对平衡。考虑到不同批次的样本之间的特征分布可能发生细微的偏移,生成样本时控制合成样本近似于当前批次中的样本。最后用新生成的样本去改进目标函数,然后再对当前批次的所有样本按序进行在线迁移学习,从而提升整体分类器对少数类的分类性能。此外,还分别设计了在目标域的输入空间和特征空间过采样的方法。在目标域的输入空间生成数据点来平衡类别分布,可以提高目标函数对少数类的分类性能,但也可能生成不代表非线性可分问题的数据点,影响函数精度。所以设计了在目标域特征空间过采样的方法,与文献^[16]不同,本文的方法在特征空间生成数据点来训练在线的函数,生成少数样本会导致类别分布得更具代表性,可以克服非线性问题的局限。

1 在线迁移学习

简要介绍多源在线迁移学习算法 HomOTLMS。HomOTLMS 根据预先给出的源域数据,在离线批处理学习范式中构建 n 个源域的决策函数 $(h^{s_1}, h^{s_2}, \dots, h^{s_n})$ 。而在线部分使用被动攻击算法 (passive aggressive, PA)^[17],在目标域上构造一个以在线的方式更新的决策函数 h^T , T 为目标函数。对于当前到达的实例 \mathbf{x}_j , 计算目标域决策函数的铰链损失:

$$\ell_j = \max(0, 1 - y_j h^T(\mathbf{x}_j)) \quad (1)$$

如果决策函数遭受非零损失,则根据式(2)更新目标域函数和添加支持向量:

$$h_{j+1}^T = h_j^T + \tau_j y_j k(\mathbf{x}_j, \cdot) \quad (2)$$

式中:支持向量系数 $\tau_j = \min\{C, \ell_j / k(\mathbf{x}_j, \mathbf{x}_j)\}$, $k(\cdot, \cdot)$ 是核函数。

然后使用一个权重向量 $\mathbf{v}_j = (v_j^1, v_j^2, \dots, v_j^n)$ 和一个权重变量 w_j 去分别表示 n 个源决策函数和目标决策函数的权重。对于做出错误预测的决策函数,需要将其权重降低。对于源决策函数,令 $v_{j+1}^i = v_j^i \alpha$; 对于目标决策函数,令 $w_{j+1} = w_j \alpha$, 其中 $\alpha \in (0, 1)$ 是权重折扣参数。与此同时要保持所有决策函数前面的权重之和为 1, 所以需要归一化权重,即

$$p_j^i = v_j^i / \left(\sum_{i=1}^n v_j^i + w_j \right), q_j = w_j / \left(\sum_{i=1}^n v_j^i + w_j \right) \quad (3)$$

式中: p_j^i 和 q_j 分别是第 j 个实例到来时,第 i 个源决策函数和目标决策函数前面的权重。所以最终集成的决策函数为

$$f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^n p_j^i \text{sign}(h^{s_i}(\mathbf{x}_j)) + q_j \text{sign}(h_j^T(\mathbf{x}_j)) \right) \quad (4)$$

上述算法能够有效解决多个源域的在线迁移学习问题,但并不能应对目标域不平衡的情况。下面介绍了一种新的在线迁移学习方法,可以在在线预测的过程中,人工平衡目标域类别的分布,从而降低总体分类误差。

2 不平衡目标域的在线迁移学习

2.1 问题描述

在多源迁移学习的问题中,对于给定的 n 个源域,用 $D^s = \{D^{s_1}, D^{s_2}, \dots, D^{s_n}\}$ 表示,目标域用 D^t 表示。对于第 i 个源域 D^{s_i} ,源域数据空间用 $\mathcal{X}_{s_i} \times \mathcal{Y}_{s_i}$ 表示,其中特征空间是 $\mathcal{X}_{s_i} = \mathbf{R}^{d_i}$ 。用 $\mathcal{X} \times \mathcal{Y}$ 表示目标域的数据空间,其中特征空间是 $\mathcal{X} = \mathbf{R}^d$ 。这里,源域和目标域共享相同的标签空间 $\mathcal{Y}_{s_i} = \mathcal{Y} = \{+1, -1\}$ 。

在在线学习的部分,目标域数据 $\{(\mathbf{x}_j, y_j)\}_{j=1}^m \in \mathcal{X} \times \mathcal{Y}$ 的类别分布是不平衡的,正类样本少于负类样本。当目标数据以在线的方式到达,并且每次到达一批数据时,每批数据中正类和负类样本的分布也是不平衡的。目标域第 b 个批次的的数据可以表示为 $\{(\mathbf{x}_j, y_j)\}_{j=1}^{l_b}$ 。

目标域采用被动攻击算法(PA)学习决策函数,当目标域的数据不平衡时,目标决策函数会更加偏向于多数类。若能在在线学习的过程中,扩充每个批次少数类的样本,就可能实现目标领域对少数类更准确的分类。考虑到目标域整体的样本个数有限,可以通过先扩增每个批次的多数类,然后再扩增少数类样本至平衡,提高目标域函数的整体分类性能,从而更好地实现知识迁移。

2.2 在输入空间过采样的在线迁移学习

本节提出一种称为 OTLMS_IO(online transfer learning multi-source input space oversampling) 的算法,该算法代表在目标域的输入空间进行过采样的多源在线迁移学习。OTLMS_IO 通过增加每个批次中多数类和少数类样本的个数来提升目标域函数的性能。

目标域的数据以在线的方式分批到达,每次到达多个实例。第 b 个批次到达的实例是 $\{(\mathbf{x}_j, y_j)\}_{j=1}^{l_b}$, 对于其中每个少数类实例,都以欧氏距离(式(5))为标准计算它到前面已经到达批次的所有少数类实例的距离,得到其 k 近邻。

$$\text{dist}(\mathbf{x}_p, \mathbf{x}_q) = \sqrt{\sum_{d=1}^m (x_{p,d} - x_{q,d})^2} \quad (5)$$

式中: \mathbf{x}_p 是当前批次中的实例,称为种子; \mathbf{x}_q 是前面批次中的实例,称为邻居, m 是实例的维数。

然后将种子和邻居组合成样本对 $\{(\mathbf{x}_p, \mathbf{x}_q)\}_{j=1}^{l_b^{\min} \times k}$, 一共 $l_b^{\min} \times k$ 个, l_b^{\min} 是当前批次 b 中少数类实例的个数。以同样的方式,可以得到当前批次中多数类实例形成的样本对 $\{(\mathbf{x}_p, \mathbf{x}_q)\}_{j=1}^{l_b^{\max} \times k}$, 共 $l_b^{\max} \times k$ 个。从少数类和多数类的样本对中分别选取 \min_num 和 \max_num 个,用于生成新样本。 \max_num 的大小决定了当前批次生成样本和真实样本整体的规模, \min_num 使得当前批次类别平衡。根据式(6)在每个样本对之间的线段上生成新样本。

$$\mathbf{x}_{\text{new}} = \mathbf{x}_p + (\mathbf{x}_q - \mathbf{x}_p) \times \delta \quad (6)$$

同时,考虑到不同批次样本之间的特征分布可能会发生细微的偏移,所以控制均匀分布的随机数 $\delta \in [0, 0.5]$, 使得生成的新样本更加靠近当前批次中的样本。

对生成的一共 t 个新样本分配相应的标签,在当前批次的样本训练之前,使用新生成的样本 $\{(\mathbf{x}_j, y_j)\}_{j=1}^t$ 改进目标函数,根据式(7):

$$h_b^T(\mathbf{x}) = h^T(\mathbf{x}) + \sum_{j=1}^t \tau_j y_j(\mathbf{x}_j, \mathbf{x}) \quad (7)$$

使用在线被动攻击算法可以轻松学得用新样本改进后的分类器,即根据式(2)对将铰链损失 $\ell > 0$ 的新实例都作为支持向量添加到支持向量集中。最后再使用集成决策函数(式(4))分别训练当前批次到达的所有实例,并按照上述方法对后面所有批次进行同样的操作可以得到训练好的集成函数。

2.3 在特征空间过采样的在线迁移学习

与在输入空间过采样不同,本节提出了一种称为 OTLMS_FO(online transfer learning multi-source feature space oversampling) 的算法,该算法表示在特征空间过采样的多源在线迁移学习。目标域的函数通过核函数进行预测,所以 OTLMS_FO 能利用与 SVM 分类器相同的核技巧,合成样本利用特征空间中的点积生成而不需要知道特征映射函数 $\phi(\mathbf{x})$ 。特征空间生成数据点在高维的空间具有更好的线性可分性,可以用来改进目标函数。

OTLMS_FO 算法在目标域第 b 个批次的样本 $\{(\mathbf{x}_j, y_j)\}_{j=1}^{l_b}$ 到达时,从中挑选出少数类样本和多数类样本。然后从前面已经到达的批次中分别找到当前到达批次中少数类和多数类样本的 k 近邻。由于是在特征空间中计算样本间的距离,需要将种子 \mathbf{x}_p 和近邻 \mathbf{x}_q 映射为特征空间的 $\phi(\mathbf{x}_p)$ 和 $\phi(\mathbf{x}_q)$, 然后计算两个实例之间的距离。特征空间中,两个实例之间的距离为

$$d^{\phi}(\mathbf{x}_p, \mathbf{x}_q)^2 = \|\phi(\mathbf{x}_p) - \phi(\mathbf{x}_q)\|^2 = k(\mathbf{x}_p, \mathbf{x}_p) - 2k(\mathbf{x}_p, \mathbf{x}_q) + k(\mathbf{x}_q, \mathbf{x}_q) \quad (8)$$

根据式(8)可以找到当前批次中的每个少数类样本的 k 近邻, 种子和邻居组成的样本对构成集合 $\{(\mathbf{x}_p, \mathbf{x}_q)_{j=1}^{l_b^{\min} \times k}\}$, 一共 $l_b^{\min} \times k$ 个, 给少数类样本对分配+1 标签。然后以同样的方法生成当前批次多数类的集合 $\{(\mathbf{x}_p, \mathbf{x}_q)_{j=1}^{l_b^{\max} \times k}\}$, 并分配-1 标签。从集合中随机选择 min_num 个少数类的样本对和 maj_num 个多数类的样本对, 在特征空间中合成新的实例, 生成新实例的式子可以写成:

$$\phi(\mathbf{x}^{pq}) = \phi(\mathbf{x}_p) + \delta^{pq}(\phi(\mathbf{x}_q) - \phi(\mathbf{x}_p)) \quad (9)$$

式中: δ^{pq} 是一个 0~0.5 的随机数, 在特征空间同样控制生成的数据点更加靠近当前批次的样本, 使得扩增的样本和当前批次中的样本的特征分布更加相似。

对当前批次的样本进行训练之前, 先用生成的样本改进目标决策函数。最后使用集成决策函数(式(4))依次对当前批次的所有实例进行预测。然而, 使用式(7)生成的新少数类实例利用通常未知的特征转换函数 $\phi(\mathbf{x})$, 所以新的合成实例 $\phi(\mathbf{x}^{pq})$ 并不能具体得到。目标域通过决策函数中支持向量的核函数计算两个特征空间中实例的内积来训练, 可以将合成实例代入目标域决策函数的核函数中计算, 其中核函数的计算分为 2 种情况:

1) \mathbf{x}_j 是合成实例, \mathbf{x} 是普通实例时, 它们在特征空间的内积为

$$k(\mathbf{x}_j^{pq}, \mathbf{x}) = \phi(\mathbf{x}_j^{pq})^T \phi(\mathbf{x}) = [\phi(\mathbf{x}_j^p) + \delta^{pq}(\phi(\mathbf{x}_j^q) - \phi(\mathbf{x}_j^p))]^T \phi(\mathbf{x}) = (1 - \delta^{pq})k(\mathbf{x}_j^p, \mathbf{x}) + \delta^{pq}k(\mathbf{x}_j^q, \mathbf{x}) \quad (10)$$

2) \mathbf{x}_j 和 \mathbf{x} 都是合成样本时, 特征空间的内积:

$$k(\mathbf{x}_j^{pq}, \mathbf{x}^{lm}) = \phi(\mathbf{x}_j^{pq})^T \phi(\mathbf{x}^{lm}) = [\phi(\mathbf{x}_j^p) + \delta^{pq}(\phi(\mathbf{x}_j^q) - \phi(\mathbf{x}_j^p))]^T \times [\phi(\mathbf{x}^l) + \delta^{lm}(\phi(\mathbf{x}^m) - \phi(\mathbf{x}^l))] = (1 - \delta^{pq})(1 - \delta^{lm})k(\mathbf{x}_j^p, \mathbf{x}^l) + (1 - \delta^{pq})\delta^{lm}k(\mathbf{x}_j^p, \mathbf{x}^m) + (1 - \delta^{lm})\delta^{pq}k(\mathbf{x}_j^q, \mathbf{x}^l) + \delta^{pq}\delta^{lm}k(\mathbf{x}_j^q, \mathbf{x}^m) \quad (11)$$

使用合成实例改进目标域决策函数, 当铰链损失大于 0 时, 将合成实例作为支持向量添加到支持向量集, 并且也能保持特征空间的可分性, 即

$$h_b^T(\mathbf{x}) = h^T(\mathbf{x}) + \sum_{j=1}^t \tau_j^{pq} y_j k(\mathbf{x}_j^{pq}, \mathbf{x}) \quad (12)$$

定理 1 在目标域的特征空间中添加合成样本同样能保证类别可分。

证明 目标域函数由支持向量组成, 可以表

示为

$$h^T(\mathbf{x}) = \sum_{i=1}^N \tau_i y_i k(\mathbf{x}, \mathbf{x}_i) \quad (13)$$

假设当前批次的样本 $\{(\mathbf{x}_j, y_j)\}_{j=1}^{l_b}$ 在目标域的特征空间是线性可分的, 从而可以得到:

$$y_j h^T(\mathbf{x}_j) = y_j \sum_{i=1}^N \tau_i y_i k(\mathbf{x}_j, \mathbf{x}_i) \geq 0, j = 1, 2, \dots, l_b \quad (14)$$

将式(9)生成少数类样本 $\phi(\mathbf{x}^{pq})$ 代入目标函数:

$$h^T(\mathbf{x}^{pq}) = \sum_{i=1}^N \tau_i y_i (\phi(\mathbf{x}^{pq})^T \phi(\mathbf{x}_i)) = (1 - \delta^{pq}) \sum_{i=1}^N \tau_i y_i k(\mathbf{x}_p, \mathbf{x}_i) + \delta^{pq} \sum_{i=1}^N \tau_i y_i k(\mathbf{x}_q, \mathbf{x}_i) = (1 - \delta^{pq})h^T(\mathbf{x}^p) + \delta^{pq}h^T(\mathbf{x}^q) \geq 0 \quad (15)$$

式中: $h^T(\mathbf{x}^p)$ 和 $h^T(\mathbf{x}^q)$ 都不小于 0, \mathbf{x}^p 和 \mathbf{x}^q 都属于少数类; $\delta^{pq} \in [0, 0.5]$ 。

所以在目标域的特征空间中生成的样本同样可以保证类别可分。每批次生成的新样本都会优化目标函数在特征空间中的超平面, 提高目标函数的性能, 从而最终提高整体函数的性能。

2.4 算法描述和复杂度分析

OTLMS_IO 和 OTLMS_FO 算法的步骤近似, 下面提供 OTLMS_FO 算法的算法描述和复杂度分析。

算法 OTLMS_FO 的算法描述

输入 源分类器 $(h^{S_1}, h^{S_2}, \dots, h^{S_n})$, 初始折衷 C , 权重折扣参数 $\beta \in (0, 1)$, 每批次扩充 min_num 个少数类和 maj_num 个多数类。

初始化: $h^T(\mathbf{x}) = \emptyset$, $v^1 = v^2 = \dots = v^n = w = 1/(n+1)$ 。

1) For 循环目标域的每个批次。

① 寻找当前批次少数类和多数类样本的 k 近邻组成种子和邻居对, 分别是 $\{(\mathbf{x}_p, \mathbf{x}_q)_{j=1}^{l_b^{\min} \times k}\}$ 和 $\{(\mathbf{x}_p, \mathbf{x}_q)_{j=1}^{l_b^{\max} \times k}\}$ 。

② 随机从少数类和多数类的样本对中选取 min_num 和 maj_num 个样本对, 根据式(9)生成新样本。

③ For 循环用于生成新样本的样本对。

a. 计算损失 ℓ 和支持向量前的参数 $\tau^{pq} = \min\{C, \ell_j/k(\mathbf{x}_j^{pq}, \mathbf{x}_j^{pq})\}$ 。

b. 损失大于 0 时, 根据式(12)更新目标域函数, 其中核函数根据式(10)和式(11)。

④ For 循环当前批次的每个实例。

a. 根据式(4)预测, 其中核函数使用式(10)和式(11)。

b. 使用式(3)更新权重。

c. 使用式(2)更新目标域。

2) 输出训练好的集成决策函数 (见式 (4))。

上述算法中, ①寻找 k 近邻的时间复杂度是 $O(3m_1m_2d + 3M_1M_2d)$, 其中 m_1 、 M_1 和 m_2 、 M_2 分别是当前批次和前面批次中的少数类和多数类, d 是样本的维数。③使用新样本改进目标函数的时间复杂度是 $O(4svd)$, s 是合成样本的总数, v 是支持向量的个数。④训练当前批次真实样本的时间复杂度是 $O(2nvd)$, 一共 n 个真实样本。在输入空间训练一个批次样本的复杂度是 $O(3m_1m_2d + 3M_1M_2d + 4svd + 2nvd)$, 整个目标域一共 N 个批次, 所以总的时间复杂度是 $O(N(3m_1m_2d + 3M_1M_2d + 4svd + 2nvd))$, 可以近似为 $O(N(m_1m_2d + M_1M_2d + svd + nvd))$ 。

3 实验结果与分析

本文对提出的算法和在线迁移学习的基线算法进行了比较, 并在多个真实数据集上进行了实验: Office-Home 数据集、Office-31 数据集和 20Newsgroups 数据集。为了获得可靠的结果, 在相同参数设置的前提下, 通过更改测试实例的到达顺序来将每个实验重复 10 次。结果表明, 本文提出的算法比基线算法获得了更好的性能。

3.1 数据集介绍

3.1.1 Office-Home 数据集

Office-Home 数据集^[18]由 4 个不同领域的图像组成: 艺术图像 (Art)、剪贴画 (Clipart)、产品图像 (Product) 和现实世界图像 (Real World), 一共大约 15 500 张图像。对于每个域, 数据集包含 65 个类别的图像。在我们的实验中, 将现实世界图像域作为目标域, 其余 3 个领域作为源域。并在目标域中随机选择一个样本数小于 50 的类别作为负类 (少数类), 选一个样本数大于 80 的类别作为正类 (多数类), 3 个源域也选取这两个类别, 然后构成一个迁移学习任务。并对原始图片进行了预处理, 每张图片都对应一个 $1 \times 10\,000$ 的向量。实验一共生成了 30 组迁移学习任务。

3.1.2 Office-31 数据集

Office-31 数据集^[19]是一个用于图像分类的迁移学习数据集。其包含 3 个领域的子集: Amazon (A)、Webcam(w)、Dslr(D), 分为 31 个类别, 共有 4652 张图片。在 Office-31 数据集中, 不仅各个领域的样本总数不同, 而且各个域内部类别分布也不平衡, 所以可以通过不平衡方法处理 Office-31 数据集, 促使迁移学习效果提升。实验中, 预处理数据集, 每个图片都是 $1 \times 10\,000$ 的向量。将 Webcam 作为目标域, 其余两个域作为源域。然后选取 Webcam 中的一个样本数多的和一个样本数少

的类别构成一组迁移学习任务, 一共生成了 16 组任务。

3.1.3 20newsgroups 数据集

20newsgroups 数据集 (<http://qwone.com/~jason/20Newsgroups/>) 由大约 20 000 个不同主题的新闻组文档组成, 这些数据被组织成 20 个不同的新闻组, 每个组对应一个不同的主题, 一共 5 个主题。例如: os、ibm、mac 和 x 是 comp 主题的新闻组, crypt、electronics、med 和 space 是 sci 主题的新闻组。其中 comp 主题的新闻组标记为正例, 而 sci 主题的新闻组标记为负例, 一共构成 4 个学习任务: os_vs_crypt、ibm_vs_electronics、mac_vs_med 和 x_vs_space。随机选择一个作为目标域, 其余作为源域, 一共构成 4 组迁移任务。

3.2 基线算法和评价指标

为了评估算法的性能, 将提出的算法和最新的几种方法进行了比较。在线被动攻击 PA 算法是一种传统的在线学习算法^[17], 采用 PA 作为基线方法, 无需知识迁移。考虑到被动攻击 PA 并非针对迁移学习问题而设计, 通过使用在整个源域中训练过的分类器初始化 PA, 来实现 PA 算法的一种变体, 称为在线迁移学习的“PAIO”。还与一种著名的在线迁移学习算法 HomOTLMS 进行了比较, 该算法从多个源域迁移知识来增强目标域的性能。所有的算法均使用 Python 语言实现和运行。

为了验证算法的可靠性, 实验结果采用分类精度和 G-mean 作为评价指标。其中 G-mean 是正例准确率与负例准确率的综合指标。当数据不平衡时, 可以评价模型表现, 若所有样本都被划分为同一个类别, G-mean 值是 0。表 1 是二分类混淆矩阵, G-mean 的计算公式为

$$G\text{-mean} = \sqrt{\frac{TN}{TN+FP} + \frac{TP}{TP+FN}} \quad (16)$$

表 1 二分类混淆矩阵
Table 1 Two-classification confusion matrix

真实情况	预测结果	
	正例	反例
正例	TP(真正例)	FN(假反例)
反例	FP(假正例)	TN(真反例)

3.3 实验结果及参数设置

3.3.1 参数设置

首先将 OTLMS_IO 和 OTLMS_FO 算法与 Office-Home、Office-31 和 20newsgroups 数据集上的所有基线算法进行比较。在 3 个数据集上, 设置

所有算法的折衷参数 C 为 5, 寻找近邻的 k 都设为 3, 并且设置多个分类器的权重折扣参数 $\beta = 0.999$ 。目标域使用高斯核, 带宽 σ 搜索范围是 $[10^{-2}, 10^2]$ 。因为在不同的数据集中一些使算法达到最优的参数的参数往往是不同的, 所以各数据集上的其他参数设置如下: 在 Office-Home 数据集中, 为了使目标域整体的类别分布相对平衡, 每批次过采样的少数类和多数类样本的个数分别是 6 和 2, 其中 $\sigma = 31.6$ 。在 Office-31 数据集中, 每批次过采样 3 个少数类和 1 个多数类样本, 高斯核带宽 $\sigma = 31.6$ 。在 20newsgroups 数据集上, OTLMS_IO 和 OTLMS_FO 算法每次过采样 40 个少数类和 10 个多数类样本, 其中高斯核函数的带宽 $\sigma = 1.12$ 。

3.3.2 Office-Home 和 Office-31 数据集上的结果
表 2 和表 3 分别列出了在 Office-Home 和 Of-

fice-31 数据集上随机选取的 4 组任务的数值结果, 并从准确率和 G-mean 指标对所有算法做出评价。其中, HomOTLMS、OTLMS_IO 和 OTLMS_FO 算法都优于 PA 和 PAIO 算法, 这表明从多个源域进行知识迁移对目标域是有帮助的。从两种评价指标可以看出, OTLMS_IO 和 OTLMS_FO 算法在应对不平衡的目标域都有着比所有基线更好的性能, 这是因为目标域整体的样本量被扩充了, 尤其是少数类样本, 增加目标分类器对少数类的偏向。其中, OTLMS_FO 算法的性能要强于 OTLMS_IO, 因为 OTLMS_FO 算法在特征空间扩增的样本使类别的分布更加近似。提出的 OTLMS_FO 算法在训练当前批次的样本之前, 会根据前面几个批次中的样本生成新样本, 因为只在几个批次中就能创建新的样本, 所以提出的算法能够保持很好的实时性。

表 2 在 Office-Home 数据集上应用不同学习算法的结果 (平均±标准差)

Table 2 Results of different learning algorithms on the Office-Home dataset (mean±standard deviations)

%

对比模型	任务							
	task3		task11		task22		task25	
	准确率	G-mean	准确率	G-mean	准确率	G-mean	准确率	G-mean
PA	78.52±0.66	72.13±0.84	73.79±0.44	69.97±0.46	60.07±0.20	50.45±0.11	67.77±0.87	63.58±1.27
PAIO	81.78±0.36	72.58±0.20	73.45±0.46	71.62±0.49	63.38±0.41	56.79±0.64	69.31±0.54	66.81±0.38
HomOTLMS	84.44±0.47	80.94±0.29	78.55±0.21	69.76±0.42	68.65±0.62	59.77±0.92	78.69±0.49	74.07±0.54
OTLMS_IO	86.74±0.62	84.83±0.98	80.55±0.41	71.80±0.83	73.31±0.87	66.23±1.02	80.31±0.71	76.06±0.76
OTLMS_FO	89.41±1.24	87.66±1.97	83.17±0.46	77.72±0.66	76.15±0.68	68.79±0.88	83.15±1.00	79.35±0.95

表 3 在 Office-31 数据集上应用不同学习算法的结果 (平均 ± 标准差)

Table 3 Results of different learning algorithms to the Office-31 dataset (mean±standard deviations)

%

对比模型	任务							
	task3		task10		task14		task15	
	准确率	G-mean	准确率	G-mean	准确率	G-mean	准确率	G-mean
PA	70.00±2.13	68.37±2.49	67.56±1.12	58.12±0.59	65.64±1.26	58.89±0.70	66.58±2.89	62.19±3.24
PAIO	75.22±1.07	73.54±0.77	88.05±0.73	83.72±0.46	75.64±1.28	67.66±0.72	71.05±0.00	65.13±0.00
HomOTLMS	80.87±2.13	74.29±2.56	90.24±0.00	81.65±0.00	82.05±0.00	60.30±0.00	71.58±2.58	65.07±3.26
OTLMS_IO	81.96±1.00	78.22±0.67	92.68±0.00	86.60±0.00	82.82±1.18	62.44±3.26	78.68±2.19	75.78±3.62
OTLMS_FO	88.04±1.75	84.38±2.34	97.80±1.71	96.13±3.04	89.74±1.62	79.69±3.61	85.53±2.94	86.88±3.32

在 Office-Home 和 Office-31 数据集上分别实验了 30 和 16 组任务, 由于受空间和可观测的局限, 在图 1 和图 2 中分别给出了 Office-Home 和 Office-31 数据集上的 PA、OTLMS、OTLMS_FO 的准确率, 而忽略了其他算法的结果。在大多数任务上, 使用多源迁移的 OTLMS_IO 和 OTLMS_FO

的性能都要优于 PA。并且在特征空间对目标域过采样的 OTLMS_FO 算法性能要更好, 证明了本文提出的算法更加适用于不平衡的目标域。

图 3 给出了 PA、HomOTLMS 和 OTLMS_FO 在 G-mean 指标上的实验结果。可以看出从多个源域迁移知识的 OTLMS_FO 和 HomOTLMS 算

法在多数任务上对少数类有着更好的表现。但是 OTLMS_FO 显然更加适合不平衡的目标域,

这种过采样的方法可以从已有数据中提取更多的信息。

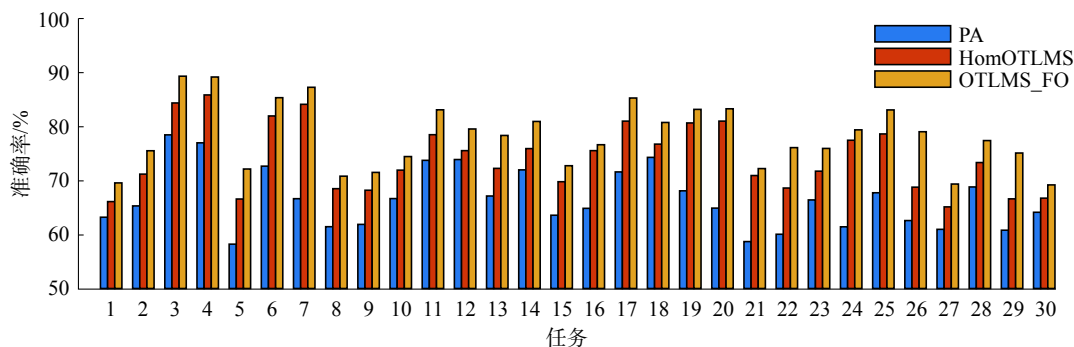


图 1 在 Office-Home 数据集的 30 组任务的准确率

Fig. 1 Accuracy of 30 sets of tasks in the Office-Home dataset

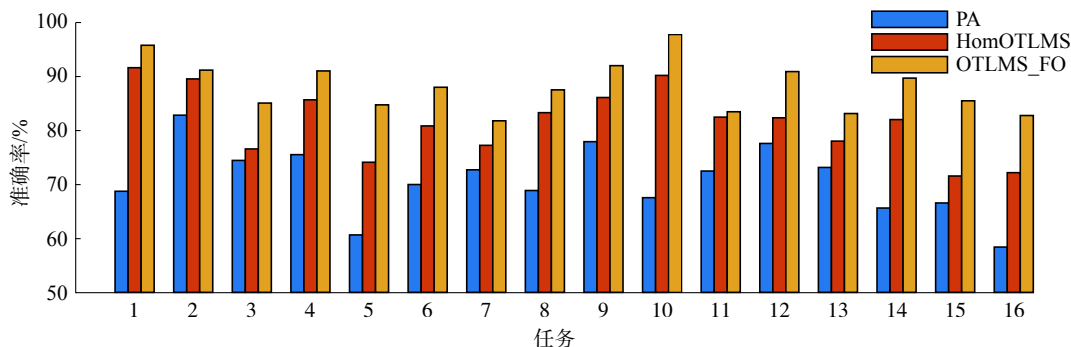


图 2 在 Office-31 数据集的 16 组任务的准确率

Fig. 2 Accuracy of 16 sets of tasks in the Office-31 dataset

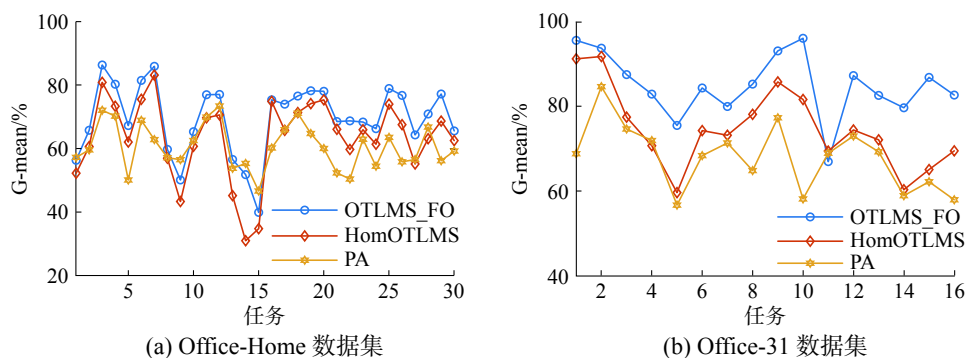


图 3 在 Office-31 和 Office-Home 数据集上各个任务的 G-mean

Fig. 3 G-mean for individual tasks on the Office-31 and Office-Home data sets

3.3.3 20newsgroups 数据集上的结果

为了更好地验证算法的性能,在 20 个新闻组的文本数据集上进行了 4 组实验。每个目标域选取 750 个样本,其中少数类占比 30%,并且每个样本的维数是 61 188,然后进行多源在线迁移。表 4 展示了文本数据集上的实验结果。与基线方法相比,我们提出的两种方法 OTLMS_IO 和 OTLMS_FO 在绝大部分任务上的性能都超过了基线。并且从实验结果可以看出 OTLMS_FO 的结果要普遍强

于 OTLMS_IO,原因是 OTLMS_FO 在核空间合成少数类,样本距离更加相似,特别是对维数较大的样本。从标准差可以看到提出的两种算法的稳定性稍弱于基线方法。因为合成样本使用了随机数 δ ,但考虑到更好的性能,牺牲一点稳定性是值得的。提出的 OTLMS_FO 算法具有很好的时效性,因为该算法只需要通过前面几个批次来扩充当前到达批次的样本,而不用在整个目标域中寻找近邻生成型样本。

表 4 在 20newsgroups 数据集上应用不同学习算法的结果 (平均 \pm 标准差)Table 4 Results of different learning algorithms to the 20newsgroups dataset (mean \pm standard deviations)

%

对比模型	任务							
	os_vs_crypt		ibm_vs_electronics		mac_vs_med		x_vs_space	
	准确率	G-mean	准确率	G-mean	准确率	G-mean	准确率	G-mean
PA	82.77 \pm 0.19	79.28 \pm 0.23	68.21 \pm 0.46	60.54 \pm 0.62	74.31 \pm 0.27	68.53 \pm 0.34	80.39 \pm 0.27	76.30 \pm 0.34
PAIO	86.72 \pm 0.19	83.98 \pm 0.23	72.79 \pm 0.34	66.41\pm0.44	78.92 \pm 0.36	74.26 \pm 0.42	85.07 \pm 0.23	81.95 \pm 0.28
HomOTLMS	88.20 \pm 0.17	84.20 \pm 0.27	76.40 \pm 0.16	50.90 \pm 0.46	79.03 \pm 0.12	68.35 \pm 0.13	86.67 \pm 0.22	79.91 \pm 0.31
OTLMS_IO	89.49 \pm 0.59	84.84 \pm 1.62	77.77 \pm 0.40	54.14 \pm 0.95	82.08 \pm 1.17	73.06 \pm 2.56	88.12 \pm 1.15	82.45 \pm 1.73
OTLMS_FO	90.29\pm0.59	87.47\pm0.93	78.44\pm0.24	56.91 \pm 0.59	85.01\pm0.60	78.92\pm0.88	89.79\pm0.45	84.89\pm0.69

3.4 时间成本

为了评估所提出算法的时间效率, 在 20newsgroups 数据集上生成多个任务对算法进行测试。实验基于 python3.7 实现, 并在具有 12 \times 2.6 GHz 的 CPU(i7-9750H) 和 16 GB 运行内存的 Windows10 专业版机器上进行。图 4 展示了 HomOTLMS、OTLMS_IO 和 OTLMS_FO 算法的平均运行时间。实验中, 对目标域样本的维数都是 61 188。从实验结果可以看出, 随着过采样样本数的增加, 两种对目标域过采样的算法所需的平均运行时间也随着增加。同时也可以发现在特征空间对目标域的样本过采样比输入空间需要花费更多的时间成本, 这是因为在特征空间中合成样本的生成需要通过多个核函数的计算才能得到。

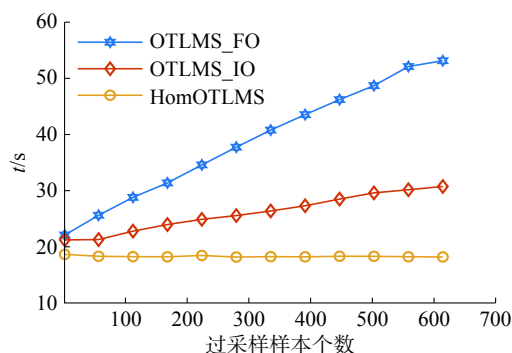


图 4 不同维数和过采样样本数的时间成本

Fig. 4 Time cost of different dimensions and oversampled sample size

4 结束语

本文提出了一种针对目标域不平衡的多源在线迁移学习算法。同时, 分别设计了在输入空间和特征空间中对目标域的样本过采样的方法。与忽略目标域类别分布的多源在线迁移学习算法相比, 提出的方法可以利用目标域已经到达的样本对当前到达的样本进行过采样, 用新生成的样本

改进目标域函数, 进而提高集成决策函数的性能, 并且时间成本的增加是可以接受的。在 3 个实际数据集上的实验结果表明, 所提出的算法与基线算法相比, 整体上实现了更好的分类性能, 也提高了对少数类预测的精度。

参考文献:

- [1] PAN S J, YANG Qiang. A survey on transfer learning[J]. *IEEE transactions on knowledge and data engineering*, 2010, 22(10): 1345–1359.
- [2] EATON E, DESJARDINS M. Selective transfer between learning tasks using task-based boosting[C]//*Proceedings of the 25th AAAI Conference on Artificial Intelligence*. San Francisco, USA, 2011.
- [3] DREDZE M, KULESZA A, CRAMMER K. Multi-domain learning by confidence-weighted parameter combination[J]. *Machine learning*, 2010, 79(1/2): 123–149.
- [4] QIAN Qi, ZHU Shenghuo, TANG Jiasheng, et al. Robust optimization over multiple domains[C]//*Proceedings of the 33rd AAAI Conference on Artificial Intelligence*. Hawaii, USA, 2019: 4739–4746.
- [5] HOFFMAN J, MOHRI M, ZHANG Ningshan. Algorithms and theory for multiple-source adaptation[C]//*Proceedings of the 32nd International Conference on Neural Information Processing Systems*. Montréal, Canada, 2018.
- [6] PENG Xingchao, BAI Qinxun, XIA Xide, et al. Moment matching for multi-source domain adaptation[C]//*Proceedings of 2019 IEEE/CVF International Conference on Computer Vision*. Seoul, South Korea: IEEE, 2019.
- [7] KANG Zhongfeng, YANG Bo, YANG Shantian, et al. Online transfer learning with multiple source domains for multi-class classification[J]. *Knowledge-based systems*, 2020, 190: 105149.
- [8] XIANG E W, PAN S J, PAN Weike, et al. Source-selection-free transfer learning[C]//*Proceedings of the 22nd In-*

- ternational Joint Conference on Artificial Intelligence. Barcelona, Spain, 2011: 2355.
- [9] GENTILE C. A new approximate maximal margin classification algorithm[J]. Journal of machine learning research, 2001, 2: 213–242.
- [10] CRAMMER K, DREDZE M, PEREIRA F. Confidence-weighted linear classification for text categorization[J]. The journal of machine learning research, 2012, 13(1): 1891–1926.
- [11] 王晓初, 包芳, 王士同, 等. 基于最小最大概率机的迁移学习分类算法 [J]. 智能系统学报, 2016, 11(1): 84–92.
- WANG Xiaochu, BAO Fang, WANG Shitong, et al. Transfer learning classification algorithms based on min-max probability machine[J]. CAAI transactions on intelligent systems, 2016, 11(1): 84–92.
- [12] ZHAO Peilin, HOI S C H. OTL: a framework of online transfer learning[C]//Proceedings of the 27th International Conference on International Conference on Machine Learning. Haifa, Israel: Omnipress, 2010.
- [13] WU Qingyao, WU Hanrui, ZHOU Xiaoming, et al. Online transfer learning with multiple homogeneous or heterogeneous sources[J]. IEEE transactions on knowledge and data engineering, 2017, 29(7): 1494–1507.
- [14] CHAWLA N V, BOWYER K W, HALL L O, et al. SMOTE: synthetic minority over-sampling technique[J]. Journal of artificial intelligence research, 2002, 16: 321–357.
- [15] 左鹏玉, 周洁, 王士同. 面对类别不平衡的增量在线序列极限学习机 [J]. 智能系统学报, 2020, 15(3): 520–527.
- ZUO Pengyu, ZHOU Jie, WANG Shitong. Incremental online sequential extreme learning machine for imbalanced data[J]. CAAI transactions on intelligent systems, 2020, 15(3): 520–527.
- [16] MATHEW J, PANG C K, LUO Ming, et al. Classification of imbalanced data by oversampling in kernel space of support vector machines[J]. IEEE transactions on neural networks and learning systems, 2017, 29(9): 4065–4076.
- [17] CRAMMER K, DEKEL O, KESHET J, et al. Online passive-aggressive algorithms[J]. The journal of machine learning research, 2006, 7: 551–585.
- [18] VENKATESWARA H, EUSEBIO J, CHAKRABORTY S, et al. Deep hashing network for unsupervised domain adaptation[C]//Proceedings of 2017 IEEE Conference on Computer Vision and Pattern Recognition. Honolulu, USA, 2017: 5385–5394.
- [19] SAENKO K, KULIS B, FRITZ M, et al. Adapting visual category models to new domains[C]//Proceedings of the 11th European Conference on Computer Vision. Heraklion, Greece, 2010: 213–226.

作者简介:



周晶雨, 硕士研究生, 主要研究方向为人工智能、模式识别。



王士同, 教授, 博士生导师, 主要研究方向为人工智能与模式识别。发表学术论文近百篇。