迭代式的深度 PU 学习与类别先验估计框架

赵昀睿1) 许倩倩2) 姜阳邦彦3).4) 黄庆明1).2).5).6)

¹⁾(中国科学院大学计算机科学与技术学院 北京 101408)
 ²⁾(中国科学院计算技术研究所智能信息处理重点实验室 北京 100190)
 ³⁾(中国科学院信息工程研究所信息安全国家重点实验室 北京 100093)
 ⁴⁾(中国科学院大学网络空间安全学院 北京 100049)
 ⁵⁾(中国科学院大学大数据挖掘与知识管理重点实验室 北京 101408)
 ⁶⁾(鹏城实验室 广东 深圳 518055)

摘 要 近年来,深度学习在诸多任务上展现了优异的性能,其一般基于海量数据并采用有监督的学习方式,依赖于完整的数据标签信息.然而在现实应用场景中,收集大量标签往往成本高昂.因此,如何利用未经充分标注的数据进行学习,成为了当下的主要挑战.二分类问题中的从正例和无标签(Positive-Unlabeled,PU)样本数据进行学习,简称 PU 学习,即为其一.当前主流的 PU 学习算法需要准确无误的类别先验知识,但实际上类别先验通常难以获得,需要估计.已有的类别先验估计算法则主要面向传统的机器学习分类器进行设计,无法直接运用在大规模数据集上,因而不利于发挥深度学习在大规模数据集上的优势.为克服以上问题,本文提出了一个基于无监督混合模型的迭代式深度 PU 学习与类别先验估计框架.它利用了深度神经网络对正例和负例给出的预测分数具有不同的分布这一特性,使用双高斯成分的混合模型近似拟合预测分数的混合分布.其中,各个高斯分量分别代表了正类和负类的条件概率分布,混合权重系数代表了类别先验.结合半监督学习中的平均教师和温度锐化技术,所提框架在类别先验未知以及数据缺失负例监督的条件下,估计类别先验的同时进行 PU 数据上的深度学习,二者相互促进. 在基准数据集 MNIST、Fashion-MNIST、CIFAR-10 和实际应用数据集 Alzheimer 上的实验结果验证了所提框架的有效性,准确率分别为 94.66%、95.16%、89.98%和 73.20%,该结果不仅超越了现有基于类别先验估计的 PU 学习算法,更可与基于真实类别先验的最前沿算法相媲美.

关键词 PU 学习;类别先验估计;半监督学习;弱监督学习;深度学习 中图法分类号 TP391 **DOI**号 10.11897/SP.J.1016.2022.02667

An Iterative Framework for Deep PU Learning and Class Prior Estimation

ZHAO Yun-Rui¹⁾ XU Qian-Qian²⁾ JIANG Yang-Bang-Yan^{3),4)} HUANG Qing-Ming^{1),2),5),6)}

¹⁾ (School of Computer Science and Technology, University of Chinese Academy of Sciences, Beijing 101408)

²⁾ (Key Laboratory of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190)
 ³⁾ (State Key Laboratory of Information Security, Institute of Information Engineering, Chinese Academy of Sciences, Beijing 100093)
 ⁴⁾ (School of Cyber Security, University of Chinese Academy of Sciences, Beijing 100049)

⁵⁾ (Key Laboratory of Big Data Mining and Knowledge Management (BDKM), University of Chinese Academy of Sciences, Beijing 101408) ⁶⁾ (Peng Cheng Laboratory, Shenzhen, Guangdong 518055)

Abstract Deep learning models have achieved superior performance with a large amount of data. Such success usually depends on complete label information in a fully-supervised training style. However, collecting all the labels can be very expensive. Consequently, weakly supervised learning,

收稿日期:2022-02-09;在线发布日期:2022-08-27.本课题研究由科技创新 2030-"新一代人工智能"重大项目(2018AAA0102000)、国家 自然科学基金项目(U21B2038,61931008,6212200758,61976202)、中央高校基本科研业务费专项资金、中科院青促会会员项目以及中国 科学院战略性先导科技专项(XDB28000000)资助. 赵昀睿,硕士研究生,中国计算机学会(CCF)学生会员,主要研究方向为弱监督学习. E-mail: zhaoyunrui20@mails.ucas.ac.cn. 许倩倩(通信作者),博士,研究员,中国计算机学会(CCF)高级会员,主要研究领域为统计机器 学习及其在多媒体领域的应用.E-mail: xuqianqian@ict.ac.cn. 姜阳邦彦,博士研究生,中国计算机学会(CCF)学生会员,主要研究方向 为机器学习与计算机视觉.黄庆明(通信作者),博士,讲席教授,中国计算机学会(CCF)会士,主要研究领域为多媒体计算、图像处理、计 算机视觉和模式识别.E-mail: qmhuang@ucas.ac.cn.

which aims at learning with incomplete, inexact, or inaccurate supervision, has attracted the machine learning community in the past decade. One of these real-world applications could be Positive-Unlabeled (PU) learning, where we have to train binary classifiers from a few positive examples with much more unlabeled data. Current state-of-the-art PU methods rely on the ground—truth class prior (i.e., the proportion of the positive samples to the unlabeled data), which is hard to obtain and needs to be estimated in practice. While previous studies of class prior estimation mainly focused on traditional machine learning models, few could deal with a relatively large-scale dataset or be applied to deep learning scenarios. To solve such problems, we propose an iterative framework for deep PU learning and class prior estimation utilizing an unsupervised mixture model in the paper. Specifically, positive and negative classes are supposed to have distinct predicted score distributions intuitively. We investigate and demonstrate our intuition and approximate the score distributions by a Gaussian Mixture Model with two components. Each component represents a relevant class-conditional distribution, and the positive weight could be approximate to the ground truth class prior. We further incorporate techniques such as mean teacher and temperature sharpening from semi-supervised learning to stabilize the whole process and boost performance. Our proposed framework could estimate the class prior and learn from PU data simultaneously, achieving well-matched performance with other PU competitors based on the ground-truth prior. Experiments on three benchmark datasets (i. e., MNIST, Fashion-MNIST, and CIFAR-10) and one practical application (i. e., Alzheimer) validate the effectiveness of our framework. Our algorithm reaches the accuracy of 94.66%, 95.16%, 89.98%, and 73. 20%, respectively, with limited labels and unknown class prior.

Keywords positive-unlabeled learning; class prior estimation; semi-supervised learning; weakly supervised learning; deep learning

1 引 言

近年来,随着互联网与信息技术的发展,人类进 入了大数据时代.机器学习,特别是基于海量数据的 深度学习,受到了广泛关注并取得了突破性进展,在 诸多任务上的表现^[1-3]达到甚至超越了人类水平.然 而,深度学习算法通常采用有监督的训练方式,其取 得的优异性能依赖于海量标注数据,尤其离不开完 整类别标签信息的指导.当数据中的标签数目较少 时,深度监督学习算法容易产生过拟合现象,难以发 据隐含在数据背后的规律,模型的泛化性能因之大 打折扣.另一方面,对于许多实际应用场景而言,数 据标注代价高昂^[4].在人力和物力有限的条件下,人 们仅能获取少量数据的标签.鉴于此,对标注数据依 赖更少的学习方式成为当下热门,半监督学习即为 其中之一.

半监督学习最早由 Merz 等人^[5]提出并应用在 分类问题上,其基于少量的标注数据,旨在挖掘无标 签数据中蕴含的信息,用以辅助提升模型的性能.尽

管半监督学习在一定程度上减轻了对标注数据的依 赖,但是一般的半监督学习算法要求标注数据中同 时包含各个类别,限制了其普适性.对于部分二分类 的实际应用问题,除了数据标注任务本身较为困难 之外,还存在负类标签难以获取的问题.此时,常规 的半监督学习算法将因缺失负例监督而失效.例如 在阿尔兹海默症(Alzheimer's Disease, AD)的诊断 问题^[6]中,病理特征明显的患者容易被临床确诊,可 以看作有标签的正例;而早期 AD 病患与正常人无 殊,难以被临床确诊.也就是说,未被确诊的人仍旧 存在患病可能,属于无标签样本.类似情况还存在于 恶意 URL 检测^[7]、虚假评论检测^[8]、冷冻电镜的粒 子拾取^[9]等任务上.由此可见,在二分类半监督学习 中,进一步考虑仅利用正例和无标签样本进行学习 (Positive-Unlabeled Learning, PU), 简称为 PU 学 习,具有重要的研究价值,近年来受到了机器学习社 区的广泛关注[10].

如图 1 所示,相较于一般的二分类半监督学习, PU 学习的标注数据缺失负例信息,因而更具挑战 性.PU 学习算法根据其对无标签数据的处理方式,

2669

可分为两大流派,分别是两阶段式的 PU 学习和基 于风险重构的 PU 学习.其中,两阶段式的 PU 学习 符合思维惯性,该类方法^[11-13]首先找出可靠的负 例,从而将 PU 问题转换为常规的二分类半监督学 习问题后,再运用(半)监督学习算法.然而,此类方 法依赖于不同类别数据分布的严格可分性和较为理 想的标签分配机制.因此,基于风险重构的 PU 学 习^[6,14-16]成为主流方法,其将二分类监督学习的风 险重构为关于正例和无标签样本的风险,从而无需 事先在无标签数据中识别出可靠的负例.nnPU^[16] 成功将基于风险重构的 PU 学习扩展到了深度场 景,Self-PU^[6]则结合自步(Self-Paced)学习的训练 策略和深度模型的自校准(Self-Calibration)与自蒸 馏(Self-Distillation)技术,在 PU 学习的基准数据 集上取得了最先进的分类性能.



图 1 二分类半监督学习与 PU 学习的联系与区别

尽管基于风险重构的 PU 学习取得了巨大成 功,但该类方法严重依赖于真实的类别先验知识,否 则将无法构造 PU 学习的无偏风险估计器.对于多 数实际问题而言,类别先验往往未知,这大大降低了 此类 PU 学习算法的实用价值.那么,可否事先随意 设置一个类别先验值,再应用此类算法呢?答案是否 定的.实验表明(如图 2 所示),不准确的类别先验预 设值将会严重损害算法性能.基于风险重构的 PU 学习与类别先验的依赖关系,驱动了类别先验估计 的相关研究^[17-20].但这些方法大多仅能在小规模数



图 2 基于风险重构的 PU算法 nnPU¹⁶和 Self-PU⁶在 CIFAR-10上的准确率随类别先验预设值的变化曲线(黑色虚 线为类别先验真实值(0.4).总体来看,类别先验预设 值越偏离真实值,算法的准确率越差)

据集和传统机器学习方法上达到较优水平,无法直接应用在基于海量数据的深度学习方法上,需要额外进行主成分分析(Principal Components Analysis, PCA)和随机采样等操作^[19].

在二分类深度神经网络过拟合前,其预测分数的分布关于正、负类别存在差异,它们的混合分布呈现双峰状^[22].基于上述观察,本文提出了一种迭代式的深度 PU 学习与类别先验估计框架,其由类别先验估计模块、PU 学习模块和迭代过程稳定模块构成.其中,类别先验估计模块基于双成分的高斯混合模型(Gaussian Mixture Model,GMM),对深度模型的预测分数进行无监督式建模,以估计基于风险重构的 PU 学习所必需的类别先验,且不受 PU 数据缺失负类标签的限制.此外,运用 EM 算法求解GMM 所需的时间复杂度、空间复杂度均为 O(N),占用的计算资源少.因此,所提方法可广泛应用于各种规模的数据集,进而可以有效应对深度场景下 PU 数据的类别先验未知问题.

综上所述,本文的主要贡献总结如下:

(1)本文提出了一种迭代式的深度 PU 学习与 类别先验估计框架.该框架利用了深度 PU 学习与 类别先验估计的相互促进关系,能够有效应对深度 场景下 PU 数据的类别先验未知问题.

(2)所提框架的核心思想在于使用 GMM 对深 度神经网络输出的预测分数建模,进而估计类别 先验.在深度神经网络其对训练数据过拟合之前, 正、负类的预测分数条件分布可分别近似为两条钟 形曲线,其混合分布可用 GMM 拟合.

(3) 在所提框架中,结合使用半监督学习常用的 平均教师(Mean Teacher)^[23]和温度锐化(Temperature Sharpening)技术^[24],将二者作为迭代过程稳 定模块,以平稳类别先验的估计过程,进而提升模型 的分类性能.

(4)本文在三个基准数据集 MNIST、F-MNIST 和 CIFAR-10 以及实际应用数据集 Alzheimer 上进 行了详细的对比实验、消融实验和超参数敏感度分 析,验证了所提框架的有效性及其设计的合理性.

2 相关工作

2.1 PU 学习

PU 学习方法主要可以分为两大流派,分别是 基于风险重构的 PU 学习和两阶段式的 PU 学习. 当前主流方法基于风险重构的思想,在优化过程中 构造了 PU 数据关于二分类监督学习的风险无偏估 计器,其基本形式如下:

$$R = \pi_{P} R_{P}^{+} + R_{U}^{-} - \pi_{P} R_{P}^{-} \tag{1}$$

其中, π_P =Pr(y=1)表示(正)类别先验,简称类别 先验, R_{P}^{+} 表示将正类样本预测成正类的期望风险, 同理 R_p 表示正类样本预测成负类的期望风险, R_u 则表示将无标签样本预测成负类的期望风险.Du Plessis 等人于 2014 年首先提出了 PU 学习的无偏风 险估计(unbiased PU Risk Estimator, uPU)^[14],并 进一步指出训练过程中使用满足特定条件的代理损 失函数(Surrogate Loss Function)可以有效降低计 算成本[15]. 自此以后, 机器学习社区涌现出许多致 力于增强或改进 uPU 的工作. 深度神经网络以其强 大的拟合能力著称,那么直接优化 PU 学习的无偏 经验风险可能导致 R_{N} 的无偏估计值小于 0. 为了防 止深度神经网络过度优化将无标签样本预测成负类 的风险,即 R_U^- ,Kiryo 等人在 uPU 中关于 R_N^- 的无 偏估计器上增加了非负约束,由此提出了 PU 学习的 非负风险估计器(non-negative PU Risk Estimator, nnPU)^[16].此外,Self-PU^[6]在nnPU的基础上引入 了自监督学习(Self-Supervision Learning)任务,通入 过导师网络(Mentor Net)^[25]和平均教师(Mean Teacher)^[23]的架构,实现深度模型的自校准和自蒸 馏,并结合自步学习的训练策略,在 PU 学习的基准 数据集上取得了最优性能. PUSB^[26]提出了次序不 变假设(Invariance of Order),即样本被标注的概率 与该样本是正例的概率正相关.并基于该假设,从理 论上证明了当 uPU 使用对数损失函数时,其贝叶斯 最优解满足预测分数与正类条件概率正相关. aPU^[27] 进一步拓展到训练数据与测试数据的正类条件分布 不一致的情况. ImbPU^[28]基于风险重构的思想构造 了等价于对少数类进行过采样的优化目标,以应对 PU 数据的类别不平衡问题.

PU学习的另一大流派是两阶段式的 PU 学习, 这类方法符合人们的惯性思维,首先识别出一定数 目的可靠负例后,便可在第二阶段运用(半)监督学 习算法.由此可见,两阶段式 PU 学习的研究重心和 主要区别在于可靠负例的识别机制.文献[29]从被 标注正例中随机选取一定数目的间谍(Spy),将之 加入无标签样本集合中;这样一来,可靠的负例即为 后验概率小于任一间谍的无标签样本.文献[30]基 于图进行标签传播,利用高斯核构建邻接矩阵,欧式 空间上特征相似的样本将被划分至同一类别.文献 [31]使用朴素贝叶斯对样本被标注的条件概率进行 建模,再依据上述条件概率熵的大小,筛选出可靠的 正例或负例. 文献 [32] 使用 k-means 进行聚类, 可靠 负例即来自距离被标注正例最远的样本簇. RP^[33]通 过计算各个样本的模型预测置信度得分,以筛选出 高置信度的样本,它们的类别即为模型的预测结果. 文献[34]从生成学习(Generative Learning)的角度, 在生成对抗网络(Generative Adversarial Network, GAN)的框架下联合使用了正例生成器和负例生成 器,以及正例判别器、无标签样本判别器和负例鉴 别器.其中,无标签样本判别器用以保证生成数据的 合理性.KLDCE^[35]首先将 PU 学习问题转化为标 签噪声问题,然后通过对有噪声的负例集合进行质 心估计来消减标签噪声的副作用.基于实例依赖 (Instance Dependent)的 PU 学习^[36]将类别标签视 为由样本输入特征确定的隐变量,并借助 EM 的算 法框架学习分类器. 两阶段式的 PU 学习依赖于所 选样本的数量及其可靠性,因此要求不同类别的数 据分布严格可分和较为理想的标签分配机制.

2.2 类别先验估计

绝大多数的 PU 学习方法依赖于类别先验知 识,即π_P.然而,在实际应用中,类别先验往往未知. 为此,针对 PU 数据的类别先验估计成为重要的研 究课题.其中,PE^[17]基于部分匹配准则,最小化无标 签样本分布和被标注正例分布之间的 Pearson 距 离.其中,类别先验即为最优匹配时有标签正例分布 的标度因子、KM1 和 KM2^[18]提出了一种基于核嵌 入的类别先验估计方法,将正类分布从总体分布中 分离出来.类别先验即为最小负例比例低于给定阈 值的最大值. TIcE^[19]指出当完全随机选择(Selected Completely At Random, SCAR) 假设成立时,给定 任一数据子域,在该子域中样本被标记的频率是正 类样本标签频率的下确界.通过决策树归纳,该方法 旨在找到令样本被标记频率最大的子域,以逼近真 实的标签频率,进而估计类别先验.CDMM^[20]基于 逻辑回归,选取其参数对数似然的上界凹函数,并通 过循环坐标下降(Cyclic Coordinate Descent, CD)和 MM 优化算法(Minimization-Majorization, MM)近 似求解上述上界凹函数,从而估计正类标签频度,进 而估计类别先验.

已有的类别先验估计算法^[17-20]往往采用非深 度的传统机器学习算法,尽管在小规模数据集上成 效卓著,但无法直接应用于深度学习场景下的大规 模数据集.此类算法的核心思想是混合比例估计,即 无标签样本的概率分布是正类样本分布和负类样本 分布的线性组合,其中的混合系数表示类别先验.当 不可约(Irreducibility)假设^[21]成立时,正类样本的 支撑集(Support)不能包含于负类样本的支撑集中, 那么类别先验可以由混合概率密度与正类样本概率 密度比值的最小值估计得到.由此可见,基于混合比 例估计的类别先验估计算法具备较为可靠的理论保 障,并且在低维、小规模数据上得到了广泛验证.然 而,由于"维数灾难"和计算资源有限等问题,高维、 大规模数据上的概率密度估计十分困难,若是简单 地对数据进行降维/采样处理,则可能会破坏不可约 假设,最终致使算法失效.考虑到深度模型强大的特 征提取能力以及模型训练与类别先验估计之间存在 相互促进的关系,所提框架在模型输出的一维预测 分数上进行类别先验估计,有效降低了计算复杂度. 另外,尽管预测分数浓缩了输入样本的类别信息,但 预测分数的分布难以满足不可约假设.为此,所提框 架使用高斯混合模型对预测分数建模,而高斯混合 模型可以有效应对不可约假设不成立的情况.最后, 根据半监督学习的平滑假设和聚类假设,所提框架 进一步引入了平均教师和温度锐化技术,有效改善 了方法的有效性和稳定性.

2.3 半监督学习

半监督学习算法常利用平滑假设、聚类假设和 流形假设三大基本假设[37-38]来建立无标签数据的 学习器.近期的深度半监督学习工作主要从平滑假 设和聚类假设着手. 平滑假设要求数据稠密区域内 距离相近的样本具有相似的类别标签.基于平滑假 设的半监督学习方法在优化目标中添加无监督分量 对无标签数据进行正则化.这类方法简单有效,通过 对数据或模型参数进行扰动,利用无标签数据引导 网络牛成信息更丰富的抽象表示. Rasmus 等人^[39] 首先提出梯形网络,它将前馈网络作为自动编码器 (Auto-Encoder)的编码器,添加解码器,加入重构损 失;伪标签集成^[40]、Ⅱ-Model^[41]等对神经网络模型本 身进行扰动,计算不同网络输出之间的一致性损失. 平均教师[23] 对模型的历史参数进行时序上的集成, 取得了明显的性能提升.聚类假设要求分类决策边界 附近的数据分布稀疏.根据聚类假设,模型在无标签 数据上的输出熵应该尽可能小.半监督深度学习的先 进方法 MixMatch^[42]、UDA^[24]和 ReMixMatch^[43]均 使用了温度锐化技术,间接降低了模型的输出熵.流 形假设指高维数据可以被映射至的低维流形中,那 么局部邻域内的样本应当具有相似的类别标签[44]. 现有的深度半监督模型甚至可在部分基准数据集上 用较少的标注数据达到与使用完整标注数据训练相 近的性能^[45].

除了半监督学习之外,主流的无标签数据处理 方法还包括主动学习(Active Learning)、迁移学习 (Transfer Learning)等.其中,主动学习严重依赖于 专家知识,需要通过专家与模型之间的频繁交互额 外进行人工标注;迁移学习则依赖于迁移前后任务 或数据分布之间的相关性,且要求标注数据同时包 含正负例.相较于上述两者,PU学习作为二分类半 监督学习的特例,无需额外的人工标注,不依赖于其 它任务或模型,在诸如 AD 诊断问题等负类标签难 以获取的特定应用场景下更具潜力.

3 方 法

3.1 问题定义

C

PU学习是二分类问题中的一类特殊情况,因此,首先对二分类问题进行回顾.在二分类问题中,输入空间是一个 d 维实数向量子空间,即 $\mathcal{X} \subseteq \mathbb{R}^{d}$;标签空间 $\mathcal{Y} = \{0,1\}$,其中 0 表示负类,1 表示正类.设 p(x,y)为真实数据分布(\mathcal{X},\mathcal{Y})的联合概率密度,p(x)为关于某一样本其输入向量x的边缘分布, $y \in \mathcal{Y}$ 为其类别标签.那么,二分类问题中的正例集合和负例集合可以分别表示为

$$\boldsymbol{X}_{P} = \{\boldsymbol{x}\}^{n_{P}} \sim p_{P}(\boldsymbol{x})$$
(2)

$$\boldsymbol{X}_{N} = \{\boldsymbol{x}\}^{n_{N}} \sim p_{N}(\boldsymbol{x}) \tag{3}$$

其中, n_p 和 n_x 分别为正例和负例集合的基数,即对应 集合的样本个数; $p_P(\mathbf{x}) = \Pr(\mathbf{x} | y=1)$ 和 $p_N(\mathbf{x}) =$ $\Pr(\mathbf{x} | y=0)$ 分别表示正类和负类的条件概率.基于 上述定义,全体样本的输入向量集合 $\mathbf{X} = \mathbf{X}_P \cup \mathbf{X}_N$ 可 以表示为

$$\boldsymbol{X} = \{\boldsymbol{x}\}^n \sim \boldsymbol{p}(\boldsymbol{x}) \tag{4}$$

$$p(\mathbf{x}) = \pi_P \cdot p_P(\mathbf{x}) + \pi_N \cdot p_N(\mathbf{x}) \tag{5}$$

其中,样本总数 $n=n_P+n_N;\pi_P=\Pr(y=1)$ 表示(正) 类别先验, $\pi_N=1-\pi_P$ 表示负类别先验.

给定深度神经网络结构 f,二分类监督学习的目标是学习模型的参数 Θ ,使得分类器 $f(\cdot;\Theta):\mathbb{R}^d \to \mathbb{R}$ 在真实数据分布(\mathcal{X}, \mathcal{Y})上取得最小的预测误差期望,即

$$\boldsymbol{\Theta}^{*} = \underset{\boldsymbol{\Theta}}{\operatorname{arg\,min}} R_{PN},$$

$$R_{PN} = \underset{\boldsymbol{x} \sim \rho(\boldsymbol{x})}{\mathbb{E}} \left[\ell_{0-1}(f(\boldsymbol{x}; \boldsymbol{\Theta}), \boldsymbol{y}) \right] \qquad (6)$$

其中, $\ell_{0-1}(\hat{y}, y) = \frac{1 - sign(\hat{y} \cdot y - 1/2)}{2}$ 表示 0-1 损

失函数. 然而, 0-1 损失函数不可导, 致使 R_{PN} 难以

优化.为此,通过引入代理损失函数ℓ,近似替代式(6) 中的优化目标:

$$R_{PN}^{\ell} = \mathbb{E}\left[\ell(f(\boldsymbol{x};\boldsymbol{\Theta}), \boldsymbol{y})\right]$$
(7)

相较于二分类监督学习,PU 学习的数据集 $X_{PU} = X_L \cup X_U Q 由少量正例和大量的无标签数据构成,其中 <math>X_L$ 表示有标签的正例集合, X_U 表示无标签样本集合.基于 SCAR 假设^[10],即给定样本所属类别 y,该样本被标注的概率与其输入特征 x 无关,以及单一训练集(Single-Training-Set)设置^[10],被标注的 正样本服从真实的正类条件分布,而全体样本服从

$$\boldsymbol{X}_{L} = \{\boldsymbol{x}\}^{n_{L}} \sim p_{P}(\boldsymbol{x}) \tag{8}$$

$$\boldsymbol{X}_{\mathrm{PU}} = \{\boldsymbol{x}\}^n \sim p(\boldsymbol{x}) \tag{9}$$

3.2 流程框架

本文提出了一种迭代式的深度 PU 学习与类别 先验估计框架,以克服深度场景下类别先验未知问 题.所提框架由三大模块组成,分别是基于非负 PU 损失的 PU 学习模块、基于 GMM 的类别先验估计 模块以及基于平均教师与温度锐化的迭代过程稳定 模块,其流程架构如图 3 所示.



图 3 迭代式的深度 PU 学习与类别先验估计框架

首先,PU 数据的全体样本 X_{PU}被同时输入至两 个网络结构相同而参数不同的学生模型 $f(\cdot; \Theta)$ 和 教师模型 $f(\cdot; \Theta')$ 中,教师模型的参数 Θ' 由学生 模型的历史参数进行时序上的指数窗口滑动平均 (Exponential Moving Average, EMA)运算得到;由 此得到学生模型和教师模型各自对 X_{PU}的预测分数 $S = f(X_{PU}; \boldsymbol{\Theta})$ 和 $S' = f(X_{PU}; \boldsymbol{\Theta}');$ 之后,类别先验估 计模块对 S'使用 GMM 进行无监督式建模,并通过 EM 算法迭代多步直至收敛,以求解 GMM 的参数,由 此得到类别先验估计值 $\hat{\pi}_{P}$; PU 学习模块再借助 $\hat{\pi}_{P}$ 计算非负 PU 损失 $\hat{R}^{\ell_{\text{sigmoid}}}$;最后结合迭代过程稳定 模块中教师模型和学生模型预测分数之间的一致性 损失 L_{con},以及温度锐化损失 L_{sharpen},进而得到了框 架的总体优化目标 L,并使用梯度反向传播算法更 新 Θ. 上述步骤迭代反复执行,直至 L 收敛或达到 预设迭代次数,接下来将依次对所提框架的各个模 块进行详细介绍.

3.3 PU 学习模块

PU 学习模块负责从训练数据中学习其潜在的 分类规则. 它包含一个基于深度神经网络 *f* 的分类 器,并定义了 *f* 在 PU 数据上进行训练的优化目标. 式(7)给出了二分类监督学习的优化目标 *R*^ℓ_{PN},其可 $R_{PN}^{\ell, +} = \underset{x \in \mathbf{X}_{L}}{\mathbb{E}} \left[\ell(f(\mathbf{x}; \mathbf{\Theta}), 1) \right]$ 表示在正例集合中 将样本预测为正类的期望损失, $R_{N}^{\ell, -} = \underset{x \sim P_{N}(\mathbf{x})}{\mathbb{E}} \left[\ell(f(\mathbf{x}; \mathbf{\Theta}), 1) \right]$ 表示在正例集合中 将样本预测为正类的期望损失, $R_{N}^{\ell, -} = \underset{x \sim P_{N}(\mathbf{x})}{\mathbb{E}} \left[\ell(f(\mathbf{x}; \mathbf{\Theta}), 0) \right]$ 表示在负例集合中将样本预测为负类的期望 损失. 在 SCAR 假设下, 式(8) 成立, 也就是说, \mathbf{X}_{L} 与 \mathbf{X}_{P} 具有相同的类条件分布 $p_{P}(\mathbf{x})$. 因此, 根据经验风 险最小化原则(Empirical Risk Minimization, ERM), $\hat{R}_{L}^{\ell, +} = \underset{x \in \mathbf{X}_{L}}{\mathbb{E}} \left[\ell(f(\mathbf{x}; \mathbf{\Theta}), 1) \right]$ 是关于 $R_{P}^{\ell, +}$ 的无偏估计. 根据式(9), \mathbf{X}_{PU} 与 \mathbf{X} 同分布, 于是以下等式关系成立:

重构为不同类别期望损失的线性组合:

$$R_{PU}^{\ell,-} = \underset{\boldsymbol{x} \sim p(\boldsymbol{x})}{\mathbb{E}} \left[\ell(f(\boldsymbol{x};\boldsymbol{\Theta}), -1) \right]$$
$$= \pi_{P} R_{P}^{\ell,-} + \pi_{N} R_{N}^{\ell,-}$$
(11)

由此可见, $\pi_N R_N^{\ell,-} = R_{PU}^{\ell,-} - \pi_P R_P^{\ell,-}$ 可由 $\hat{R}_{PU}^{\ell,-} = \underset{x \in \mathbf{X}_{PU}}{\mathbb{E}} \left[\ell(f(\mathbf{x}; \boldsymbol{\Theta}), 0) \right]$ 句 弟 $\hat{R}_L^{\ell,-} = \underset{x \in \mathbf{X}_L}{\mathbb{E}} \left[\ell(f(\mathbf{x}; \boldsymbol{\Theta}), 0) \right]$ 间接估计^[14]. 综上所述,可以得到 R_{PN}^{ℓ} 关于PU数据的无偏估计,即:

$$\hat{R}_{\rm PU}^{\ell} = \pi_P \hat{R}_L^{\ell,+} + \hat{R}_{\rm PU}^{\ell,-} - \pi_P \hat{R}_L^{\ell,-}$$
(12)

若代理损失函数 $\ell(\bullet, \bullet)$ 满足 $\ell(\bullet, 1) + \ell(\bullet, 0) = 1, 那$ 么 $\hat{R}_{L}^{\ell, -} = 1 - \hat{R}_{L}^{\ell, +}$ 成立,将之代入式(12),有:

$$\hat{R}_{\rm PU}^{\ell} = \pi_P \hat{R}_L^{\ell,+} + \hat{R}_{\rm PU}^{\ell,-} - \pi_P + \pi_P \hat{R}_L^{\ell,+}$$
(13)

其中, \hat{R}_{PU}^{ℓ} 仅需计算 $\hat{R}_{L}^{\ell,+}$ 和 $\hat{R}_{PU}^{\ell,-}$,相较于式(12)不需 要计算 $\hat{R}_{L}^{\ell,-}$,因而有效降低了计算复杂度^[15].

由于深度神经网络具有强大的拟合能力,在实际 训练过程中,模型过度优化 $\hat{R}_{PU}^{\ell,-}$,以至于关于 $R_N^{\ell,-}$ 的 无偏估计小于 $0^{[16]}$,与真实情况下 $R_N^{\ell,-} \ge 0$ 相悖.为此,需要向式(13)引入非负约束:

 $\hat{R}_{nnPU}^{\ell} = \pi_P \hat{R}_L^{\ell,+} + \max(0, \hat{R}_{PU}^{\ell,-} + \pi_P \hat{R}_L^{\ell,+} - \pi_P) (14)$ 所提框架的代理损失函数选取 PU 学习中常用的 sigmoid 损失,即:

$$\ell_{\text{sigmoid}}(\hat{y}, y) = 1 - \frac{1}{1 + \exp[\hat{y} \cdot (2y - 1)]} \quad (15)$$

显然, $\ell_{sigmoid}$ 满足 $\ell_{sigmoid}(\bullet, 1) + \ell_{sigmoid}(\bullet, 0) = 1$.因此,所提框架中 PU 学习模块的优化目标为

$$\hat{R}_{nnPU}^{\ell_{sigmoid}} = \pi_{P} \hat{R}_{L}^{\ell_{sigmoid},+} + \max(\Theta, \hat{R}_{PU}^{\ell_{sigmoid},-} + \pi_{P} \hat{R}_{L}^{\ell_{sigmoid},+} - \pi_{P})$$
(16)

由上式可知,类别先验 π_P是构造关于 PU 数据 无偏风险估计器的关键. 然而,实际应用场景的 π_P 往往未知,这将导致 PU 学习模块无法正常运作. 为 此,本文设计了类别先验估计模块.

3.4 类别先验估计模块

类别先验估计模块利用不同类别分布的可分性 和深度神经网络的抽象表征能力,通过由两个高斯 成分的 GMM 对全体样本的预测分数进行无监督式 地建模.对于 GMM 中均值较大的高斯成分,它的权 重即为 PU 学习模块必需的类别先验.其中,样本 *x* 关于模型 ❷ 的预测分数是对条件概率 Pr(*y*=1|*x*) 的建模,其定义如下:

s=sigmoid[f(x; ④)]∈[0,1] (17) 混合模型是一类无监督建模方法^[46-48],其中, GMM^[48]的应用最为广泛.使用双高斯成分的 GMM 对预测分数建模的依据基于以下观察.在深度神经 网络发生过拟合现象之前,正类的预测分数与负类 的预测分数呈现出不同的分布.如图 4 所示,正类的 预测分数集中分布在分值较高的区间,而负类的预



测分数集中分布在分值较低的区间,两者分别构成 了两条中间高、两端低的钟形曲线;预测分数的混合 分布由上述两条钟形曲线加权叠加而成.由此可见, 不同类别的预测分数分布可以近似为不同均值和方 差的高斯分布;预测分数的混合分布可近似为上述 高斯分布的线性组合.因此,双高斯成分的 GMM 对预测分数 S 的建模如下:

$$S = \{s\}^{n} \stackrel{i.i.d}{\sim} \pi_{P} \operatorname{Pr}(s \mid y=1) + \pi_{N} \operatorname{Pr}(s \mid y=0)$$
$$= \pi_{P} \mathcal{N}_{P}(s) + \pi_{N} \mathcal{N}_{N}(s)$$
(18)

 $\mathcal{N}_{P}(s) = \mathcal{N}(s | \mu_{P}, \sigma_{P}), \mathcal{N}_{N}(s) = \mathcal{N}(s | \mu_{N}, \sigma_{N}) (19)$ 其中,类别标签 y 是隐变量; π_{P} 既是类别先验,同时 也是 GMM 中 $\mathcal{N}_{P}(s)$ 的混合系数;同理, $\pi_{N} = 1 - \pi_{P}$ 既 是负类别先验,同时也是 $\mathcal{N}_{N}(s)$ 的混合系数; $\mathcal{N}_{P}(\cdot)$, $\mathcal{N}_{N}(\cdot)$ 分别表示正类和负类预测分数各自所服从的 高斯分布, $\mathcal{N}(s | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(s-\mu)^{2}}{2\sigma^{2}}\right), \mu, \sigma$ 分别表示其均值和标准差; 因为正类的预测分数整

体上应当大于负类,所以有 μ_N<μ_P. 接下来需要对 GMM 中的类别先验 π_P进行求 解.考虑到 GMM 中包含隐变量 y,选择最大期望算 法(Expectation Maximization,EM)^[49]进行求解. 3.4.1 EM 算法

EM算法^[49]主要通过期望(Expectation,E)步和最大化(Maximization,M)步反复迭代直至似然函数收敛至局部最优解. 令双高斯成分的所有参数 $\boldsymbol{\Phi} = (\mu_P, \sigma_P, \mu_N, \sigma_N), 根据极大似然优化准则,首先为其构造对数似然函数:$

 $Q_{p}^{(t_{in})}(s) = \Pr(y=1|s; \boldsymbol{\phi}^{(t_{in})})$ $= \frac{\pi_{p}^{(t_{in})} \mathcal{N}_{p}^{(t_{in})}(s)}{\pi_{p}^{(t_{in})} \mathcal{N}_{p}^{(t_{in})}(s) + (1-\pi_{N}^{(t_{in})}) \mathcal{N}_{N}^{(t_{in})}(s)} (21)$ $Q_{N}^{(t_{in})}(s) = 1 - Q_{p}^{(t_{in})}(s) \qquad (22)$

进而得到关于
$$\mathcal{L}(S|\boldsymbol{\Phi},\pi_{P},\pi_{N})$$
的一阶近似:

$$\mathcal{L}^{(t_{\text{in}})}(\boldsymbol{S}|\boldsymbol{\Phi}) = \sum_{s \in \boldsymbol{S}} \left[Q_{p}^{(t_{\text{in}})}(s) \log(\mathcal{N}_{p}^{(t_{\text{in}})}(s)) + Q_{N}^{(t_{\text{in}})}(s) \log(\mathcal{N}_{N}^{(t_{\text{in}})}(s)) \right]$$
(23)

图 4 CIFAR-10 上 GMM 拟合预测分数混合分布实例图

接下来,EM 算法通过 M 步最大化 $\mathcal{L}^{(t_{in})}(S \mid \boldsymbol{\Phi})$,以 得到第 t_{in} +1次迭代后的参数:

$$\boldsymbol{\Phi}^{(t_{\text{in}}+1)} = \arg \max \mathcal{L}^{(t_{\text{in}})} \left(\boldsymbol{S} \middle| \boldsymbol{\Phi} \right)$$
(24)

$$\pi_{p}^{(t_{\text{in}}+1)} = \mathbb{E}_{\boldsymbol{s} \in \boldsymbol{s}} \left[\boldsymbol{Q}_{p}^{(t_{\text{in}})} \right] \tag{25}$$

$$\pi_{\nu in}^{(t_{in}+1)} = 1 - \pi_{\nu}^{(t_{in}+1)}$$
(26)

其中,式(24)关于**Φ**的参数更新方程可转换为以下 形式^[50]:

$$\mu_{p}^{(t_{in}+1)} = \frac{1}{\pi_{p}^{(t_{in}+1)}} \mathop{\mathbb{E}}_{s \in S} \left[s Q_{p}^{(t_{in})} \right],$$

$$\sigma_{p}^{(t_{in}+1)} = \frac{1}{\pi_{p}^{(t_{in}+1)}} \mathop{\mathbb{E}}_{s \in S} \left[(s - \mu_{p}^{(t_{in}+1)})^{2} Q_{p}^{(t_{in})} \right],$$

$$\mu_{N}^{(t_{in}+1)} = \frac{1}{\pi_{N}^{(t_{in}+1)}} \mathop{\mathbb{E}}_{s \in S} \left[s Q_{N}^{(t_{in})} \right],$$

$$\sigma_{N}^{(t_{in}+1)} = \frac{1}{\pi_{N}^{(t_{in}+1)}} \mathop{\mathbb{E}}_{s \in S} \left[(s - \mu_{N}^{(t_{in}+1)})^{2} Q_{N}^{(t_{in})} \right]$$
(27)

EM 算法求解 GMM 的伪代码如算法 1 所描述, E 步和 M 步反复迭代执行, 收敛后可得类别先验 π_p 的估计值 $\hat{\pi}_p$.

算法1. 类别先验估计.

输入:

 $S: PU 数据的预测分数, S = [s_1; s_2; \dots; s_n]$

Φ⁽⁰⁾:双高斯成分的初始参数

 $\pi_P^{(0)}, \pi_N^{(0)}$: GMM 的初始混合系数

输出:

 $\hat{\pi}_{P}$:类别先验估计值

过程:

#初始化

```
t_{\rm in} = 0
REPEAT
```

```
#E 步
```

```
根据式(21)和(22)计算 Q_{p}^{(t_{in})}(s), Q_{N^{in}}^{(t_{in})}(s).

#M 步

根据式(27)计算 \boldsymbol{\Phi}^{(t_{in}+1)}.

根据式(25)和(26)计算 \pi_{p}^{(t_{in}+1)}, \pi_{N}^{(t_{in}+1)}.

t_{in} = t_{in} + 1

UNTIL(收敛)

\hat{\pi}_{p} = \pi_{p}^{(t_{in})}
```

GMM 拟合预测分数的效果如图 4 所示,其概 率密度函数(Probability Density Function, p.d.f.) 大致符合正类和负类预测分数各自所服从分布的叠 加.尽管 PU 学习模块和类别先验估计模块的协同 工作基本达成了本文目标,即在类别先验未知条件 下进行深度 PU 学习,但它无法确保其迭代过程的 稳定性,*π*_P可能随着迭代轮次增长而发生震荡,进而 影响模型最终的分类性能.为此,本文设计了迭代过 程稳定模块,以平稳所提框架的迭代过程.

3.5 迭代过程稳定模块

迭代过程稳定模块基于平滑假设和聚类假设, 通过平均教师^[23]的架构对模型历史参数进行集成, 以提升深度模型的抗扰动能力;另外结合使用温度 锐化技术,加强深度模型对可信样本的预测置信度, 图 4 中不同类别的钟形曲线因之更可区分,从而平 稳了类别先验估计的迭代过程.

3.5.1 平均教师

平均教师^[23]由深度神经网络结构相同的学生 模型和教师模型组成.在所提框架中,学生模型为 PU学习模块中的深度神经网络,其参数为 Ø,由 PU学习模块训练得到;教师模型的参数 Ø'不直接 在 PU 数据上学习,而是由学生模型的历史参数经 过时序上的 EMA 运算得到,即:

给定样本 x,教师模型的输出 $f(x; \Theta')$ 可以看 作对学生模型输出 $f(x; \Theta)$ 的轻微扰动. 根据平滑假 设,若模型的输出高度可信,那么 $f(x; \Theta')$ 和 $f(x; \Theta)$ 应当具有一致性. 因此,平均教师将教师模型的输出 作为监督信息,指导学生模型的训练. 首先将样本 x关于学生模型 Θ 的输出置信度 c 定义为由学生模 型给出的最大类条件概率,即:

$$c = \max(s, 1-s) \tag{29}$$

结合基于置信度的掩码技术^[24],并将置信度阈值设 为 r,那么学生模型和教师模型的输出一致性损失 可定义为

$$L_{\rm con} = \frac{1}{n_{\mathbf{x} \in \mathbf{X}_{\rm PU}}} \mathbb{1}(c \ge \tau) \| f(\mathbf{x}; \boldsymbol{\Theta}) - f(\mathbf{x}; \boldsymbol{\Theta}') \|_{2}^{2} (30)$$

其中,1(•)是指示函数,当满足条件(•)时函数值取 1,反之取 0.

3.5.2 温度锐化

温度锐化通过调整预测分数的温度(Temperature,T),降低预测类别条件分布的信息熵.其理论 依据是聚类假设,即分类决策面不应穿过数据密集 的区域.因此,样本 x 经过模型 $f(\cdot; \Theta)$ 的预测类别 条件分布具有较低的信息熵.给定预测分数 s,温度 锐化使用锐化函数来降低上述信息熵:

Sharpen(s,T) =
$$\frac{s^{1/T}}{s^{1/T} + (1-s)^{1/T}}$$
 (31)

结合平均教师中提及的基于置信度的掩码技术,只 对高置信度的预测分数进行温度锐化,那么温度锐 化的优化目标可写作:

$$L_{\text{sharpen}} = \frac{1}{n} \sum_{x \in \mathbf{X}_{\text{PU}}} \mathbb{1}(c \ge \tau) \| s - Sharpen(s, T) \|_{2}^{2} (32)$$

3.6 优化目标

综上所述,所提框架由 PU 学习模块、类别先验 估计模块和迭代过程稳定模块组成,其单次迭代步 骤中的整体优化目标为

$$L = \hat{R}_{nnPU}^{\ell_{sigmoid}} + \lambda_1 L_{con} + \lambda_2 L_{sharpen}$$
(33)

其中, λ_1 , λ_2 为超参数.所提框架反复迭代执行深度 PU学习和类别先验估计,直至 L 收敛或迭代次数 达到预设值,其算法伪代码如算法 2 所描述.

算法 2. 迭代式的深度 PU 学习与类别先验估计. 输入:

 X_{PU} :输入特征矩阵, $X_{PU} = [x_1; x_2; \cdots; x_n].$

 X_L :被标注正例的输入特征矩阵, $X_L = [x_1, x_2; \dots; x_{n_L}]$ α :式(28)中的平滑系数.

τ:式(30)和(32)中的置信度阈值.

λ1,λ2:优化目标(33)中一致性损失和温度锐化的超参数.

 $f(\bullet; \bullet)$:所选深度神经网络结构对应的函数映射.

O(0), **O**(0): 学生模型和教师模型的初始参数.

输出:

 $\hat{\pi}_{P}$:正类先验估计值.

∅,**∅**′:算法结束时学生模型和教师模型的参数.

过程:

t = 0

REPEAT

 $S, S' = sigmoid [f(X_{PU}; \boldsymbol{\theta}_{(t)})], sigmoid [f(X_{PU}; \boldsymbol{\theta}_{(t)})]$ 对 S'运用算法 1,得到 $\hat{\pi}_P$. 根据式(16)计算 $\hat{R}^{estimated}$.

依据式(10) / 异 Λ_{nnPU} .

根据式(30)计算 L_{con}.

根据式(32)计算
$$L_{\text{sharpen}}$$
.

 $L = \hat{R}_{nnPU}^{\ell_{sigmoid}} + \lambda_1 L_{con} + \lambda_2 L_{sharpen}$

对 L 运用梯度反向传播算法,得到 $\Theta_{(t+1)}$.

$$\boldsymbol{\Theta}_{(t+1)}^{\prime} = \alpha \boldsymbol{\Theta}_{(t)}^{\prime} + (1-\alpha) \boldsymbol{\Theta}_{(t+1)}$$

$$t = t + 1$$

UNTIL(收敛或迭代次数达到预设值)

 $\boldsymbol{\Theta}, \boldsymbol{\Theta}' = \boldsymbol{\Theta}_{(t)}, \boldsymbol{\Theta}'_{(t)}$

3.7 设计选项讨论

在对所提框架进行设计的过程中,存在多个设 计选项,分别是:(1)在算法流程方面,是迭代式地 进行深度 PU 学习和类别先验估计,抑或是基于单 次类别先验估计的深度 PU 学习;(2)类别先验的估 计方式,是使用 GMM 拟合预测分数的混合分布, 从而间接估计类别先验,抑或直接统计伪标签的正 类频率;(3)迭代过程稳定模块的必要性.接下来对 上述设计选项进行探讨.

3.7.1 迭代式 vs 单次估计

为解决类别先验未知问题,以往工作^[17-20]通常 首先进行类别先验估计,进而基于类别先验估计值 进行 PU 学习.然而,单次估计的随机性较强.况且 训练初期深度模型的学习不够充分,基于预测分数 进行单次估计的质量更加难以得到保障.此外,在深 度神经网络过拟合前,深度模型的分类性能的提升 一般意味着其提取的高层级特征更具代表性.上述 增益无法为单次估计所利用.因此,所提框架选用迭 代式地进行深度 PU 学习和类别先验估计的算法流 程,即随着深度模型的训练,预测分数愈发可靠,从 中估计的类别先验将更加准确;更准确的估计值进 一步促使模型分类性能的提升.深度 PU 学习和类 别先验估计二者相互促进.

3.7.2 GMM vs 伪标签频率

关于如何利用预测分数估计类别先验,其最为 简单直接的做法是统计伪标签中的正类频率:

$$\hat{\pi}_P = \frac{1}{n} \sum_{\mathbf{x} \in \mathbf{X}_{\text{PU}}} \mathbb{1}(s \ge 0.5) \tag{34}$$

尽管基于伪标签的估计方法符合常理,且易于实现, 但是它对预测分数的分布要求较为严苛.预测分数 不仅需要具有较好的类别可分性,且其分类面需位 于 *s*=0.5 附近.然而,上述要求往往难以满足.考虑 当预测分数的分值整体偏大或偏小时,伪标签的正类 频率将大幅高于或低于真实的类别先验.图5展示了



图 5 CIFAR-10 上不同方法(伪标签、GMM、本文方法)的 $\hat{\pi}_{p}$ 随训练 epoch 数的变化曲线(其中,彩色实线为重 复 5 次实验的均值; 阴影部分由 5 次实验的最小值和最大值决定. 黑色虚线表示 $\pi_{p} = 0.4$)

CIFAR-10上基于伪标签频率预测的 $\hat{\pi}_P$ 随着训练轮次的变化趋势.可以发现,训练初期偏低的 $\hat{\pi}_P$ 导致了模型对正类的拟合不足,于是其输出的预测分数整体偏低,进而造成了更低的 $\hat{\pi}_P$,由此形成了确认偏差(Confirmation Bias)^[51],即模型过拟合于训练初期不准确的预测,模型的预测分数和 $\hat{\pi}_P$ 逐渐向 0 靠近.此外,消融实验结果分析(见 4.5 节)进一步说明了基于伪标签的 $\hat{\pi}_P$ 不可行.相较于伪标签频率法,基于GMM 的类别先验估计方法通过将预测分数的分布 建模为双高斯成分的混合分布,于是无论预测分数 整体偏大或偏小,只要其分布大致符合两条钟形曲线加权叠加的特征,就能较为准确地估计类别先验, 从而应对预测分数分布整体偏移问题.图 5 和表 4 表明,迭代式的 GMM 基本达成本文目标,并取得 了良好的分类性能.

3.7.3 迭代过程稳定模块的必要性。

实验发现,尽管迭代式的 GMM 基本能够估计 类别先验,并且取得了良好的分类性能,但是其类别 先验估计过程不够稳定.如图 5 左二所示, 施 積 训 练 轮次的变化而震荡,5次实验的差异较大.为此, 本文特别设计了迭代过程稳定模块,通过引入平均 教师^[23]和温度锐化^[24]以平稳迭代过程.如图 5 右 一、右二所示,迭代过程稳定模块的效果明显,*π*_P的 震荡现象得到有效缓解,5次实验的标准差明显缩 小,这说明了迭代过程稳定模块使得所提框架更加 稳定可靠.另外,平均教师中学生模型与教师模型的 输出一致性为模型训练提供了训练标签以外的自监 督信息,模型最终的分类性能因之得到了进一步提 升(见表 4).

4 实 验

4.1 数据集与主干网络

本文的实验设置整体上与论文^[6]保持一致,在 三个基准数据集(手写数字识别数据集 MNIST^[52]、 服饰图片分类数据集 F-MNIST^[53]、普适物体识别 数据集 CIFAR-10^[54])和一个实际应用数据集(阿尔 茨海默病神经成像数据集 Alzheimer^①)上进行了实 验.各个数据集的基本信息和训练时使用的主干网 络如表1所示,接下来对各个数据集进行详细介绍.

表 1 数据集与主干网络概述

数据集	训练集大小	测试集大小	输入尺寸	π_P	标签个数	正类/负类	主干网络
MNIST	60000	10000	28×28	0.49	1000	奇数/偶数	多层感知机(6层)
$MNIST^2$	60000	10000	28×28	0.48	1000	$0 \sim 4/5 \sim 9$	多层感知机(6层)
F-MNIST	60000	10000	28×28	0.40	500	上装/非上装	多层感知机(6层)
F-MNIST ²	60000	10000	28×28	0.30	500	1,4,7/其它	多层感知机(6层)
CIFAR-10	50000	10000	$3 \times 32 \times 32$	0.40	1000	交通工具/动物	卷积神经网络(13 层)
CIFAR-10 ²	10000	2000	$3 \times 32 \times 32$	0.50	2000	猫/狗	卷积神经网络(13 层)
Alzheimer	5121	1279	$3\!\times\!224\!\times\!224$	0.50	769	痴呆/健康	Resnet-18

4.1.1 MNIST

MNIST^[52]共有 70000 个手写阿拉伯数字图像, 其中训练集和测试集分别包含 60000 和 10000 张图 片.每个图片的输入尺寸为 28×28.将奇数 1、3、5、7、 9 当作正类,偶数 0、2、4、6、8 作为负类.经统计,类别 先验 π_P 为 0.49.同时,为符合 PU 学习的实验设定, 本文随机选择了训练集中的 1000 个正样本,保留其 标签,并将剩余未被选择的样本作为无标签数据. MNIST 数据集使用的主干网络是一个 6 层的多层感 知机,其隐层由 4 个无偏置的线性层构成,每一个线 性层都紧接着批归一化^[55]操作和 ReLU^[56]激活函 数.MNIST²将小于 5 的数当作正类,反之则为负类. 经统计,类别先验 π_P 为 0.48.

4.1.2 F-MNIST

F-MNIST^[53]数据集是服饰图片分类任务下的

MNIST 数据集,输入图片尺寸,训练、测试样本数及 类别数的设定同经典的 MNIST 相同. F-MNIST 将 上装类:0(T恤)、2(套衫)、4(外套)、6(汗衫)视作正 类,将非上装类:1(裤子)、3(裙子)、5(凉鞋)、7(运动 鞋)、8(包)、9(踝靴)视作负类,则 π_P为 0.40.为符合 PU学习的实验设定,本文随机选择了训练集中的 500个正样本,保留其标签,并将剩余未被选择的样 本作为无标签数据. F-MNIST 使用的主干网络与 MNIST 相同.类似地,F-MNIST²将 1、4、7 视作正 类,π_P为 0.30.

4.1.3 CIFAR-10

CIFAR10^[54]是一个常用于目标识别的图像基

① Alzheimer's Dataset. https://www.kaggle.com/tourist55/ alzheimers-dataset-4-class-of-images

准数据集,由 60000 个 $3 \times 32 \times 32$ 彩色图像组成,包 含 10 类对象,其中,每个类有 6000 个图像. CIFAR-10 将交通工具类:0(飞机)、1(汽车)、8(船)和 9(卡 车)视作正类,将动物类:2(鸟类)、3(猫)、4(鹿)、5 (狗)、6(蛙)、7(马)视作负类.可知, π_P 为 0.40.为符 合 PU 实验设定,本文随机选择了训练集中的 1000 个正样本,保留其标签,并将剩余未被选择的样本作 为无标签数据. CIFAR10 采用的主干网络为全卷积 深度神经网络^[57]. CIFAR10²则将 3(猫)视作正类,5 (狗)视作负类, π_P 为 0.50,标签个数为 2000.

4.1.4 Alzheimer

Alzheimer 数据集是关于阿尔兹海默症的多分 类数据集,总共包含了4类共5760张脑部磁共振成 像(Magnetic Resonance Imaging, MRI),训练集和 测试集的大小分别为 5121 和 639. 如图 6 所示, AD 有4个病程发展阶段,包括非痴呆(Non-Demented)、 轻微痴呆(VeryMild Demented)、轻度痴呆(Mild Demented)和中度痴呆(Moderate Demented).其 中,非痴呆与轻微痴呆的 MRI 相近. 值得注意的是, 正常衰老与 AD 之间存在一种过度状态,称作轻度 认知功能障碍(Mild Cognitive Impairment, MCI). 该状态下的 MRI 与常人无异,然而部分 MCI 被临/ 床诊断为 AD^[6]. 因此,训练集中的中度痴呆和轻度 痴呆样本均可被临床诊断,视作正样本.而轻微痴呆 的 MRI 与健康人相差无几,共同组成了无标签数 据.本文采用 ImageNet 预训练后的残差网络 Resnet-18^[58]作为 Alzheimer 数据集的主干网络,其输 入图像尺寸被统一裁剪为 3×224×224.



图 6 AD 的 4 个病程发展阶段对应的 MRI(非痴呆的 MRI(左上)与轻微痴呆的 MRI(右上)十分相似)

4.2 评价指标

二分类问题可根据样本真实标签的正负以及模

型预测的正确与否,划分为4种基本情况,分别是样本被:(1)正确地预测成正类(TruePositive, TP); (2)正确地预测成负类(TrueNegative, TN);(3)错误地预测成正类(False Positive, FP)和(4)错误地预测成负类(False Negative, FN).基于上述4种情况定义以下4类评价指标,分别是准确率(Accuracy, Acc)、精确率(Precision)、召回率(Recall)和 F1-score(f_1).

(1)准确率,反映算法预测的总体正确程度,即 全部预测中,预测正确的数目占比.

$$Acc = \frac{TP + TN}{TP + TN + FP + FN} \tag{35}$$

(2)精确率,即预测结果为正类中,预测正确的 数目占比.

$$precision = \frac{TP}{TP + FP} \tag{36}$$

(3) 召回率,即真实标签为正类中,预测正确的 数目占比.

$$recall = \frac{TP}{TP + FN}$$
 (37)

(4) F1-score,精确率和召回率的调和均值,用 于综合评估精确率和召回率,是 PU 学习的重要评 价指标.

$$f_1 = \frac{2 \times precision \times recall}{precision + recall}$$
(38)

4.3 实现细节

所提方法基于深度学习框架 Pytorch 实现. 所 有实验均在GeForce RTX 3090上进行.为确保GMM 的初次拟合结果,首先对所提框架的学生模型进行 5个轮次的热身训练(warm up),以使得初始预测分 数的类别条件分布基本可分.式(33)中超参数 λ1 在 MNIST、F-MNIST、CIFAR-10 和 Alzheimer 的取 值分别为 0.05、0.05、0.3 和 0.3; 而超参数 λ₂ 采取 余弦退火的策略,即随着训练轮次的增加,其值逐渐从 0 增至 1. 式(28)中的平滑系数 α 设为 0. 999,式(32) 中的温度 T 均设为-10,置信度阈值 τ 均设为 0.8.基 准数据集的初始步长均设置为 5e-4, Alzheimer 的初 始步长则为 1e-4; 在训练过程中使用 Adam 优化 器^[59]和余弦退火策略器 Cosine-AnnealingLR^[60]对 步长进行调整;其中,Adam的权重衰减系数为 5e-3. 基准数据集的训练批大小(Batch Size)均设置为 256, 而 Alzheimer 的训练批大小为 128. 为了保证实验的 准确性和可复现性,程序的每次运行都会指定一个 整数作为 NumPy、Pytorch、random 标准库、CUDA 等模块的随机种子.所有实验在随机种子 0~4 下重 复 5 次进行,并记录教师模型或学生模型在验证集 上准确率最高时对应各项指标的均值和标准差.

4.4 对比实验

4.4.1 对比方法

(1) supervised,即有监督学习方法,将所有无标 签样本当作负例.由于 PU 学习中无标签样本数目往 往远大于被标注的正例数目,即 $n_U \gg n_L$,深度模型将 过拟合于负例.因此,对无标签数据进行下采样,即 从集合 $X - X_L$ 中无放回地随机抽取 n_L 个样本当作 负例,与已有的 n_L 个正例构成二分类监督学习训练 集.显然,上述操作可能会将无标签数据中的正例当 作负例,从而引入了标签噪声(Label Noise).

(2) uPE 是 PE^[17] 与 uPU^[15] 的结合. PE^[17] 通过 实现正类条件密度与混合密度之间的部分匹配估计 类别先验. 其中, Pearson 散度用于衡量上述密度间 的差异.

(3) nnPE 是 PE 与 nnPU^[16]的结合.

(4) uKM1 是 KM1^[18] 与 uPU 的结合. KM1 和 KM2^[18] 基于再生核希尔伯特核空间. 在该空间下, 正类条件分布是混合分布核的重新加权. 通过最小 化正类条件分布的核嵌入距离,找到最佳的权重和 类别先验估计值.

(5) nnKM1 是 KM1 与 nnPU 的结合.

(6) uKM2 是 KM2 与 uPU 的结合.

(7) nnKM2 是 KM2 与 nnPU 的结合.

(8) uTIcE 是 TIcE^[19] 与 uPU 的结合. TIcE^[19] 指出任一数据子域的正类标签频率是真实正类标签 频率的下界,并通过自顶向下的决策树,搜索使得正 类标签频率最大的数据子域,从而间接逼近真实的 正类标签频率,进而估计类别先验.

(9) nnTIcE 是 TIcE 与 nnPU 的结合.

(10) uCDMM是CDMM^[20]与uPU的结合.CD-MM^[20]结合循环坐标下降和 MM 优化框架,优化逻 辑回归的参数似然函数,估计正类标签频度,进而估 计类别先验.

(11) nnCDMM 是 TIcE 与 nnPU 的结合.

(12)Ours,即本文所提框架.

4.4.2 实验结果及分析

表 2 展示了不同算法在各个数据集上的对比实 验结果.总体来看,所提方法在多数评价指标上取得 了最优性能.

具体而言,Ours 在 MNIST 上取得了最优的准确率、召回率和 F1-score,相较于次优结果分别平均

提升了 8.33%、3.54% 和 7.84%; 在 MNIST²上取 得了最优的准确率和 F1-score,较次优结果分别平 均提升了 0.64%和 0.42%. 在 FMNIST 上取得了 最优的准确率、召回率和 F1-score,较次优结果分别 平均提升了 1.77%、2.79% 和 2.27%;在 FMNIST² 上取得了最优的准确率、召回率和 F1-score,较次 优结果平均分别提升 0.52%、2.21%和 1.16%. 在 CIFAR10 上取得了最好的准确率、召回率和 F1-score, 较次优结果分别平均提升了 3.40%、4.40%和 4.55%; 在CIFAR-10²上取得了最优的准确率和 F1-score, 较次优结果分别平均提升了 7.51%和 3.39%.由于 Alzheimer 的主干网络是 ImageNet 预训练后的模 型,这减轻了类别先验估计值不准确的负面影响. Ours的优势尽管有所削减,但仍在主要评价指标准 确率和 F1-score 上取得了最优性能,较次优结果分 别平均提升了 1.25% 和 0.86%. 此外, Alzheimer 上的实验结果表明所提框架有助于 AD 的早期诊 断,具有实用价值. Ours 与 supervised 在 Alzheimer 上针对病患各个病程的检出率如图7所示,鉴于 Alzheimer 训练集中所有的中度痴呆或轻度痴呆均 为有标签的正例,而轻微痴呆缺失标签,Ours 对轻 微痴呆的检出率相较 supervised 具有 40.40%的显 著提升,对中度痴呆也有近17.87%的提升.此外, 由于测试集中属于轻度痴呆的样本数仅有 12 个, Ours 和 supervised 都对轻度痴呆达到了 100%的检 出率.考虑到轻微痴呆和正常人的 MRI 相近,现 实条件下轻微痴呆难以被专家诊断, Ours 相较于 supervised 在轻微痴呆检出率上的显著提升,说明 了所提方法有助于 AD 病患尽早诊断和治疗,进而 利于保障患者的生命健康.

Ours 的先进性很大程度上可以归因于其类别 先验估计准确,且稳定可靠.以 MNIST、FMNIST、 CIFAR-10 和 Alzheimer 为例,不同类别先验估计算 法在各个数据集上的对比实验结果如图 8 所示,观 察可知,Ours 的类别先验估计均值在 CIFAR-10 和 Alzheimer 上最接近真实的类别先验,其在 MNIST、 F-MNIST 上的 $\hat{\pi}_p$ 与 π_p 也仅平均相差 0.01 和 0.03. 此外,Ours 随机 5 次实验的 $\hat{\pi}_p$ 标准差较小,最高不 超过 0.02. 由此可见,Ours 兼具较小的类别先验估 计误差和标准差,因而具有更好且更稳定的分类性 能.相反的是,其它类别先验估计算法无法兼顾类别 先验估计的准确性和稳定性.例如 KM1 在 MNIST 上的估计均值最接近真实值,但其 5 次实验的差异

表 2 对比实验结果(单元格记录了均值(标准差),最优性能以加粗黑体标出,次优性能以倾斜黑体标出.根据配对 t 检验的 结果,分别用√(X)表示 Ours 明显优于(更差)的相应方法,置信水平为 95%) (单位:%)

数据集	评价指标	supervised	uPE	nnPE	uKM1	nnKM1	uKM2	nnKM2	uTIcE	nnTIcE	uCDMM	nnCDM-M	Ours
	1.00	85.89√	86.07 🗸	86.33 \	81.88√	85.44 🗸	$64.27\checkmark$	64.33 🗸	77.74	79.92	66.89√	$57.74\checkmark$	94.66
	ACC	(0.43)	(0.75)	(0.35)	(8.48)	(1.59)	(18.54)	(18.61)	(14.37)	(11.69)	(12.19)	(15.62)	(0.37)
-		89.35 🗸	85.83√	83.48√	86.50 🗸	83.70 🗸	95.09	94.97	94.74	94.29	84.87	97.13	<i>95.41</i>
数据集 	precision	(1.69)	(1.15)	(1.22)	(2.05)	(3.02)	(6.78)	(7.09)	(1.68)	(1.28)	(10.78)	(6.42)	(1.05)
MNIST -		81.09	85.93√	<i>90.14</i>	74.73	87.77	31.97	32.28	58.53	63.24	43.39 🗸	17.06	93.68
	recall	(2.47)	(1.71)	(2.27)	(19.87)	(4.97)	(43.84)	(44.22)	(31.66)	(25.84)	(30.89)	(38.09)	(0.94)
-		84.98	85.87	86.66 √	78.93	85.56 🗸	33.43√	33.55 🗸	67.18	73.06	50.30 🗸	17.12	94.53
	<i>F</i> -score	(0.70)	(0.81)	(0.52)	(13.81)	(1.81)	(45.79)	(45.94)	(29.97)	(19.33)	(31.48)	(38.18)	(0.37)
	100	81.84 🗸	91.22 🗸	92.70V	74.11	75.23	87.80√	89.67√	62.83 🗸	66.09	68.55 🗸	71.23	93.34
	ACC	(0.40)	(0.33)	(0.32)	(23.29)	(24.30)	(2.11)	(1.88)	(15.50)	(17.91)	(11.87)	(12.51)	(0.75)
-		90.91	92.89√	91.92	96.30	95.51	92.35	93.31 🗸	96.19	97.04	94.69	94.38	95.25
10000	precision	(1.04)	(1.94)	(1.12)	(3.46)	(4.21)	(2.12)	(1.44)	(4.45)	(2.82)	(2.17)	(2.50)	(1.51)
MNIST ² -		71.86 🗸	89.87	94.09X	53.14	56.38	83.26	86.15 🗸	29.38	36.01 🗸	40.80 🗸	47.20√	91.64
	recall	(1.30)	(2.35)	(1.04)	(48.55)	(51.48)	(5.41)	(4.79)	(31.62)	(37.10)	(23.82)	(26.81)	(1.03)
-	n	80.26	91.31 🗸	92.98V	54.66	55.93	87.45√	89.50√	37.36 🗸	42.87 🗸	53.82 🗸	58.79√	93.40
	<i>F</i> -score	(0.57)	(0.41)	(0.28)	(49.90)	(51.06)	(2.58)	(2.24)	(37.32)	(41.12)	(24.32)	(25.89)	(0.72)
-	100	90.87√	93.39V	93.32√	86.87	87.39	86.43	86.48	62.53 🗸	66.42 🗸	88.33	86.58	95.16
	ACC	(5.09)	(0 . H)	(0.15)	(15.02)	(13.72)	(14.78)	(14.80)	(5.48)	(8.65)	(11.52)	(14.87)	(0.30)
-		91.06√	90.99√	90.29	93.72	93.18	92.46	92.80	98.98×	95.50	92.37	92.51	92.09
E MNICT	precision	(0.39)	(0.57)	(1.21)	(3.68)	(3.26)	(4.38)	(4.21)	(2.28)	(6.47)	(2.43)	(4.39)	(0.52)
F-MINIST -	11	85.58√	92.65 🗸	93.38	73.47	75.10	73.78	73.49	6.67	$17.59\checkmark$	77.91	74.13	96.17
	recall	(1.31)	(0.60)	(1.83)	(41.10)	(37.98)	(41.28)	(32.28)	(14.49)	(23.38)	(32.31)	(41.44)	(0.90)
_	E	88.23√	<i>91.81</i> √	91.79√	73.58	76.23	73.10	73.13	10.00 🗸	$24.64\checkmark$	80.19	73.30	94.08
	1-score	(7.30)	(0.13)	(0.30)	(41.13)	(35.07)	(40.87)	(33.55)	(21.53)	(29.12)	(26.17)	(40.98)	(0.38)
– F-MNIST ² –	ACC	86.55 🗸	89.59√	89.66 🗸	85.77 🗸	88.70	86.55 🗸	88.98√	$73.62\checkmark$	73.96 🗸	89.28√	89.95V	90.47
	ACC	(1.01)	(0.40)	(0.47)	(2.59)	(1.43)	(2.36)	(0.70)	(4.51)	(6.21)	(0.36)	(0.22)	(0.34)
	pracision	$79.64\checkmark$	84.28	83.11	$86.41 \times$	$86.66 \times$	83.72	86.34 imes	97.48 $ imes$	97.66X	85.80	84.07	83.87
	precision	(2.66)	(2.15)	(2.23)	(3.03)	(1.79)	(2.34)	(1.93)	(2.74)	(3.22)	(4.35)	(1.27)	(1.81)
	recall	$74.24\checkmark$	80.43 🗸	$82.36\checkmark$	$62.81\checkmark$	$73.67\checkmark$	68.49 \checkmark	75.31	$12.50\checkmark$	$14.03\checkmark$	$77.51\checkmark$	82.10	84.57
-		(2.23)	(3.20)	(1.72)	(12.35)	(4.75)	(9.18)	(4.58)	(15.48)	(21.96)	(4.73)	(2.24)	(1.28)
	<i>F</i> -score	76.81 🗸	$82.24\checkmark$	82.70√	$72.06\checkmark$	$79.58\checkmark$	75.08	80.34	$19.48\checkmark$	$19.79\checkmark$	$81.25\checkmark$	83.04 🗸	84.20
		(1.65)	(0.97)	(0.49)	(7.43)	(3.09)	(5.77)	(2.02)	(23.54)	(28.83)	(0.83)	(0.65)	(0.30)
	ACC	84.70 🗸	82.86√	83.78√	75.40	84.92 🗸	84.31 🗸	86.58	81.84	82.57	74.81	79.27	89.98
-	nee	(0.52)	(2.87)	(0.82)	(19.84)	(0.95)	(4.92)	(2.57)	(12.37)	(12.64)	(13.69)	(10.86)	(0.83)
	precision	83.64 🗸	$79.42\checkmark$	80.45 🗸	73.01	80.60 🗸	$82.26\checkmark$	84.04 🗸	88.82	90.54	90.33	87.29	87.36
CIFAR-10-	Freedom	(2.28)	(6.00)	(2.26)	(18.71)	(2.64)	(3.44)	(2.26)	(6.41)	(5.32)	(8.99)	(7.27)	(0.82)
011111110	recall	76.98 🗸	77.97 🗸	78.66 🗸	83.23	82.24 🗸	77.31	82.01	65.19	65.18	45.78	59.44	87.63
-		(4.21)	(5.07)	(1.33)	(9.51)	(2.46)	(12.96)	(6.09)	(36.79)	(36.50)	(42.03)	(33.44)	(1.79)
	F-score	80.06	78.47	79.51	75.49	81.35 🗸	79.31	<i>82.94</i> √	66.90	67.74	47.92	63.07	87.49
F-MNIST -		(1.42)	(3.15)	(0.69)	(10.35)	(0.98)	(8.32)	(3.71)	(37.52)	(37.88)	(43.86)	(35.33)	(1.13)
	ACC	62.47	63.87	55.71	65.74	66.46	53.28	50.00	65.22	53.86	64.10	66.30	73.81
-		(1.01)	(8.50)	(3.82)	(9.05)	(10.66)	(7.33)	(0.00)	(3.11)	(3.85)	(8.27)	(9.35)	(1.50)
	precision	61.46	60.41	53.29	78.19	76.91	95.11×	100.00×	60.94	52.21	76.47	76.55	71.96
CIFAR-10 ² -		(1.31)	(7.13)	(2.29)	(12.35)	(13.15)	(10.94)	(0.00)	(3.55)	(2.27)	(13.29)	(13.29)	(0.97)
	recall	67.36	89.06×	96.20×	50.34	54.72	9.70	0.00	86.90	97.16×	48.20	55.80	78.00
-		(7.35)	(7.36)	(3.23)	(28.50)	(33.66)	(21.69)	(0.00)	(5.26)	(3.28)	(27.48)	(31.58)	(3.04)
	F-score	64.07	71.43	68.52	53.93	55.30	11.81	0.00	$71.45 \checkmark$	67.85	51.94	56.10	74.84
		(3.04)	(2.96)	(1.25)	(30.29)	(32.21)	(26.42)	(0.00)	(0.89)	(1.12)	(29.29)	(31.47)	(1.79)
	ACC	$61.47 \checkmark$	$69.74\checkmark$	$70.74 \checkmark$	$70.45 \checkmark$	$71.95 \checkmark$	$69.85 \checkmark$	$71.62 \checkmark$	$70.27 \checkmark$	$71.62 \checkmark$	$63.77 \checkmark$	$65.86 \checkmark$	73.20
-		(4.44)	(1.00)	(1.37)	71.00	(0.39)	71.00	(0.74)	71.50	70.44	(3.49)	(1.00)	(0.33)
	precision	(8.04√ (14.0°)	(2.43	(1.77)	(1.02	(1.56)	(2.14)	$(0.30 \checkmark$	(1.58	(0.44√	15.66	(5.76)	(0.72)
Alzheimer -		20 00 /	(2.13)	79.00	(4.44)	(1.07)	(2.14)	74.04	(4.44)	74.40	(0.03)	(3.70)	75 20
	recall	38.09√ (17.97)	$03.98 \lor$	(9.07)	09.48 (6.17)	75.52 (4 47)	(6.56)	(2.21)	$07.01 \lor$	(2.80)	$43.00 \checkmark$	$\frac{41.10}{(8.82)}$	(2.45)
F-MNIST ² CIFAR-10 ²		47 57.7	67 79. /	71 19	70.04	72 96	68 21. /	79 / 9. /	60.26.7	79.90	52 60. /	57.01.7	72 73)
	F-score	(12.96)	(3.41)	(3.67)	(2.15)	(1.29)	(2.76)	(1.00)	(2.11)	(1.10)	(12.45)	(5.30)	(0.95)



图 7 Ours 与 supervised 在 Alzheimer 上患者各个病程(轻 微痴呆、轻度痴呆和中度痴呆)的检出率



图 8 不同类别先验估计方法在各个数据集的均值-标准差箱 型图(黑色虚线表示各个数据集的 π_P)

过大,因此 uKM1 和 nnKM1 的平均分类性能欠佳, 且具有较大的标准差;PE 在 MNIST 上的类别先验估 计最稳定,但其 $\hat{\pi}_p$ 过于偏离 π_p ,因此 uPE 和 nnPE 的 平均分类性能欠佳.

4.4.3 已知 π_P的对比实验结果与分析

本文致力于解决 PU 学习中深度场景下类别先 验未知问题. 实验发现,即使所提框架不利用 π_p ,也 能达到与已知 π_p 的 PU 算法相媲美的分类性能, 甚至在部分结果上略有超越. 表 3 展示了 MNIST、 FMNIST、CIFAR-10 和 Alzheimer 数据集上 Ours 与基于真实类别先验 PU 算法的对比实验结果. 其 中,uPU 是 PU 学习的无偏风险估计器; nnPU 是 uPU 针对深度模型的扩展; Self-PU 在 nnPU 的基 础上引入自步学习训练策略和深度模型的自校准与 自蒸馏,代表了 PU 学习的最先进水平. Ours 整体 性能最佳,在 MNIST 上的准确率、精确率和 F1score 上分别平均超出次优结果 0.80%、2.17%和

表 3 Ours 与已知 π_p 的 PU 算法对比实验结果(根据配对 t 检验的结果,分别用 \checkmark (\times)表示 Ours 明显优于(更

差)的相应方法,置信水平为95%) (单位:%										
数据集	评价指标	uPU	nnPU	Self-PU	Ours					
	ACC	91.57 🗸	93.26√	93.86V	94.66					
	ACC	(1.27)	(0.38)	(0.30)	(0.37)					
	pracision	92.50 🗸	93.24V	93.12√	95.41					
MNIST	precision	(1.25)	(0.84)	(0.89)	(1.05)					
MINIST	recall	90.24 🗸	93.08√	94. 54 $ imes$	93.6 8					
		(2.90)	(1.28)	(0.99)	(0.94)					
4	Fescoro	91.33	93.15	93.82	94. 53					
V	1 Score	(1.44)	(0.42)	(0.31)	(0.37)					
	ACC	93.66 🗸	94.39√	94.75 V	95.16					
	ACC	(0.31)	(0.42)	(0.25)	(0.30)					
F-MNIST	bracision	92. 78	92.88	91.73√	92.09					
	precision	(1.91)	(0.94)	(0.80)	(0.52)					
		91.33 🗸	93.13√	95.50V	96.17					
	recutt	(2.06)	(0.84)	(0.61)	(0.90)					
	Facero	92.02 🗸	93.00 🗸	93.57V	94.08					
	r-score	(0.40)	(0.52)	(0.28)	(0.38)					
	ACC	88.56 🗸	89.07√	89.28 V	89.98					
	Acc	(0.52)	(0.38)	(0.72)	(0.83)					
	precision	87.35	86.35	86.16√	87.36					
CIFAR-10	precision	(2.57)	(1.68)	(0.78)	(0.82)					
	recall	83.66√	86.39	87.21	87.63					
	100011	(2.98)	(2.06)	(2.35)	(1.79)					
	F-score	85.39√	86.34 🗸	86.67V	87.49					
	I score	(0.72)	(0.52)	(1.06)	(1.13)					
	ACC	71.09 🗸	72 . 12√	72. 59V	73.20					
	ACC	(0.70)	(0.46)	(0.64)	(0.55)					
	hussision	71.55	69.05 🗸	69.83 🗸	72.25					
Alaboimor	precision	(2.55)	(1.79)	(1.03)	(0.73)					
Aizheimer	racall	70.49	80.44	79.56	75.30					
	<i>recall</i>	(6.10)	(4.18)	(3.36)	(2.45)					
	F-score	70.81 🗸	74.22	74.34	73.72					
	1 30010	(1.93)	(0.79)	(1.13)	(0.95)					

0.71%;在 F-MNIST 上的准确率、召回率和 F1score 超出次优结果 0.41%、0.67%和 0.51%;在 CIFAR-10 上的所有指标超出次优结果 0.60%、 0.01%、0.42%和 0.82%;在 Alzheimer 上的准确 率和精确率超出次优结果 0.61%和 0.70%.接下来 的消融实验结果进一步揭示了所提框架基于类别先 验估计,但优于基于真实类别先验 PU 算法的原因.

4.5 消融实验分析

为评估所提框架中各个设计的有效性与合理 性,以及探究对比实验结果中 Ours 性能优异的 原因,分别对 GMM、平均教师(L_{con})和温度锐化 (L_{sharpen})在基准数据集上进行了消融实验,其结果 如表4所示.其中,无GMM 的类别先验估计为迭代 式地统计伪标签的类别频率.总体而言,GMM 用于 估计类别先验,为深度 PU 学习奠定了基础;平均教 师大幅提升了分类性能,温度锐化则进一步巩固并 提升了分类性能;Ours 有效结合了之者的优势,性 能最佳.下面将结合实验结果对各个设计的作用进 行分析.

(1) GMM. 如表 4 所示, 无 GMM 的方法在各

个基准数据集上的评价指标均显著落后于 GMM (序号 0 vs 1、2 vs 4、3 vs 5、6 vs 8). 以准确率为例, 1比0在MNIST、F-MNIST和CIFAR-10上分别 平均提升了 1.76%、3.85%和 22.72%;4 比 2 在各个 基准数据集上分别提升了 3.17%、0.08%和 16.82%:5 比3分别提升了1.58%、3.46%和20.81%;8比6 分别提升了 3.42%、1.18%和 17.15%. 值得一提的 是,1 仅仅基于 GMM, 便取得了接近基于真实类别 先验的 PU 算法性能. 另外,无 GMM 的标准差远高 于 GMM 的标准差. 以准确率为例, MNIST、F-MNIST 和 CIFAR-10 上无 GMM 的标准差最大可 分别达到 6.29%、7.15%和 12.06%, 而 GMM 的标 准差最高不超过 2.12%. 无 GMM 部分标准差较大 的原因在于其类比先验估计过程出现了如图 5 左二 所示的情况,即训练初期偏低的 $\hat{\pi}_{p}$ 导致了模型对正 类的拟合不足,于是其输出的预测分数整体偏低,进 而造成了更低的 $\hat{\pi}_{P}$,最终导致算法失效.由此可见, GMM 的表现更加稳固可靠. 综上所述, GMM 能够 更好地估计类别先验,为所提框架的分类性能带来 了大幅提升.

表 4 基准数据集上的消融实验结果(\ 表示包含对应设计)

(单位:%)

数据集	序号	GMM	平均教师	伪标签	温度锐化	acc	precision	recall	f-score
	0					91.84(1.22)	94.02(1.53)	89.17(3.38)	91.48(1.42)
-	1	\checkmark				93.60(0.53)	92.98(0.73)	94. 13(1. 35)	93.54(0.58)
	2		\checkmark			90.83(6.05)	94.51(0.46)	86.37(12.78)	89.84(7.85)
	3				\checkmark	92.05(1.16)	94.18(1.53)	89.44(3.26)	91.71(1.38)
数据集 MNIST F-MNIST CIFAR-10	4	\checkmark	\checkmark			94.00(0.46)	94.13(0.42)	93.67(1.26)	93.89(0.51)
	5	\checkmark			\checkmark	93.63(0.48)	94.29(0.71)	92.68(0.72)	93.48(0.49)
	6		\checkmark		\checkmark	91.24(6.29)	94.76(0.75)	86.99(13.18)	90.27(8.11)
-	7	\checkmark	\checkmark	\checkmark		94.72(0.15)	95. 31 (1.05)	93.92 (1.12)	94.60(0.16)
	8(Ours)	\checkmark	\checkmark		\checkmark	94.66(0.37)	95.41 (1.05)	93.68(0.94)	94. 53 (0. 37)
F-MNIST	0					90.86(6.90)	94. 36 (1. 49)	82.30(19.20)	86.68(12.52)
	1	\checkmark				94.71(0.35)	91.61(1.43)	95.57(1.71)	93.53(0.43)
	2		\checkmark			94.89(0.24)	92.22(0.79)	95.26(1.07)	93.71(0.31)
	3				\checkmark	91.42(7.15)	93.13(2.22)	85.22(20.65)	87.56(12.95)
	4	\checkmark	\checkmark			94.97(0.33)	92.17(0.90)	95.58(1.21)	93.83(0.42)
	5	\checkmark			\checkmark	94.88(0.36)	92.08(0.44)	95.41(0.82)	93.71(0.46)
	6		\checkmark		\checkmark	93.98(0.83)	93. 22 (1. 38)	91.64(2.80)	92.39(1.16)
	7	\checkmark	\checkmark	\checkmark		95.05(0.36)	92.22(0.96)	95. 73 (1. 01)	93.93 (0.43)
	8(Ours)	\checkmark	\checkmark		\checkmark	95. 16(0. 30)	92.09(0.52)	96. 17(0. 90)	94. 08(0. 38)
	0					66.64(10.34)	32.54(44.64)	21.05(31.74)	24.95(35.84)
	1	\checkmark				89.36(0.55)	87.19(1.59)	86.12(2.33)	86.62(0.80)
	2		\checkmark			72.79(12.06)	79.64(8.82)	39.07(33.48)	47.24(31.22)
	3				\checkmark	67.15(11.40)	33.16(45.56)	21.86(33.36)	25.68(37.20)
CIFAR-10	4	\checkmark	\checkmark			89.61(0.65)	86.89(1.01)	87.18(0.89)	87.03(0.80)
	5	\checkmark			\checkmark	87.96(2.12)	82.89(2.41)	88. 11 (3. 82)	85.4(2.70)
_	6		\checkmark		\checkmark	72.73(11.95)	79.73(8.91)	38.78(32.96)	47.13(31.04)
MNIST F-MNIST CIFAR-10	7	\checkmark	\checkmark	\checkmark		89.95(0.54)	86.86(0.64)	88. 23 (1. 72)	87. 53(0. 79)
	8(Ours)	\checkmark	\checkmark		\checkmark	89.98(0.83)	87.36(0.82)	87.63(1.79)	87.49 (1.13)

(2) 平均教师. 仍以准确率为例,仅仅对比是否 使用平均教师,2 比 0 在 F-MNIST 和 CIFAR-10 上 分别提升了 4.03%和 6.15%;4 比 1 在各个基准数 据集上分别提升了 0.40%、0.26%和 0.25%;6 比 3 在 F-MNIST 和 CIFAR-10 上分别提升了 2.56%和 5.58%;7 比 5 在各个数据集上分别提升了 1.03%、 0.28%和 2.02%. 由此可见,平均教师基于平滑假 设为模型训练引入了额外的监督信息,进一步提升 了分类性能. 此外,8 比 5 的准确率标准差分别降低 了 0.11%、0.06%和 1.29%,这说明了平均教师有 助于所提框架的稳定性.

(3) 温度锐化. 以准确率为例, 仅对比是否使用 温度锐化, 3 比 0 在所有基准数据集分别提升了 0.21%、0.56%和 0.51%; 5 比 1 在 MNIST 和 F-MNIST 上分别提升了 0.03%和 0.17%; 6 比 2 在 MNIST 上提升了 0.41%; 8 比 4 在各个数据集上分 别提升了 0.66%、0.19%和 0.37. 由此可见, 温度锐 化进一步提升了分类性能.

目前,半监督学习的流行做法是对无标签数据 打上伪标签.根据聚类假设,伪标签和温度锐化的作 用都是使得模型在无标签数据上的输出信息熵应该 尽可能小,那么它们的实际效果应当相似.消融实验 通过7 vs 8 对二者进行了比较.实验结果显示,温度 锐化与伪标签的各项指标十分接近.

4.6 超参数敏感性分析

所提框架的超参数包括平均教师中式(28)里的

平滑系数 α 、温度锐化中式(32)里的温度 T、置信度 阈值 τ 、优化目标中式(33)里的 λ_1 和 λ_2 .鉴于 α 、T和 τ 的设置策略已被广泛研究^[24,42-43],接下来重点对 所提框架的主要参数 λ_1 和 λ_2 展开分析.

 λ_1 和 λ_2 分 别 控 制 平 均 教 师 (L_{con}) 和 温 度 锐 化 (L_{sharpen}) 的重要程度.图 9 展示了基准数据集上的超 参数敏感度分析结果,可以发现所提方法对超参数 的微小扰动不敏感,具有良好的稳定性.具体而言, 在 MNIST 上, λ_1 和 λ_2 的取值区间分别为[0.01, 0.2]和[0.2,5],准确率均值为94.21%,标准差仅 为 0.85%;在 F-MNIST 上 λ_1 和 λ_2 的取值区间分别 为[0.01,0.2]和[0.2,5],准确率均值为 94.89%, 标准差仅为 0.37%;在 CIFAR-10 上 λ_1 和 λ_2 的取值 区间分别为「0.1,0.7〕和「0.2,5〕,准确率均值为 89.95%,标准差仅为 0.26%.此外,合理地设置 λ₁ 和λ2有利于所提方法分类性能的进一步提升.例 如, MNIST 在 $\lambda_1 = 0.01$ 和 $\lambda_2 = 5$ 处取得了最佳性 能,对应准确率为 95.91%; F-MNIST 在 $\lambda_1 = 0.2$ 和 $\lambda_2 = 5$ 处取得了最佳性能,准确率为 95.40%. CIFAR- $10 在 \lambda_1 = 0.5 和 \lambda_2 = 2 处取得了最佳性能,准确率$ 为 90.45%.

所提方法关于 λ₁ 和 λ₂ 的稳定性一定程度上可 归因于基于置信度的掩码技术.在置信度掩码的作 用下,只有高度可信的样本参与了 L_{con}和 L_{sharpen}的计 算,大多数有害的监督信息被滤除.因此,λ₁和 λ₂的 微小扰动不会对算法性能造成严重的负面影响.



图 9 基准数据集上的超参数敏感度分析

4.7 主干网络的影响

在此前的对比实验中,不同的数据集采用了不同的神经网络结构.本节探究了不同主干网络对实验结论的影响,结果如表 5 所示.可以发现,所提方法在相同的主干网络设置下均取得了最优的准确率和 F1-score,进一步验证了所提方法的有效性.总体

来看,在相同的 PU 算法设置下,模型复杂度越高, 对应的分类性能通常更好.然而,对于所提框架而 言,Resnet-18 的分类性能与 CNN-13 接近.以准确 率和 F1-score 为例,Resnet-18 较 CNN-13 分别平 均仅提升 0.09%和 0.04%.这说明 CNN-13 足以胜 任 CIFAR-10 的普适物体分类任务. 表 5 CIFAR-10 上不同的主干网络对算法性能的影响

2683

	水。 CHAR 10 工作时的工作网络对身体任能的影响 (平位:											<u>.</u> :/0/	
主干网络	评价指标	supervised	uPE	nnPE	uKM1	nnKM1	uKM2	nnKM2	uTIcE	nnTIcE	uCDMM	nnCDMM	Ours
MLP-6	ACC	78.69	75.90	76.35	76.31	79.72	77.14	80.11	74.47	75.63	70.07	74.22	81.81
	АСС	(1.08)	(1.14)	(0.45)	(1.46)	(0.75)	(1.05)	(0.74)	(8.12)	(8.76)	(6.66)	(7.97)	(0.12)
		71.69	67.33	64.84	68.90	71.55	71.28	72.39	81.15	82.16	86.48	82.61	78.44
	precision	(2.58)	(1.95)	(0.50)	(3.57)	(1.71)	(3.30)	(1.58)	(10.74)	(10.17)	(8.44)	(9.94)	(0.95)
		77.63	77.53	89.31	75.11	82.02	72.41	81.38	52.46	55.05	32.44	49.38	75.22
	recall	(3.68)	(4.57)	(0.64)	(4.24)	(2.04)	(4.43)	(1.76)	(29.45)	(31.12)	(22.67)	(27.87)	(1.39)
		74.45	71.98	75.13	71.72	76.40	71.68	76.60	56.41	58.26	42.17	55.12	76.78
	F1	(1.04)	(1.63)	(0.42)	(1.08)	(0.55)	(0.97)	(0.62)	(31.56)	(32.62)	(26.62)	(30.86)	(0.30)
-	ACC	84.70	82.86	83.78	75.40	84.92	84.31	86.58	81.84	82.57	74.81	79.27	89.98
		(0.52)	(2.87)	(0.82)	(19.84)	(0.95)	(4.92)	(2.57)	(12.37)	(12.64)	(13.69)	(10.86)	(0.83)
	precision	83.64	79.42	80.45	73.01	80.60	82.26	84.04	88.82	90.54	90.33	87.29	87.36
		(2.28)	(6.00)	(2.26)	(18.71)	(2.64)	(3.44)	(2.26)	(6.41)	(5.32)	(8.99)	(7.27)	(0.82)
CNN-13	magall	76.98	77.97	78.66	83.23	82.24	77.31	82.01	65.19	65.18	45.78	59.44	87.63
	recatt	(4.21)	(5.07)	(1.33)	(9.51)	(2.46)	(12.96)	(6.09)	(36.79)	(36.50)	(42.03)	(33.44)	(1.79)
	F-score	80.06	78.47	79.51	75.49	81.35	79.31	82.94	66.90	67.74	47.92	63.07	87.49
	<i>P</i> -score	(1.42)	(3.15)	(0.69)	(10.35)	(0.98)	(8.32)	(3.71)	(37.52)	(37.88)	(43.86)	(35.33)	(1.13)
	4.00	85.35	89.36	86.05	88.89	88.10	89.07	88.25	81.04	82.43	76.02	79.11	90.07
	ACC	(0.47)	(0.52)	(0.73)	(0.27)	(0.54)	(0.33)	(0.36)	(12.29)	(12.68)	(10.03)	(11.56)	(0.44)
		84.68	86.98	80.05	85.70	83.52	86.42	83.36	90.90	90.12	90.81	90.35	87.93
D (10	precision	(0.11)	(1.13)	(3.17)	(2.10)	(3.36)	(1.74)	(1.41)	(5.14)	(5.53)	(5.76)	(5.93)	(0.96)
Resilet 10		77.39	86.34	87.24	86.84	87.88	86.32	88.34	60.38	65.31	46.54	55.86	87.17
	recall	(1.59)	(1.07)	(6.73)	(2.78)	(4.08)	(2.17)	(2.66)	(35.36)	(36.96)	(30.04)	(34.54)	(1.69)
	F 1	80.86	86.65	83.27	86.21	85.51	86. 33	85.74	64.86	67.53	55.00	61.31	87. 53
	F1	(0.81)	(0.62)	(1.64)	(0.46)	(0.55)	(0.45)	(0.65)	(36.84)	(37.87)	(32.62)	(35.59)	(0.65)

5 总 结

本文提出了一种迭代式的深度 PU 学习与类别 先验估计框架.该框架能准确且稳定地估计类别先 验,并利用所得先验估计值训练深度模型,因此无需 已知数据的真实先验分布,从而更利于解决实际应 用场景下的 PU 学习问题.本文所提出的迭代框架 包括以下核心步骤:(1)将网络的预测分数建模为 GMM,从而估计类别先验;(2)基于类别先验的估 计值,计算非负 PU 风险;(3)结合平均教师和温度 锐化技术,提高模型性能及稳定性.本文通过消融实 验验证了上述核心设计的合理性,并在多个数据集 上进行了大量对比实验,验证了该框架的先进性.其 中,在 Alzheimer 上的实验结果显示该框架能应用 于基于 MRI 的 AD 自动识别任务,且效果优于现有 方法,进而表明所提框架兼具科学价值和实用价值.

尽管本文为解决大规模 PU 数据上的类别先验 估计问题做出了早期探索与尝试,但所提框架仍然 存在一些不足之处.所提框架的主要局限性在于其 学习过程较长,超参数个数较多.为得到更好的分类 性能,选取最优的超参数较为繁琐.此外,所提框架 中的深度模型输出预测分数是一个黑盒过程,其可 解释性在理论方面相对欠缺.后续将考虑在算法的 简洁性、理论性与实用性方面做出改进.

参考文献

- [1] Hestness J, Narang S, Ardalani N, et al. Deep learning scaling is predictable, empirically. arXiv preprint arXiv: 1712.00409, 2017
- [2] Raffel C. Shazeer N, Roberts A, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. arXiv preprint arXiv:1910.10683, 2019
- [3] Mahajan D, Girshiek R, Ramanathan V, et al. Exploring the limits of weakly supervised pretraining//Proceedings of the European Conference on Computer Vision. Munich, Germany, 2018: 181-196
- [4] Zhou Z H. A brief introduction to weakly supervised learning. National Science Review, 2018, 5(1): 44-53
- [5] Merz C J, Clair D C S, Bond W E. Semi-supervised adaptive resonance theory (smart2)//Proceedings of the International Joint Conference on Neural Networks. Baltimore, USA, 1992, 3: 851-856
- [6] Chen X, Chen W, Chen T, et al. Self-PU: Self boosted and calibrated positive-unlabeled training//Proceedings of the International Conference on Machine Learning. Virtual Event, 2020: 1510-1519
- [7] Zhang Y L, Li L, Zhou J, et al. POSTER: A PU learning based system for potential malicious URL detection//Proceedings of the ACM SIGSAC Conference on Computer and Communications Security. Dallas, USA, 2017: 2599-2601
- [8] Li H, Chen Z, Liu B, et al. Spotting fake reviews via collective positive-unlabeled learning//Proceedings of the International

Conference on Data Mining. Shenzhen, China, 2014: 899-904

- [9] Bepler T, Morin A, Rapp M, et al. Positive-unlabeled convolutional neural networks for particle picking in cryo-electron micrographs. Nature Methods, 2019, 16(11): 1153-1160
- [10] Bekker J, Davis J. Learning from positive and unlabeled data: A survey. Machine Learning, 2020, 109(4): 719-760
- [11] He F, Liu T, Webb G I, et al. Instance-dependent pu learning by Bayesian optimal relabeling. arXiv preprint arXiv: 1808. 02180, 2018
- [12] Ienco D, Pensa R G. Positive and unlabeled learning in categorical data. Neurocomputing, 2016, 196: 113-124
- Liu L, Peng T. Clustering-based method for positive and unlabeled text categorization enhanced by improved TFIDF.
 Journal of Information Science and Engineering, 2014, 30(5): 1463-1481
- [14] Du Plessis M C, Niu G, Sugiyama M. Analysis of learning from positive and unlabeled data//Advances in Neural Information Processing Systems, Montreal Canada, 2014, 27: 703-711
- [15] Du Plessis M, Niu G, Sugiyama M. Convex formulation for learning from positive and unlabeled data//Proceedings of the International Conference on Machine Learning. Lille, France, 2015: 1386-1394
- [16] Kiryo R, Niu G, Du Plessis M C, et al. Positive-unlabeled learning with non-negative risk estimator//Advances in Neural Information Processing Systems. Long Beach, USA, 2017, 1674-1684
- [17] Du Plessis M C, Sugiyama M. Class prior estimation from positive and unlabeled data. IEICE Transactions on Information and Systems, 2014, 97(5): 1358-1362
- [18] Ramaswamy H, Scott C, Tewari A. Mixture proportion estimation via kernel embeddings of distributions//Proceedings of the International Conference on Machine Learning. New York, USA, 2016: 2052-2060
- [19] Bekker J, Davis J. Estimating the class prior in positive and unlabeled data through decision tree induction//Proceedings of the AAAI Conference on Artificial Intelligence. New Orleans, USA, 2018: 2712-2719
- [20] Łazęcka M, Mielniczuk J, Teisseyre P. Estimating the class prior for positive and unlabelled data via logistic regression// Advances in Data Analysis and Classification, 2021, 15(4): 1039-1068
- [21] Blanchard G, Lee G, Scott C. Semi-supervised novelty detection. Journal of Machine Learning Research, 2010, 11: 2973-3009
- [22] Arazo E, Ortego D, Albert P, et al. Unsupervised label noise modeling and loss correction//Proceedings of the International Conference on Machine Learning. Long Beach, USA, 2019: 312-321
- [23] Tarvainen A, Valpola H. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. arXiv preprint arXiv: 1703.01780, 2017

- [24] Xie Q, Dai Z, Hovy E, et al. Unsupervised data augmentation for consistency training//Advances in Neural Information Processing Systems. Virtual Event, 2020, 33: 6256-6268
- [25] Jiang L, Zhou Z, Leung T, et al. MentorNet: Learning datadriven curriculum for very deep neural networks on corrupted labels//Proceedings of the International Conference on Machine Learning. Stockholm, Sweden, 2018: 2304-2313
- [26] Kato M, Teshima T, Honda J. Learning from positive and unlabeled data with a selection bias//Proceedings of the International Conference on Learning Representations. New Orleans, USA, 2019
- [27] Hammoudeh Z, Lowd D. Learning from positive and unlabeled data with arbitrary positive shift//Advances in Neural Information Processing Systems. Virtual Event, 2020, 33: 13088-13099
- [28] Su G, Chen W, Xu M. Positive-unlabeled learning from imbalanced data//Proceedings of the International Joint Conference on Artificial Intelligence. Virtual Event, 2021: 2995-3001
- [29] Liu B, Lee W S, Yu P S, et al. Partially supervised classification of text documents//Proceedings of the International Conference on Machine Learning. Sydney, Australia, 2002, 2(485): 387-394
- [30] Zhou D, Bousquet O, Lal T N, et al. Learning with local and global consistency//Advances in Neural Information Processing Systems. Montréal Canada, 2004: 321-328
- Li X L, Liu B, Ng S K. Learning to identify unexpected instances in the test set//Proceedings of the International Ioint Conference on Artificial Intelligence. Hyderabad, India, 2007; 2802-2807
- [32] Du Plessis M C, Sugiyama M. Semi-supervised learning of class balance under class-prior change by distribution matching. Neural Networks, 2014, 50: 110-119
- [33] Northcutt C G, Wu T, Chuang I L. Learning with confident examples: Rank pruning for robust classification with noisy labels. arXiv preprint arXiv:1705.01936, 2017
- [34] Hsieh Y G, Niu G, Sugiyama M. Classification from positive, unlabeled and biased negative data//Proceedings of the International Conference on Machine Learning. Long Beach, USA, 2019: 2820-2829
- [35] Gong C, Shi H, Liu T, et al. Loss decomposition and centroid estimation for positive and unlabeled learning. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2019, 43(3): 918-932
- [36] Gong C, Wang Q, Liu T, et al. Instance-dependent positive and unlabeled learning with labeling bias estimation. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2021, 44(8): 4163-4177
- [37] Zhu X, Goldberg A B. Introduction to semi-supervised learning. Synthesis Lectures on Artificial Intelligence and Machine Learning, 2009, 3(1): 1-130
- [38] Zhu X J. Semi-supervised learning literature survey. University of Wisconsin-Madison, Madison, USA: Technical Reports: TR1530, 2005

- [39] Rasmus A, Berglund M, Honkala M, et al. Semi-supervised learning with ladder networks//Advances in Neural Information Processing Systems. Montréal, Canada, 2015, 28: 3546-3554
- [40] Bachman P, Alsharif O, Precup D. Learning with pseudo-ensembles//Advances in Neural Information Processing Systems. Montréal, Canada, 2014, 27: 3365-3373
- [41] Laine S, Aila T. Temporal ensembling for semi-supervised learning. arXiv preprint arXiv:1610.02242, 2016
- [42] Berthelot D, Carlini N, Goodfellow I, et al. MixMatch: A holistic approach to semi-supervised learning//Proceedings of the International Conference on Neural Information Processing Systems. Long Beach, USA, 2019: 5049-5059
- [43] Berthelot D, Carlini N, Cubuk E D, et al. ReMixMatch: Semi-supervised learning with distribution matching and augmentation anchoring//Proceedings of the International Conference on Learning Representations. New Orleans, USA, 2019
- [44] Niyogi P. Manifold regularization and semi-supervised learning: Some theoretical analyses. Journal of Machine Learning Research, 2013, 14: 1229-1250
- [45] Sohn K, Berthelot D, Carlini N, et al. FixMatch: Simplifying semi-supervised learning with consistency and confidence// Advances in Neural Information Processing Systems. Virtual Event, 2020, 33: 596-608
- [46] Stauffer C, Grimson W E L. Adaptive background mixture models for real-time tracking//Proceedings of the IEEE/CVF Computer Vision and Pattern Recognition Conference. Fort Collins, USA, 1999: 246-252
- [47] Ma Z, Leijon A. Bayesian estimation of beta mixture models with variational inference. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2011, 33(11): 2160-2173
- [48] Permuter H, Francos J, Jermyn I. A study of Gaussian mixture models of color and texture features for image classification and segmentation. Pattern Recognition, 2006, 39(4): 695-706
- [49] McLachlan G J, Krishnan T. The EM Algorithm and Exten-

ZHAO Yun-Rui, M. S. candidate. His research interest is weakly supervised learning. sions. Hoboken, USA: John Wiley & Sons, 2007

- [50] Shental N, Bar-Hillel A, Hertz T, et al. Computing Gaussian mixture models with EM using equivalence constraints// Advances in Neural Information Processing Systems. Vancouver, Canada, 2004, 16(8): 465-472
- [51] Arazo E, Ortego D, Albert P, et al. Pseudo-labeling and confirmation bias in deep semi-supervised learning//Proceedings of the International Joint Conference on Neural Networks. Virtual Event, 2020; 1-8
- [52] LeCun Y, Bottou L, Bengio Y, et al. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 1998, 86(11): 2278-2324
- [53] Xiao H, Rasul K, Vollgraf R. Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms. arXiv preprint arXiv:1708.07747, 2017
- [54] Krizhevsky A, Hinton G. Learning Multiple Layers of Features from Tiny Images [M. S. dissertation]. University of Toronto, Toronto, Canada, 2009
- [55] Ioffe S, Szegedy C. Batch normalization: Accelerating deep network training by reducing internal covariate shift//Proceedings of the International Conference on Machine Learning. Lille, France, 2015: 448-456
- [56] Nair V, Hinton G E. Rectified linear units improve restricted Boltzmann machines//Proceedings of the International Conference on International Conference on Machine Learning. Haifa, Israel, 2010: 807-814
- [57] Springenberg J T, Dosovitskiy A, Brox T, et al. Striving for simplicity: The all convolutional net. arXiv preprint arXiv: 1412.6806, 2014
- [58] He K. Zhang X, Ren S, et al. Deep residual learning for image recognition//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. Las Vegas, USA, 2016: 770-778
- [59] Kingma D P, Ba J. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014
- [60] Loshchilov I, Hutter F. Fixing weight decay regularization in Adam. arXiv preprint arXiv:1711.05101, 2017

XU Qian-Qian, Ph. D., professor. Her research interests include statistical machine learning and its applications in multimedia.

JIANG Yang-Bang-Yan, Ph. D. candidate. Her research interests include machine learning and computer vision.

HUANG Qing-Ming, Ph. D., chair professor. His research areas include multimedia computing, image processing, computer vision and pattern recognition.

Background

With the advent of big data, deep neural networks have attracted extensive attention and their performance has reached or even surpassed the human level in various tasks. Specifically, such great success usually relies on full supervision by a large amount of labelled data. However, it is hard to obtain intact label information in many real-world applications, even those with only binary options. For example, observed interactions between users and items in recommendation systems are labelled positives. Since many unconcerned factors like lack of exposure or other coincidences could account for missing interactions, we cannot view all the unobserved interactions as negatives. Similar scenarios include Alzheimer's disease recognition, malicious URL detection, and particle picking in cryo-electron micrographs, where we only have access to a few labelled positives with plenty of unlabeled data. Such a great demand motivates us to learn from positive and unlabeled data, also known as PU learning.

While PU learning based on unbiased risk estimators has achieved the state-of-the-art performance on several benchmarks, it relies on the knowledge of class prior, which might be unknown in reality. To this end, there emerge some class prior estimation algorithms designed for PU data. However, these algorithms are found difficult in handling large-scale data.

In this paper, we propose an iterative framework for deep PU learning and class prior estimation utilizing an unsupervised mixture model. Specifically, positive and negative classes have distinct predicted score distributions intuitively. We investigate and demonstrate our intuition and approximate the score distributions by a Gaussian Mixture Model with two components each representing the relevant class-conditional distribution. Accordingly, class prior is estimated through the weights of the Gaussian Mixture Model. With mean teacher and temperature sharpening incorporated, our proposed framework could estimate the class prior and learn from PU data simultaneously, achieving competent performance with other PU competitors based on the ground-truth class prior.

This work was supported in part by the National Key R&D Program of China under Grant (No. 2018AAA0102000), in part by National Natural Science Foundation of China (Nos. U21B2038, 61931008, 6212200758, and 61976202), in part by the Fundamental Research Funds for the Central Universities, in part by the Strategic Priority Research Program of Chinese Academy of Sciences (No. XDB28000000), in part by the Youth Innovation Promotion Association CAS.