



南京大學

研究生畢業論文 (申請博士學位)

論文題目 開放環境下的半監督學習研究

作者姓名 郭蘭哲

學科、專業名稱 計算機科學與技術

研究方向 機器學習與數據挖掘

指導教師 李宇峰副教授、黎銘教授

2022年5月20日

学 号： DG1933007

论文答辩日期： 2022 年 5 月 06 日

指 导 教 师：

(签字)

Semi-Supervised Learning for Open Environment

by

Lan-Zhe Guo

Supervised by

Associate Professor Yu-Feng Li and Professor Ming Li

A dissertation submitted to
the graduate school of Nanjing University
in partial fulfilment of the requirements for the degree of
DOCTOR OF PHILOSOPHY
in
Computer Science and Technology



Department of Computer Science and Technology
Nanjing University

April 20, 2022

南京大学研究生毕业论文中文摘要首页用纸

毕业论文题目： 开放环境下的半监督学习研究
计算机科学与技术 专业 2019 级博士生姓名： 郭兰哲
指导教师(姓名、职称)： 李宇峰副教授、黎铭教授

摘 要

半监督学习旨在利用无标注数据提升泛化性能，可以显著降低机器学习对数据标注的需求，是近年来机器学习领域的热门研究方向。既有半监督学习研究大多针对封闭环境，依赖一致的数据分布、静态的训练数据、充分的先验知识、平衡的类别比例，然而真实世界环境往往是开放的，存在“数据分布失配”、“数据动态流式”、“先验知识不足”、“类别比例失衡”等挑战，既有方法泛化性能时好时坏，难以适应开放环境。本文从数据输入-模型构建-标注输出三个层面，针对开放环境下半监督学习面临的四种典型挑战展开研究，提出了一套更稳健适应开放环境的半监督学习解决方案，具体包括如下创新成果：

- 1. 适于数据分布失配的稳健半监督学习。**针对数据分布失配的挑战，本文提出了一种稳健使用无标注样本的方法，通过加权经验风险最小化进行模型参数寻优并利用双层优化进行无标注样本赋权，避免分布外样本导致泛化性能下降。理论上，证明了该方法的收敛性并从经验风险和泛化风险角度分析了其稳健性，实验上，提升既有半监督学习方法对分布失配的稳健性超过 20%。
- 2. 适于数据动态流式的稳健半监督学习。**针对数据动态流式的挑战，本文提出了一种稳健适应流式数据的方法，通过影响力机制进行子集选择以适应受限的资源和变化的分布，避免数据动态变化导致泛化性能下降。本文方法可与任意半监督学习方法结合，一致有效提升既有方法的稳健性，大量实验结果表明，相比于既有方法，本文方法性能提升 20% 以上。
- 3. 适于先验知识不足的稳健半监督学习。**针对先验知识不足的挑战，本文提出了一种稳健选择半监督模型的方法，通过优化潜在最坏情况下的性能提升，避免模型选择错误导致泛化性能下降。理论上，分析了半监督学习实现稳健

性的条件, 相比以往理论结果更容易满足, 实验上, 既有方法均出现相比监督学习泛化性能下降的问题, 而本文方法在所有场景中都优于监督学习性能。

4. **适于类别比例失衡的稳健半监督学习。** 针对类别比例失衡的挑战, 本文提出了一种稳健处理少数类标注的方法, 面向单标注和多标注数据分别设计 AUC 优化和标注分离策略, 避免少数类样本过少导致泛化性能下降。本文方法在现实工业界应用场景, 网约车智能评价和智能判责任务中成功落地转化, 显著提升了业界既有做法的稳健性。

关键词: 机器学习; 开放环境; 半监督学习; 稳健性; 分布失配; 流式数据; 先验不足; 类别失衡

南京大学研究生毕业论文英文摘要首页用纸

THESIS: Semi-Supervised Learning for Open Environment

SPECIALIZATION: Computer Science and Technology

POSTGRADUATE: Lan-Zhe Guo

MENTOR: Associate Professor Yu-Feng Li and Professor Ming Li

ABSTRACT

Semi-supervised learning (SSL) provides a powerful framework for leveraging unlabeled data. Due to its ability to reduce the label requirement of machine learning, SSL has attracted a lot of attention in recent years. Traditional SSL typically assumes closed environments, which require identical distribution, static offline data, adequate prior knowledge, and balanced classes. However, for many practical applications, SSL faces open environments, which are characterized by distribution mismatch, dynamic streaming data, inadequate prior knowledge, and imbalanced classes. To overcome these challenges, this thesis studies SSL for the open environment and proposes a robust SSL scheme. The main contributions are summarized as follows.

1. **Robust SSL for distribution mismatch.** To alleviate the problem of distribution mismatch, this thesis proposes a robust unlabeled example utilization method, which optimizes the model by weighted empirical risk minimization and learns the example weight by bi-level optimization. Theoretically, we prove the convergence rate and analyze the robustness from both empirical and generalization risk. Empirically, our method can improve the robustness of existing methods to distribution mismatch by more than 20%.
2. **Robust SSL for dynamic streaming data.** To alleviate the problem of dynamic streaming data, this thesis proposes a robust dynamic distribution adaptation method, which selects the subset of examples by an influence-based mechanism to satisfy the constrained memory resources and track the changed distributions. The proposal

can be combined with any SSL algorithm. Compared with existing methods, our method can improve the performance by more than 20%.

3. **Robust SSL for inadequate prior knowledge.** To alleviate the problem of inadequate prior knowledge, this thesis proposes a robust SSL model selection method, which integrates multiple candidates and optimizes the performance gain in the worst case. Theoretically, sufficient conditions of robustness are presented, which are much easier to satisfy than previous results. Empirically, all existing methods suffer performance degradations, while our method always achieves robustness.
4. **Robust SSL for imbalanced classes.** To alleviate the problem of imbalanced classes, this paper proposes a robust minority classes exploitation method, which directly optimizes the AUC measure for single-label data and adopts a label separation technique for multi-label data. The proposal improves the robustness of existing methods significantly and has been successfully applied to real-world industrial tasks: ride-sharing intelligent comment and ride-sharing liability judgment.

KEYWORDS: machine learning, open environments, semi-supervised learning, robustness, distribution mismatch, streaming data, inadequate prior knowledge, class imbalance

目 录

中文摘要	i
英文摘要	iii
目 录	v
1 绪论	1
1.1 引言	1
1.2 研究背景	3
1.3 有待研究的问题	10
1.4 本文工作	11
2 适于数据分布失配的稳健半监督学习	15
2.1 引言	15
2.2 相关工作	17
2.3 本文工作	19
2.4 理论分析	29
2.5 实验验证	33
2.6 小结	37
3 适于数据动态流式的稳健半监督学习	39
3.1 引言	39
3.2 相关工作	41
3.3 本文工作	42
3.4 实验验证	48
3.5 小结	59
4 适于先验知识不足的稳健半监督学习	61
4.1 引言	61
4.2 相关工作	62
4.3 本文工作	65
4.4 理论分析	76
4.5 实验验证	79
4.6 小结	94

5 适于类别比例失衡的稳健半监督学习	97
5.1 引言	97
5.2 相关工作	99
5.3 本文工作	101
5.4 实验验证	111
5.5 小结	125
6 结束语	127
6.1 本文工作总结	127
6.2 未来研究展望	128
参考文献	131
致 谢	147
A 攻读博士学位期间学术成果及参加科研项目情况	149
B 攻读博士学位期间获奖及学术活动情况	153

第一章 绪论

1.1 引言

机器学习 (Machine Learning) 致力于研究如何利用经验 (数据) 改善计算机系统性能, 是人工智能 (Artificial Intelligence) 领域的核心技术之一 [166]。二十一世纪以来, 机器学习与人工智能技术得到了越来越多的关注。一方面, 在学术界, 人工智能相关的国际顶级会议参会人数和投稿数量节节攀升, 各高校陆续建立人工智能学院, 开展人工智能与机器学习相关的学科教学。另一方面, 在工业界, 人工智能与机器学习技术促进各传统行业融合创新, 全面提升了制造、农业、物流、金融、商务、家居等领域的智能化水平, 同时也催生出大量新兴产业, 如智能机器人、无人驾驶汽车等。全球主要国家陆续发布文件部署启动人工智能与机器学习相关工作。例如, 2017 年, 我国国务院印发的《新一代人工智能发展规划》中提到“人工智能是引领未来的战略性技术, 世界主要发达国家把发展人工智能作为提升国家竞争力、维护国家安全的重大战略”, 将人工智能的发展提升到国家科技战略层面; 2018 年, 美国发布《国防部人工智能战略》将人工智能与机器学习技术作为未来国防军事的重要发展方向。

典型的机器学习过程包括“输入-模型-输出”三个层面: 输入训练数据 (Training Data), 利用学习算法 (Algorithm) 构建机器学习模型 (Model), 在未见的测试数据 (Testing Data) 上输出正确的预测结果。训练数据通常由特征 (Feature) 和标注 (Label) 两部分组成^①, 例如, 我们希望训练一个“猫”和“狗”的图像分类模型, 则特征表示输入的图片像素, 标注表示每张图片属于猫还是狗。目前机器学习的研究主要关注在监督学习 (Supervised Learning) 场景中, 即收集大量有标注数据作为训练集提供给机器学习算法, 算法通过拟合训练数据得到训练好的模型, 最终在未见环境中测试部署学习得到的模型。然而, 在现实任务中, 数据标注往往耗费大量的人力、物力和财力, 导致大量的数据标注是难以获取的。例如, 在计算机辅助医学图像分析任务中 [167], 如果希望医学专家把影像

^① “Label” 一词也可以翻译成“标记”、“标签”, 与标注同义。

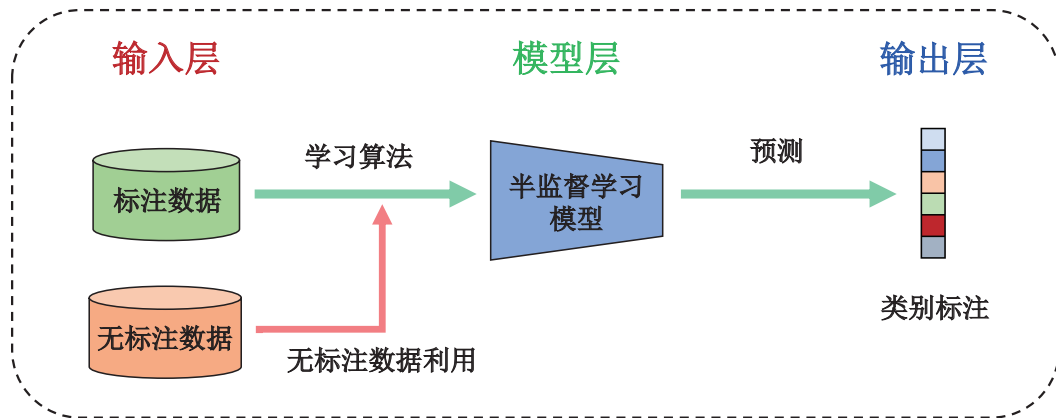


图 1-1: 半监督学习基本流程。

中的病灶全都标识出来是几乎不可能的。但是，与此同时，无标注数据往往是非常容易获得的，可以从医院收集到大量无标注的医学影像。因此，当标注数据不足时，如何利用易于获取的无标注数据提升机器学习泛化性能成为机器学习领域一个关键性的问题。

半监督学习 (Semi-Supervised Learning, SSL) 是利用无标注数据提升机器学习泛化性能的代表性技术。图 1-1 展示了半监督学习的基本流程，输入标注数据和无标注数据，利用半监督学习算法训练模型，使其能够在未见的测试数据上预测出准确的类别标注。自二十世纪九十年代提出以来，半监督学习一直是机器学习领域热门的研究方向。尤其是近年来，随着深度学习的发展，模型复杂度逐渐增加，对标注数据的需求量与日俱增，半监督学习的研究吸引了越来越多学术界和工业界的关注。例如，2017 年以来，在 ICML/NeurIPS/ICLR/CVPR 等人工智能和机器学习顶级会议上，半监督学习相关论文热度高居前三名；在《国家科学评论》(National Science Review, NSR) 2018 年 1 月份出版的机器学习专题期刊中，南京大学周志华教授发表题为《A brief introduction to weakly supervised learning》的综述论文 [157]，指出能够利用无标注数据辅助模型训练的半监督学习是机器学习领域亟待解决的问题；2019 年谷歌首席科学家 Vincent Vanhoucke 发文称“能够利用少量标注数据及大量无标注数据训练机器学习模型的半监督学习革命已经到来”。由此可见，研究半监督学习算法，降低机器学习模型对数据标注的需求，已经得到了学术界和工业界的广泛关注，是人工智能与机器学习领域发展的关键研究问题。

尽管半监督学习近年来取得了巨大的进展，但现有工作大多针对静态、封

闭的学习环境进行研究，依赖一致的数据分布、离线的训练数据、充分的先验知识、平衡的类别比例等假设，然而在现实应用中，机器学习模型所面临的环境往往是动态、开放的，存在数据分布失配、数据动态流式、先验知识不足、类别比例失衡等挑战，从而导致现有半监督学习方法在开放环境中性能不稳健，泛化性能时而提升时而下降，在某些情况下甚至不如只利用少量标注数据的简单监督学习方法，这无疑违背了半监督学习利用无标注数据提升泛化性能的出发点。周志华教授在 2016 年中国计算机大会中发表题为“机器学习：发展与未来”的特邀报告，指出开放环境下的机器学习是机器学习走向实际应用需要解决的共性问题，而稳健性（robustness）是其中的关键要求。AAAI 学会主席 Thomas G. Dietterich 教授在 2016 年国际人工智能大会发表“通往稳健人工智能”（Step towards Robust AI）的主席报告，指出人工智能需要具备稳健性才能进一步应用于现实任务，特别是具有高风险的应用场景 [36]。由此可见，研究开放环境下具有高稳健性的半监督学习理论与方法，是当前半监督学习发展的重要方向。

本章后续内容组织如下：第 1.2 节介绍介绍研究背景，包括半监督学习的基本概念、代表工作，以及开放环境的特点、研究进展；然后，第 1.3 节介绍目前有待研究的问题；最后在第 1.4 节介绍本文工作以及论文的组织结构。

1.2 研究背景

本节介绍相关研究背景，主要包括半监督学习基础知识和代表工作，以及开放环境的特点和当前研究进展。

1.2.1 半监督学习简介

半监督学习是机器学习的一个重要子领域，在给出半监督学习的定义之前，首先回顾一下机器学习的定义。

定义 1-1 （机器学习（Machine Learning） [103, 166]）。对于任务 T 和性能度量 P ，如果一个计算机程序在 T 上其性能 P 随着经验 E 而自我完善，那么我们称这个计算机程序从经验 E 中学习。

上述定义是被机器学习社区广泛采纳的经典定义。该定义表明，一个机器

学习问题由三个重要的部分组成：任务 T 、经验 E 和性能度量 P 。例如，对于图像分类任务 (T ，如猫狗图像分类)，一个机器学习程序能够通过大量标注的图像数据 (E ，如 ImageNet 数据集 [34])，提升图像分类的准确率 (P)；对于语音识别任务 (T ，如手机智能语音助手)，机器学习程序通过在大量采集到的用户语音数据 (E) 上训练，能够提升语音内容识别的准确率 (P)。

传统的机器学习任务，如上述例子提到的，需要大量的标注数据作为监督信息。然而，正如引言中介绍的，在很多现实应用中，大规模标注数据的获取是非常困难的，甚至是不可能的。半监督学习是机器学习领域的一个子研究方向，关注的是当标注数据很少时，如何利用大量容易获取的无标注数据提升机器学习模型的性能。半监督学习的定义如下所示：

定义 1-2 (半监督学习 (Semi-Supervised Learning))。半监督学习是机器学习问题 (由 E 、 T 和 P 定义) 中一个常见的场景，其中 E 包含和任务 T 相关的少量标注数据和大量无标注数据。

例如，对于医学影像分析任务 (T)，我们可以获取大量无标注的医学影像数据，但是标注医学影像需要领域知识，只能请医学专家标注很少一部分，半监督学习通过少量标注的医学影像图片和大量无标注的医学影像图片 (E) 构建机器学习模型，提升医学影像分析任务中病灶检测的准确率 (P)；对于商品推荐任务 (T)，需要请用户标记出对商品是否感兴趣，但很少有用户愿意花时间来提供标注，对于绝大多数的商品，无法得知用户是否感兴趣，半监督学习利用少量有标注的商品数据和大量无标注的商品数据 (E) 构建机器学习模型，提升商品推荐后用户的点击率 (P)。

半监督学习的研究从上个世纪九十年代就已经开始展开 [119]，并随着现实任务中利用无标注数据的需求增加而蓬勃发展，相关的研究工作对机器学习及相关领域产生了许多深远的影响，得到了广泛的关注。例如，国际机器学习大会 (ICML) 从 2008 年开始评选“十年最佳论文”，半监督学习代表性工作，由 Avrim Blum 和 Tom Mitchell 发表于 1998 年国际机器学习理论会议 (COLT) 的文章“Combining labeled and unlabeled data with co-training” [10]，Joachims Thorsten 发表在 1999 年国际机器学习大会 (ICML) 的文章“Transductive inference for text classification using support vector machines” [70] 和 Xiao-Jin Zhu, Zou-Bin

Ghahramani, John D. Lafferty 发表在 2003 年国际机器学习大会 (ICML) 的文章 “Semi-supervised learning using gaussian fields and harmonic functions” [164] 在短短六年中, 先后于 2008 年、2009 年和 2013 年获奖。

接下来, 我们给出半监督学习的形式化描述, 然后介绍代表性的半监督学习算法, 并概述目前进展。

半监督学习：形式化

半监督学习主要考虑同时利用标注数据和无标注数据构建学习模型。形式化而言, 给定标注数据集 $D_l = \{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$ 和无标注数据集 $D_u = \{\mathbf{x}_{n+1}, \mathbf{x}_{n+2}, \dots, \mathbf{x}_{n+m}\}$, 其中 $\mathbf{x} \in \mathbb{R}^d$ 为样本的特征向量 (Feature), $\mathbf{y} \in \{0, 1\}^K$ 为其对应的类别标注 (Label), d 和 K 分别表示样本特征的维度和类别空间的维度。通常来说, $n \ll m$, 即标注样本数量远小于无标注样本。半监督学习的目标是通过拟合给定的数据集 D_l 和 D_u , 训练一个机器学习模型 $f(\mathbf{x}; \theta) : \mathcal{X} \rightarrow \mathcal{Y}$, 使得模型 f 在未知的测试数据上能够输出正确的类别标注, 其中 \mathcal{X} 表示输入特征空间, \mathcal{Y} 表示输出标注空间, θ 表示模型参数。

半监督学习：代表工作

半监督学习的发展经历了两个重要的时期, 统计学习时期 (即, 浅层学习) 和深度学习时期。统计学习时期的半监督学习工作可以分为四大范式: 1) 生成式半监督学习; 2) 半监督支持向量机; 3) 图半监督学习; 4) 基于分歧的半监督学习。下面分别介绍其对应的内容:

- **生成式半监督学习 (Generative Semi-Supervised Learning Methods)**。生成式半监督学习假设标注数据和无标注数据都是通过一个潜在的模型生成的, 因此可以通过潜在模型的参数将无标注数据与类别标注进行关联, 将无标注数据的标注看作模型的缺失参数, 通过期望最大化算法 (Expectation-Maximization Algorithm, EM) 进行极大似然估计求解。生成式半监督学习对潜在模型有不同的假设, 如混合专家模型 [102], 朴素贝叶斯模型 [107] 等, 不同的模型假设将产生不同的半监督学习算法。
- **半监督支持向量机 (Semi-Supervised Support Vector Machine, S3VM)**。半

监督支持向量机基于低密度分割假设 (Low-Density Separation Assumption), 亦称大间隔假设 (Large-Margin Assumption), 即分类器的决策边界应当穿过数据密度较低的区域。因此, 给定标注数据和无标注数据, 半监督支持向量机试图找到能够将有标注样本分开, 且穿过数据低密度区域的分类超平面。半监督支持向量机中最具代表性的是 TSVM (Transductive Support Vector Machine) [70], TSVM 考虑二分类问题, 在对无标注数据所有可能的标注指派中, 选择一个能够使得间隔最大化的模型。

- **图半监督学习 (Graph-based Semi-Supervised Learning, GSSL)**。图半监督学习基于平滑假设 (Smooth Assumption), 即相似的样本应该具有相似的标注。具体而言, 给定一个数据集, 可以将其映射为一个图, 图中的每一个节点代表数据集中的一个样本, 如果两个样本之间的相似度很高, 则对应的节点之间连接一条边, 得到图结构之后, 便可以采用标注传播算法 (Label Propagation) [164] 将标注从有标注数据沿着图中的边传播到无标注数据。
- **基于分歧的半监督学习 (Disagreement-based Semi-Supervised Learning)**。基于分歧的半监督学习采用多个学习器同时进行训练, 首先, 每个学习器基于有标注数据训练得到一个初始模型, 然后, 每个学习器挑选自己置信度最高的无标注样本赋予伪标注 (Pseudo-Label), 并将伪标注样本提供给另一个学习器作为新增的标注样本用于模型更新。这个互相学习、共同进步的过程一直迭代进行, 直到两个学习器都不再发生变化, 或者达到预定的轮数。

近年来, 随着深度学习在各种领域取得了越来越多的成功 [83], 在经典四大半监督学习范式的基础上结合深度模型的深度半监督学习得到了越来越多的关注, 大量深度半监督学习算法被提出, 主要包括:

- **熵最小化方法 (Entropy Minimization)**: 此类方法是低密度假设在深度学习模型上的扩展。低密度假设要求分类器的决策边界不应该穿过数据密度较高的区域, 为了实现该目标, Grandvalet 和 Bengio [55] 提出熵最小化正则项, 显式地优化模型在无标注样本 \mathbf{x} 上的预测结果 $f(\mathbf{x})$ 的熵值, 并从理论上分析熵最小化有助于产生低密度的决策边界; Lee 等人 [84] 进一步提出 Pseudo-Label 算法, 在无标注样本上选择具有高置信度的预测结果作为伪标注, 即预测概率最大的类别标注, 隐式地实现熵最小化的目标。
- **一致性正则方法 (Consistency Regularization)**: 一致性正则是在平滑假设在深

度学习模型上的扩展，通过将数据增广技术应用到无标注数据中，鼓励模型在原始数据的不同增广上产生相似的预测结果。形式化来说，对于无标注样本 \mathbf{x} ，一致性正则最小化如下损失：

$$\|f(\text{Augment}(\mathbf{x}); \theta) - f(\text{Augment}'(\mathbf{x}); \theta)\|_2^2 \quad (1-1)$$

其中 $\text{Augment}(\mathbf{x})$ 和 $\text{Augment}'(\mathbf{x})$ 表示对样本 \mathbf{x} 进行不同的数据增广，例如，对于图像数据，常用的数据增广策略包括旋转、裁剪、随机翻转等。Temporal Ensembling 算法 [118, 80] 直接将上式与标注数据上的监督损失相结合进行优化；Mean Teacher 算法 [128] 进一步将式 1-1 中其中一项的模型输出用基于 EMA (Exponential Moving Average) 策略得到的集成模型输出替代；VAT (Virtual Adversarial Training) 算法 [105] 将数据增广部分替换为对抗的数据增广，即寻找能够最大化输出分布变化的特征增广；UDA 算法 [137] 指出对于不同的数据类型，如图像、文本、语音等，应结合先验知识采用领域相关的增广策略以进一步提升半监督学习性能。

- **混合式方法 (Holistic Methods)**：此类方法综合考虑熵最小化和一致性正则的思想，代表算法如：MixMatch 算法 [7]，首先对样本做基于 MixUp [151] 的数据增广，然后引入锐化函数 (Sharpening)，产生熵最小化的模型预测结果，并将该预测结果作为增广后无标注样本的真实标注进行训练；ReMixMatch 算法 [8] 在 MixMatch 的基础上提出标注分布对齐和强增强、弱增强的数据增广策略，进一步提升了 MixMatch 算法的性能；FixMatch [124] 利用模型在弱增强样本上产生的高置信度的预测结果作为伪标注，并将强增强样本和对应的伪标注作为标注数据，以监督学习的方式进行训练。

1.2.2 开放环境简介

上述半监督学习算法取得成功的前提条件是运行环境是静态的、封闭的，满足数据分布一致、数据静态离线、先验知识充分、类别比例平衡等假设，而在实际问题中，半监督学习所面临的往往是“开放的”、“动态的”环境，这给半监督学习带来了严峻的挑战。具体而言，开放环境的难点主要包括如下几个方面：

数据分布失配：封闭环境下的半监督学习假设标注数据和无标注数据是独

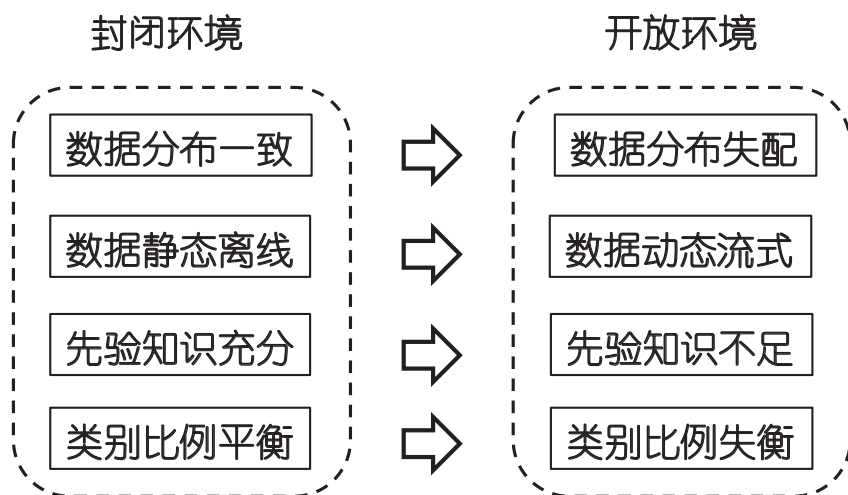


图 1-2: 对比封闭环境, 开放环境下半监督学习的难点。

立同分布采样得到的。然而在实际的开放环境中, 由于无标注样本的收集过程不可能有大量的人工监督, 否则就与半监督学习降低人力需求的出发点相违背, 很容易出现标注数据与无标注数据分布不一致的问题。例如, 在进行猫狗图像分类任务的模型训练时, 无标注数据可以在互联网上根据关键词进行爬取, 在这个过程中, 很容易包含一些任务无关的图片, 如, 狼、老虎等类别的图片, 这些分布外无标注样本严重影响半监督学习的性能。因此, 开放环境下需要半监督学习能够稳健利用无标注数据, 避免分布外样本导致泛化性能下降。

数据动态流式: 封闭环境下的半监督学习假设训练数据是在模型训练开始阶段就已经全部给定的, 模型基于训练数据离线训练。然而在实际的开放环境中, 数据往往以流的方式不断收集得到, 随时间不断积累, 其数据分布随环境不断动态变化, 并且存储资源受限导致不能将所有的样本全部收集起来。例如, 在社交网络情感分析任务中, 每天在社交网络上活跃的用户都会产生大量的数据, 这些数据通常是无标注的, 数据分布会随时间的变化而变化, 考虑到数据的海量性, 显然也无法将所有数据全部存储下来后再进行离线模型训练。因此, 开放环境下需要半监督学习能够在资源受限的条件下稳健适应数据分布的变化, 避免数据动态变化导致泛化性能下降。

先验知识不足: 封闭环境下的半监督学习假设有充分的领域先验知识, 在标注数据规模不足以对模型性能进行准确评估时, 可以依赖先验知识进行稳健的模型选择。然而在实际的开放环境中, 通常难以获取充分的先验知识, 导致无法对多个不确定的半监督学习模型进行准确的选择, 而错误的模型选择非但

不能发挥无标注数据的作用，反而会使模型性能退化，出现性能“不安全”的问题^①。例如，在半监督支持向量机中，可以获得多个满足大间隔假设的学习器，准确的模型选择依赖与数据分布信息和类别标注信息相关的先验知识，在开放环境中，由于环境的动态变化，类似的知识往往是无法获取的。因此，开放环境下需要半监督学习能够在先验知识不充分的条件下稳健选择模型，避免模型选择错误导致泛化性能下降。

类别比例失衡：封闭环境下的半监督学习假设数据的类别比例是平衡的，即每个类别包含相似的样本数量。然而在实际的开放环境中，数据的类别比例往往表现出不平衡的现象，例如，在物种分类任务中，“猫”、“狗”等物种天然比“大熊猫”等物种更为常见；在金融防欺诈检测任务中，非欺诈类别的样本要远远多于欺诈类别的样本。传统的半监督学习算法在面对类别比例失衡的数据时往往只能在多数类上取得良好的性能而在少数类上会发生性能退化的问题。因此，开放环境下需要半监督学习能够稳健处理少数类标注，避免少数类样本过少导致泛化性能下降。

图 1-2展示了相比于封闭环境，开放环境下半监督学习的主要难点。

1.2.3 开放环境下半监督学习研究进展

目前关于半监督学习的研究工作主要集中在静态的封闭环境，对于开放环境的几个难点，只有初步的研究工作。针对数据分布失配的问题，Oliver 等人 [109] 通过实验验证现有的最先进的深度半监督学习方法在标注数据与无标注数据分布不匹配程度达到 40% 时就会发生性能退化的问题，不如只利用标注数据的监督学习方法，但是没有提出针对性的方法解决该问题。针对数据动态流式的问题，[51, 37, 52, 130, 165] 尝试处理流式到来的训练数据，但没有考虑到数据分布会逐渐发生变化的问题，与开放环境的真实应用场景不符；[39] 考虑数据流式且分布动态变化的半监督学习问题，但是没有考虑到现实环境中流式半监督学习会面临存储资源受限的问题。针对先验知识不足的问题，Cozman 等人 [29] 对于生成式半监督学习，提出当先验知识不足以帮助选择与真实数据分布相符的模型时，半监督学习会产生性能不安全的问题；Li 和 Zhou [92] 对于半监督支持向量机，提出当先验知识不足以保障在半监督支持向量机学到的多个大间隔

^①这里的“安全”指的是利用更多无标注样本后能保证性能至少不差于仅利用标注数据的监督学习。

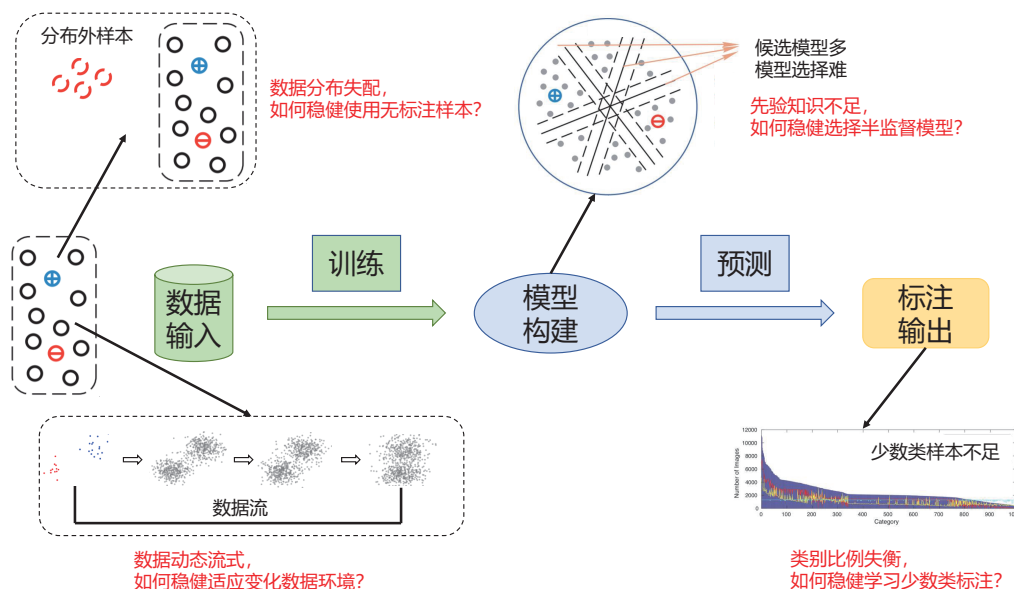


图 1-3: 从输入、模型、输出三个层面考虑开放环境下半监督学习四种典型问题。

分类器中进行可靠的模型选择时，半监督支持向量机会产生性能不安全的问题，并针对性地提出一种综合考虑多个决策边界的集成方法。但这些方法只是针对特定的半监督学习算法进行研究，不适用于通用的半监督学习场景。针对类别比例失衡的问题，目前的研究主要关注在监督学习领域 [71]，半监督学习方向的研究还比较少，并且现有研究主要关注训练数据类别不平衡的问题，忽视了开放环境下测试数据也存在类别不平衡的问题。

1.3 有待研究的问题

综上所述，半监督学习长期以来得到机器学习研究者的广泛关注。然而，以往的研究主要面向静态封闭的学习环境，现实世界的环境往往是开放动态的，既有方法泛化性能时好时坏，难以适用于开放环境。因此，亟需建立开放环境下稳健的半监督学习理论与方法，推动半监督学习在更多真实任务中的应用。

典型的机器学习流程包括数据输入、模型构建、标注输出三个步骤，本文从输入、模型、输出这三个层面作为切入点，研究四种典型问题，如图 1-3 所示。

(1) **适于数据分布失配的稳健半监督学习**: 针对开放环境下标注数据与无标注数据分布不匹配的问题，如何构建适于数据分布失配的稳健半监督学习方法，避免分布外样本导致模型泛化性能下降。

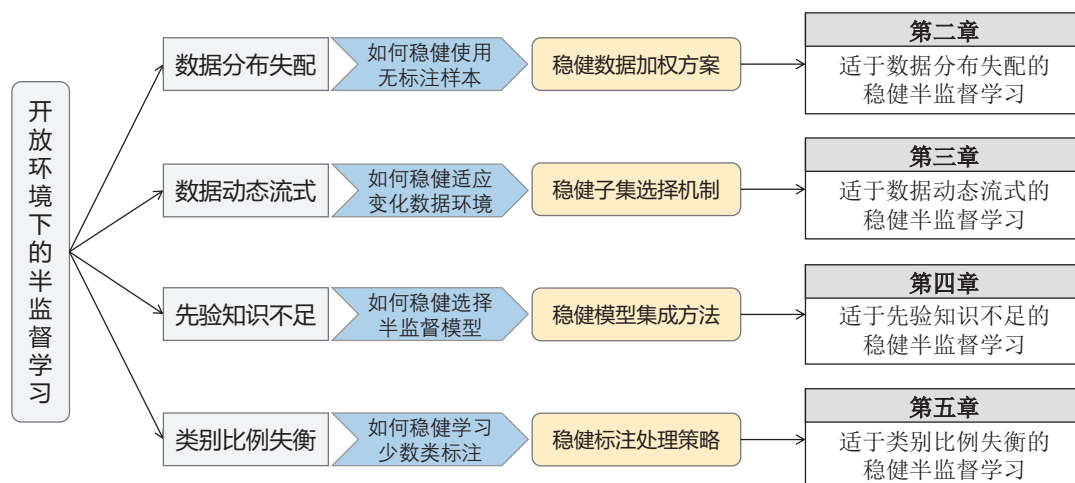


图 1-4: 从输入、模型、输出三个层面，本文对开放环境下半监督学习的四种典型问题进行研究，完成四项工作，分别对应本文的第二章至第五章。

(2) **适于数据动态流式的稳健半监督学习:** 针对开放环境下数据以流的方式不断产生，数据分布随环境不断变化，且存储资源受限的问题，如何构建适于数据动态流式的稳健半监督学习方法，避免环境变化导致泛化性能下降。

(3) **适于先验知识不足的稳健半监督学习:** 针对开放环境下先验知识不充分，不足以提供可靠的半监督学习模型选择的问题，如何构建适于先验知识不足的稳健半监督学习方法，避免模型选择错误导致泛化性能下降。

(4) **适于类别比例失衡的稳健半监督学习:** 针对开放环境下类别比例不平衡，模型在少数类上性能下降的问题，如何构建适于类别比例失衡的稳健半监督学习方法，避免少数类样本不足导致泛化性能下降。

1.4 本文工作

本文从输入、模型、输出三个层面，针对第 1.3 节所列出的开放环境下半监督学习的四种典型问题展开研究，设计适于开放环境的稳健半监督学习方法，对算法进行理论分析并应用到实际任务中，完成四项工作，分别对应本文的第二章至第五章。图 1-4 对本文结构和成果做全局概览。

在第二章中，针对开放环境下标注数据与无标注数据分布不匹配的问题，本文创新性的提出了一种基于数据加权的稳健深度半监督学习方法 DS3L。相比传统的深度半监督学习方法同等地利用所有的无标注数据，DS3L 方法通过赋予分布外样本较低的权重，降低其对模型性能的负面影响。特别地，DS3L 方法将权

重指派过程形式化为一个新颖的双层优化目标，在监督性能的指导下进行无标注样本权重的学习。理论上，DS3L 在分布失配的情况下经验风险不会比只利用标注数据的监督学习方法差，泛化风险能够以 $O(\sqrt{d_w \ln(n)/n})$ 的速率接近最优解，其中 d_w 为权重模型的参数维度，往往取值较小，该结果表明 DS3L 比依赖大模型的深度监督学习方法具有更优的泛化能力。实验上，现有深度半监督学习方法在分布不匹配程度达到 40% 时，相比监督学习就已经出现性能退化的情况，而 DS3L 在分布不匹配程度超过 60% 时依然可以取得性能提升。

在第三章中，针对开放环境下数据流式到来，分布随环境不断动态变化且存储资源受限的问题，本文首次提出了资源受限的流式半监督学习问题设置，在该场景中，标注数据只在初始阶段给出，无标注数据以流的形式不断收集得到且其数据分布逐渐发生变化，由于存储资源的约束，无法将所有时刻的无标注样本全部存储用于模型训练。针对该问题，本文提出了一种基于子集选择的稳健半监督学习方法 Record。Record 在数据流中根据存储资源的约束选择对后续模型训练最有帮助的样本进行存储，并创新性地提出了一种基于影响力机制的样本选择技术，通过计算旧数据分布下的样本在新数据分布上的影响力，选择有助于学习任务的同时与新数据分布最有关联的样本子集，使得数据分布变化能够被有效检测。Record 可结合任意半监督学习方法，实验结果表明，相比于现有方法，Record 在多个流式数据场景上平均分类正确率提升 20% 以上。

在第四章中，针对开放环境下难以获得充分的领域先验知识，无法进行稳健的半监督模型选择的问题，本文提出了一种基于模型集成的稳健半监督学习方法 SafeW。SafeW 同时利用多个模型的预测结果，根据其加权组合构建最终的预测结果，特别地，由于没有充分的先验知识，SafeW 优化潜在最坏情况下的性能提升，得到一个最大最小优化目标。理论上，本文证明了当真实标注可以由多个半监督学习模型预测结果组合得到时，SafeW 可以实现稳健性，该条件比以往理论结果更容易满足。同时，证明了对分类任务和回归任务多种常用损失函数，SafeW 可以将原最大最小优化问题转化为凸优化问题，从而高效获得全局最优解。实验上，现有半监督学习方法均出现了泛化性能下降的问题，部分方法在半数以上的场景中表现不如简单监督学习，而 SafeW 在所有的场景中都优于监督学习性能。此外，SafeW 方法易于扩展以提升更宽泛的弱监督学习稳健性，如不具体监督学习、不精确监督学习等。

在第五章中，针对开放环境下类别比例不平衡，少数类泛化性能下降的问题，针对单标记和多标记数据分别提出了基于 AUC 优化和标注分离的稳健半监督学习方法 CWSL 和 LIMI。CWSL 方法通过直接优化对类别比例不敏感的 AUC 指标，提升半监督学习对类别比例失衡的稳健性。为了解决 AUC 优化非凸非连续的问题，CWSL 引入替代损失函数，在理论上保证与优化 AUC 指标具有一致性的同时给出了高效的优化算法，相比现有方法，在不增加时间开销的前提下，AUC 性能提升 5% 以上。LIMI 方法采用两个分类器分别处理多数类标注和少数类标注，然后利用聚合网络挖掘标注相关性，最后设计无监督损失对无标注数据进行标注补全。大量实验验证 LIM I 在 Hamming Loss、Macro/Micro AUC、Macro/Micro F1 等 9 种多标记学习常用评价指标上，相比既有方法均可实现性能提升。此外，值得一提的是，CWSL 和 LIM I 方法已在现实工业界场景，网约车智能评价和网约车智能判责任务中成功应用。

第二章 适于数据分布失配的 稳健半监督学习

2.1 引言

机器学习，尤其是深度学习，近年来取得了飞速的发展，在图像、语音、文本等场景中获得了成功的应用，研究报告在某些监督学习任务 [83]，例如图像分类任务中，机器学习模型已经可以实现比人类更好的性能。这些机器学习取得成功的任务都有一个基本的条件，即，有大量用于模型训练的标注数据，比如用于图像分类任务的 ImageNet 数据集 [34] 等。然而，在很多现实场景中，标注数据的获取需要耗费大量的人力、物力和财力 [157]，收集大量标注的训练数据是几乎不可能的，这严重限制了机器学习在更多领域的应用。为此，机器学习研究人员提出半监督学习范式，试图利用大量廉价易于获取的无标注数据来帮助机器学习模型提升性能，从而减少对标注数据的需求。近年来，与深度模型结合的深度半监督学习算法已经在多种任务上取得了成功应用，如图像分类 [124]、目标检测 [69]、语义分割 [125]、文本分类 [104] 等。

上述深度半监督学习所取得成功的应用均满足一个基本假设，即，标注样本和无标注样本来自相同的数据分布。然而，由于无标注数据的收集过程缺少人工监督（否则就与半监督学习降低人力成本的出发点相违背），这样的假设在现实开放环境中是很难成立的，其中一种常见的情况是无标注数据中包含标记数据分布外的样本。我们在图 2-1 中以猫狗图像分类任务为例展示了该问题，其中标注数据中全部为与学习任务相关的猫和狗的图片，而无标注数据中既包含了分布内的猫和狗的图片，又包含了“飞机”、“轮船”等与学习任务无关的分布外样本。该现象在真实任务中非常常见，例如，在网页分类任务中 [143]，大量无标注的网页通常是关键词从互联网上爬取的，由于互联网内容的持续更新，很容易爬取到与标注数据分布不匹配的网页数据；在医学图像诊断任务中 [142]，无标注的医学影像中可能会出现标注数据中未出现过的病灶；在图像



图 2-1: 分布不匹配的半监督学习示例: 无标注数据中包含与学习任务无关的分布外样本, 如红色框所示。

分类任务中 [28], 从社交网络或互联网中获取的无标注图像通常会比人工标注的少量标注图像属于更广泛的类别。

有研究指出 [109], 尽管深度半监督学习算法在基准任务上取得了良好的性能, 但是当面对这种包含分布外样本的无标注数据时, 现有深度半监督学习算法不再有效, 甚至可能伴随着严重的性能下降, 导致利用了更多无标注数据的半监督学习模型比只利用少量标注数据的简单监督学习模型性能还要差。这种现象无疑违背了半监督学习引入无标注数据提升模型性能的目标, 严重限制了半监督学习算法在更多实际任务中的有效性。现有的深度半监督学习研究在此类问题上还没有取得良好的结果。

为了解决该问题, 本文提出了一种简单有效的深度半监督学习方法 DS3L。与现有的深度半监督学习算法不同, DS3L 并没有直接使用所有的无标注数据, 而是通过对无标注样本进行赋权, 选择性地使用, 并根据模型在标注数据上的性能优化样本赋权函数, 使模型学到更好的权重。其基本思路是, 利用更优的权重训练得到的模型应该在标注数据上取得更好的性能。整个权重学习的过程被形式为一个新颖的双层优化目标式 [5], 并且提出了高效的迭代优化方法, 证明了其收敛性和收敛速率。DS3L 方法的有效性在理论上和实验中都得到了证明, 具体而言, 在理论上, DS3L 训练得到的模型在经验风险上不会比只利用少量标注数据的监督学习模型差, 并且可以保证其泛化误差界为 $O(\sqrt{d_w \ln(n)/n})$, 比依赖大模型的深度监督学习更优, 在实验上, 当无标注数据中分布外样本比例超过 40% 时, 现有深度半监督学习方法的性能就已经不如简单的监督学习方法, 但是 DS3L 在分布外无标注样本超过 60% 的情况下仍然能实现性能提升。

此外，值得一提的是，DS3L 是通用的深度半监督学习框架，可以与任意深度半监督学习算法进行结合。

2.2 相关工作

本节介绍与本章研究的分布失配的深度半监督学习相关的工作，包括深度半监督学习算法，安全半监督学习算法和分布外样本检测算法，并讨论其与本文研究问题的联系与区别。

2.2.1 深度半监督学习

半监督学习研究当标注数据不足时如何利用无标注数据提升模型性能 [22]。半监督学习已经有很长的研究历史，包括以浅层学习模型为主的统计半监督学习和以深度学习模型为主的深度半监督学习，关于统计半监督学习的相关方法可以参考综述 [22]。本文主要基于深度半监督学习展开，深度半监督学习将经典的半监督学习框架与深度神经网络模型进行结合，在近年来取得了大量的关注 [55, 84, 80, 118, 128, 105, 137, 7, 8, 124]。深度半监督算法可以分为两类，一类方法利用模型生成伪标注，然后计算模型在无标注数据上的预测结果与伪标注之间的监督损失 [84]；另一类方法致力于设计不需要标注信息的无监督正则项，例如，基于熵最小化的正则项 [55]，鼓励模型在无标注数据上预测结果的熵值尽可能小，从而产生置信度尽可能高的预测结果；基于一致性的正则项 [80, 118, 128, 105, 137]，通过对无标注样本进行多次增广，然后鼓励模型针对同一个样本的多次增广版本产生一致的预测结果。此外，还有同时考虑以上多种思想的混合式方法，如 MixMatch [8]，ReMixMatch [7]，FixMatch [124] 等。以上算法在基准的图像分类任务中取得了良好的结果，但是，这些算法均面向封闭环境，没有考虑分布外的无标注样本，导致在数据分布不匹配的半监督学习场景中会出现性能退化的问题。

2.2.2 安全半监督学习

安全半监督学习同样关注在利用了更多无标注数据之后，如何保证模型性能不会比只利用少量有标注数据的监督学习差 [22, 92, 91, 109]。半监督学习不安

全的现象早在 2002 年就已经由 Cozman 等人 [29] 指出, 对于生成式模型, 其成因被认为是当半监督学习模型假设与真实数据分布不符时, 半监督学习会面临性能退化的问题。Loog 等人 [98] 通过优化最坏情况下的似然函数增益, 提升生成式半监督学习模型的安全性。对于半监督支持向量机, 当训练数据中存在多个低密度划分时, 学习算法有可能做出错误的选择, 导致模型性能下降。S4VM [92] 通过优化最坏情况下的性能综合利用多个低密度划分, 提升了半监督支持向量机的安全性。Balsubramani 等人 [4] 证明当无标注数据的真实标注分布位于某个特定的集合中时, 可以通过集成多个分类器的预测结果, 得到安全的半监督学习模型。Li 等人 [90] 针对图半监督学习, 指出图的质量是实现安全半监督学习的关键, 并设计了一种基于大间隔准则的图质量判断标准。但是, 安全半监督学习方法主要关注浅层模型, 并且没有考虑到无标注数据中存在分布外样本的问题, 不能直接用于解决本文研究的数据分布不匹配的半监督学习问题。

2.2.3 分布外样本检测

分布外样本检测方法致力于研究如何使机器学习模型具有识别分布外样本的能力。最简单的分布外样本检测算法直接利用模型预测结果的置信度进行检测, 如果模型在某个样本上的预测置信度低于预先设定的阈值, 则认为该样本是分布外样本 [63]。Liang 等人 [93] 在此基础上, 进一步引入输入样本预处理和输出结果温度缩放 (Temperature Scaling) 技术, 提升了基于预测置信度进行分布外样本检测的能力。Lee 等人 [85] 利用生成对抗网络 (Generative Adversarial Networks, GAN) [54] 生成与学习任务无关的分布外样本, 然后在训练过程中鼓励模型在分布外样本上产生平均的预测概率。Vyas 等人 [129] 提出利用留出法训练多个分类器, 然后对多个分类器进行集成, 利用集成结果进行分布外样本检测。Yu 等人 [147] 提出同时训练两个分类器, 利用分类器预测结果的不一致性进行分布外样本检测, 背后的思想是模型在分布外样本上更容易产生不一致的预测结果。此外, 还有基于能量函数的分布外样本检测方法 [95, 56], 基于等级关系的分布外样本检测方法 [76, 77] 等, 具体可以参考综述 [144]。分布外样本检测与本文研究的面向数据分布不匹配的半监督学习问题的区别主要体现在两个方面。首先, 分布外样本检测算法假定有标注的分布内样本是充足的, 而在半监督学习任务中, 标注数据是非常有限的; 其次, 分布外样本检测算法假设

训练数据中全部都是分布内样本，分布外样本仅出现在测试过程中，而本文研究的分布失配的深度半监督学习需要处理训练数据中出现的分布外样本。

2.3 本文工作

在本节中，我们首先形式化介绍问题设定，然后回顾现有同等利用所有无标注数据的深度半监督学习算法，最后展示我们提出的通过样本赋权，选择性利用无标注数据的深度安全半监督学习方法 DS3L。

2.3.1 问题设定

我们首先对现有深度半监督学习算法的基本概念做简单介绍。在深度半监督学习任务中，训练数据由两部分组成，包含 n 个标注样本的标注数据集 $\mathcal{D}_l = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$ ，以及包含 m 个无标注样本的无标注数据集 $\mathcal{D}_u = \{\mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+m}\}$ ，通常来说， $m \gg n$ 。其中 $\mathbf{x} \in \mathcal{X} \in \mathbb{R}^d$ ， $\mathbf{y} \in \mathcal{Y} = \{0, 1\}^K$ ， \mathcal{X} 表示样本的特征空间， \mathcal{Y} 表示样本的类别空间， d 表示特征空间的维度， K 表示类别空间的维度。半监督学习的目标是通过拟合 \mathcal{D}_l 和 \mathcal{D}_u 学习从特征空间到类别空间的映射模型 $f(\mathbf{x}; \theta) : \{\mathcal{X}; \Theta\} \rightarrow \mathcal{Y}$ ，其中 $\theta \in \Theta$ 表示模型 f 的参数。

用于半监督学习模型训练的损失函数 \mathcal{L} 通常由两部分组成，即标注数据上的监督学习损失 \mathcal{L}_s 和无标注数据上的无监督学习损失 \mathcal{L}_u ，如下所示：

$$\min_{\theta \in \Theta} \sum_{i=1}^n \mathcal{L}_s(f(\mathbf{x}_i; \theta), \mathbf{y}_i) + \sum_{i=n+1}^{n+m} \mathcal{L}_u(\mathbf{x}_i; \theta) \quad (2-1)$$

通常来说，监督学习损失 \mathcal{L}_s 可以通过在标注数据集上计算交叉熵损失 (Cross-Entropy Loss) 得到，如下所示：

$$\begin{aligned} \mathcal{L}_s &= \frac{1}{n} \sum_{i=1}^n H(f(\mathbf{y}|\mathbf{x}_i; \theta), \mathbf{y}_i) \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K -y_{i,k} \log \left(\frac{\exp(f(\mathbf{y} = k|\mathbf{x}_i; \theta))}{\sum_{j=1}^K \exp(f(\mathbf{y} = j|\mathbf{x}_i; \theta))} \right) \end{aligned} \quad (2-2)$$

其中 $f(\mathbf{y}|\mathbf{x}; \theta) \in [0, 1]^K$ 表示当模型参数为 θ 时，模型 f 在输入样本 \mathbf{x} 上的预测概率， $H(\cdot, \cdot)$ 表示交叉熵函数。

无监督损失的构造是深度半监督学习的核心内容，不同的构造方式将产生不同的学习算法。一般来说，有两种构造无监督损失 \mathcal{L}_u 的方式，一种是通过某种策略为无标注数据打上伪标注 \hat{y} ，然后优化样本与伪标注之间的监督学习损失，如下所示：

$$\mathcal{L}_u(\mathbf{x}; \theta) = H(f(\mathbf{y}|\mathbf{x}_i; \theta), \hat{y}) \quad (2-3)$$

代表性的方法如 Pseudo-Label 算法 [84]，其根据模型的预测概率为无标注样本赋予伪标注。

另一种是设计不需标注信息的无监督正则项，代表算法如一致性正则 (Consistency Regularization) [80, 118, 128, 105, 137]、熵最小化正则 (Entropy Minimization) [55, 84] 等。

一致性正则鼓励模型在原始样本的不同增广版本上产生相似的预测结果，如下所示：

$$\mathcal{L}_u(\mathbf{x}; \theta) = \|f(\text{Augment}(\mathbf{x}); \theta) - f(\text{Augment}'(\mathbf{x}); \theta)\|_2^2 \quad (2-4)$$

其中 $\text{Augment}(\mathbf{x})$ 和 $\text{Augment}'(\mathbf{x})$ 表示对样本 \mathbf{x} 进行某种数据增广，如裁剪、平移、翻转、添加随机高斯噪声等 [124]。

熵最小化正则希望模型产生的预测结果具有尽可能小的熵值，即，模型预测的置信度越高越好，如下所示：

$$\mathcal{L}_u(\mathbf{x}; \theta) = - \sum_{k=1}^K f(y = k|\mathbf{x}; \theta) \log(f(y = k|\mathbf{x}; \theta)) \quad (2-5)$$

混合方法 (Holistic Methods) 考虑将以上策略结合起来一起使用，比如 Fix-Match [124] 算法，首先对数据进行弱增广和强增广，并选择在弱增广样本上置信度尽可能高的预测结果作为伪标注，然后计算强增广样本和对应伪标注之间的交叉熵损失，这其中同时用到了熵最小化的思想、模型预测一致性的思想和伪标注的思想。

当标注数据和无标注数据是独立同分布采样得到时，上述深度半监督学习算法取得了良好的结果。例如，图片分类任务上，在大幅减少标注样本数量的情况下，深度半监督学习算法可以取得接近监督学习的分类准确率 [124]。但是，当标注数据与无标注数据分布不匹配时，现有深度半监督学习算法的性能严重

下降，甚至不如简单的只利用有标注数据的监督学习模型 [109]。

2.3.2 DS3L 算法

为了解决标注数据与无标注数据分布不匹配时深度半监督学习算法性能下降的问题，本文提出了安全的深度半监督学习算法 DS3L。不同于以往的深度半监督学习算法同等利用所有的无标注样本，DS3L 算法通过样本赋权选择性的使用无标注数据，希望通过赋予分布外样本较低的权重降低其对模型性能的负面影响。同时，DS3L 根据训练得到的模型在标注数据上的性能指导样本权重的学习，最大化模型性能的提升。

具体而言，一方面，DS3L 引入权重模型 $w(\mathbf{x}; \alpha) : \mathbb{R}^d \rightarrow \mathbb{R}$ ，其中 $\alpha \in \mathbb{B}^{d_w}$ 为权重模型参数， d_w 为权重模型参数空间的维度，输入样本 \mathbf{x} ， $w(\mathbf{x}; \alpha)$ 输出该样本对应的权重。DS3L 试图在样本赋权的条件下，在给定数据集上通过经验风险最小化 [120, 169] 训练得到最优的半监督学习模型，目标如下所示：

$$\hat{\theta}(\alpha) = \min_{\theta \in \Theta} \sum_{i=1}^n \mathcal{L}_s(f(\mathbf{x}_i; \theta), \mathbf{y}_i) + \sum_{i=n+1}^{n+m} w(\mathbf{x}_i; \alpha) \mathcal{L}_u(\mathbf{x}_i; \theta) \quad (2-6)$$

其中 $\hat{\theta}(\alpha)$ 表示当权重模型 w 的参数为 α 时，根据该权重模型训练得到的半监督学习模型 f 的参数。

另一方面，DS3L 通过跟踪模型在监督数据上的性能防止出现性能退化的问题，特别地，DS3L 要求通过上述经验风险最小化得到的模型应该最大化泛化性能，即，

$$\alpha^* = \operatorname{argmin}_{\alpha \in \mathbb{B}^{d_w}} \mathbb{E}_{(X, Y)} [\mathcal{L}_s(f(X; \hat{\theta}(\alpha)), Y)] \quad (2-7)$$

然而，在实际任务中，我们无法得知完整的数据分布信息，无法直接优化上式。采用机器学习中经典的经验风险最小化的思想，DS3L 利用模型在标注数据上的性能近似其泛化性能。具体而言，DS3L 通过优化根据权重模型参数 α 训练得到的模型 $\hat{\theta}(\alpha)$ 在标注数据集上的损失来寻找最优的权重模型参数 α ，即，

$$\hat{\alpha} = \operatorname{argmin}_{\alpha \in \mathbb{B}^{d_w}} \sum_{i=1}^n \mathcal{L}_s(f(\mathbf{x}_i; \hat{\theta}(\alpha)), \mathbf{y}_i) \quad (2-8)$$

为了简化符号，在本章剩余部分我们用 $\hat{\theta}$ 指代 $\hat{\theta}(\alpha)$ 。将上述求解模型参数 θ

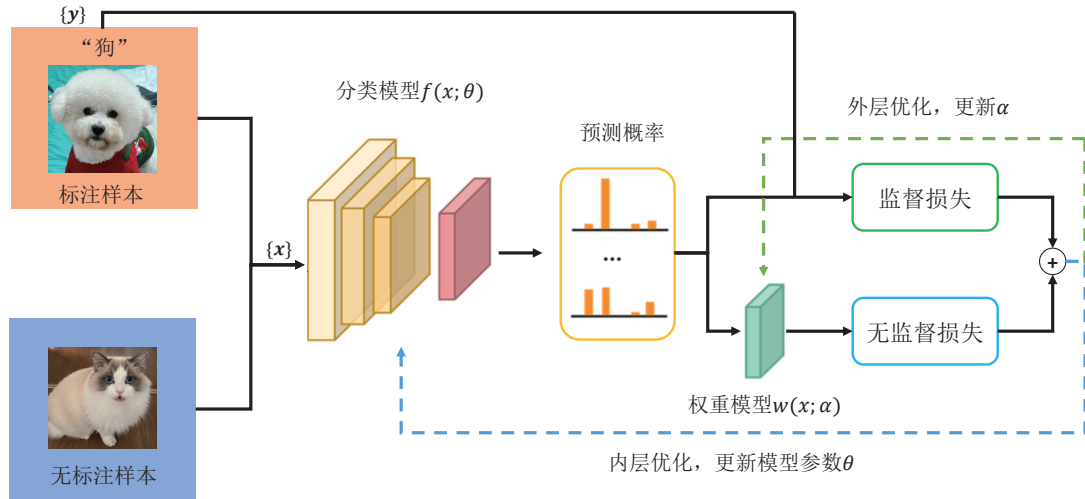


图 2-2: 半监督学习框架 DS3L 的流程。给定标注样本和无标注样本, 计算模型的预测概率, 然后通过该预测概率计算标注数据上的监督损失和样本权重, 并根据权重计算加权的无监督损失。利用监督损失和加权无监督损失优化模型参数 θ , 利用模型在标注数据上的监督损失优化权重参数 α , 直至收敛。

和权重参数 α 的优化目标综合考虑, 可以得到如下的双层优化目标式:

$$\begin{aligned} \min_{\alpha \in \mathbb{B}^{d_w}} \sum_{i=1}^n \mathcal{L}_s(f(\mathbf{x}_i; \hat{\theta}), \mathbf{y}_i) & \quad (2-9) \\ \text{s.t. } \hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \sum_{i=1}^n \mathcal{L}_s(f(\mathbf{x}_i; \theta), \mathbf{y}_i) + \sum_{i=n+1}^{n+m} w(\mathbf{x}_i; \alpha) \mathcal{L}_u(\mathbf{x}_i; \theta) \end{aligned}$$

该优化目标分为两个阶段: 首先, 给定权重模型 $w(\mathbf{x}; \alpha)$, 在内层优化过程中, DS3L 通过优化标注数据损失 $\mathcal{L}_s(f(\mathbf{x}_i; \theta), \mathbf{y}_i)$ 和加权的无标注数据损失 $w(\mathbf{x}_i; \alpha) \mathcal{L}_u(\mathbf{x}_i; \theta)$ 得到经验风险最优的模型参数 $\hat{\theta}$ 。然后, 在外层优化过程中, DS3L 在标注数据集上评估模型 $f(\mathbf{x}; \hat{\theta})$ 的性能, 并更新权重参数 α 。其背后的基本思想是, 利用更优的样本权重训练得到的模型可以在标注数据上取得更好的性能。图 2-2 进一步展示了 DS3L 的基本流程。

2.3.3 优化方法及收敛性证明

如上文所述, DS3L 的优化目标式 2-9 是一个双层优化过程 [5], 其内层优化是给定训练样本及样本权重, 寻找经验风险最小化的半监督学习模型, 外层优化是给定学习模型, 寻找使模型性能最优的权重参数 α 。相比于传统的单层优化, 在双层优化问题中外层优化过程依赖内层优化的结果, 这给双层优化目标的求解带来了挑战。在本节, 我们将介绍 DS3L 方法的优化过程, 并给出收敛性

证明和收敛速率分析。

为了描述的简洁性，我们分别使用 $\mathcal{L}^{outer}(\theta)$ 表示外层优化目标， $\mathcal{L}^{inner}(\theta, \alpha)$ 表示内层优化目标。在深度半监督学习的模型训练中，往往是通过梯度下降算法进行参数求解，因此我们无法得到模型参数 θ 的闭式解 (Closed-Form Solution)，无法直接将 θ 的取值带入外层优化进行求解。传统的求解双层优化的算法包括梯度算法、演化算法等 [123]，但是这些算法计算复杂度非常高，无法直接应用到深度学习模型的训练过程中。为了满足深度学习模型训练对优化算法效率的要求，我们在本节提出了一种近似的双层优化求解策略。

在深度学习模型训练过程中，我们通常采用随机梯度下降或者其变种算法 (如，动量随机梯度下降 [112]，AdaGrad [38]，Adam [75] 等算法) 进行参数 $\hat{\theta}$ 的优化 [53]，其优化过程可以表示为如下公式：

$$\theta_{t+1} = \theta_t - \eta_{\theta} \nabla_{\theta} \mathcal{L}^{inner}(\theta_t, \alpha) \quad (2-10)$$

其中 η_{θ} 表示求解 θ 的优化算法的学习率 (也称为步长)， t 表示当前正处在优化过程的第 t 轮。

在得到最优模型参数 $\hat{\theta}$ 之后，我们计算模型 $f(\hat{\theta})$ 在监督数据上的损失，并根据该损失优化权重模型参数 α ，如下所示：

$$\alpha_{t+1} = \alpha_t - \eta_{\alpha} \nabla_{\alpha} \mathcal{L}^{outer}(\hat{\theta}) \quad (2-11)$$

其中 η_{α} 表示求解 α 的优化算法的学习率。

然而，上述过程依赖两层循环，即，对于每一步 α_t 的求解，都需要计算最优的模型参数 $\hat{\theta}$ 。如果每一次优化的迭代轮数为 T 的话，一共需要 $T \times T$ 次迭代。对于大规模数据集和深度学习模型来说，这个计算复杂度是难以接受的。为了进一步提升优化的效率，我们采用了如下近似优化的方法，相比于利用最优的模型参数优化权重参数 α ，DS3L 迭代的进行模型参数 θ 和权重参数 α 的更新。

模型参数 θ 的更新。 在第 t 步优化过程中，给定权重参数 α_t ，模型参数 θ_{t+1} 可以利用如下梯度下降过程来进行更新：

$$\theta_{t+1} = \theta_t - \eta_{\theta} \nabla_{\theta} \mathcal{L}^{inner}(\theta_t, \alpha_t) \quad (2-12)$$

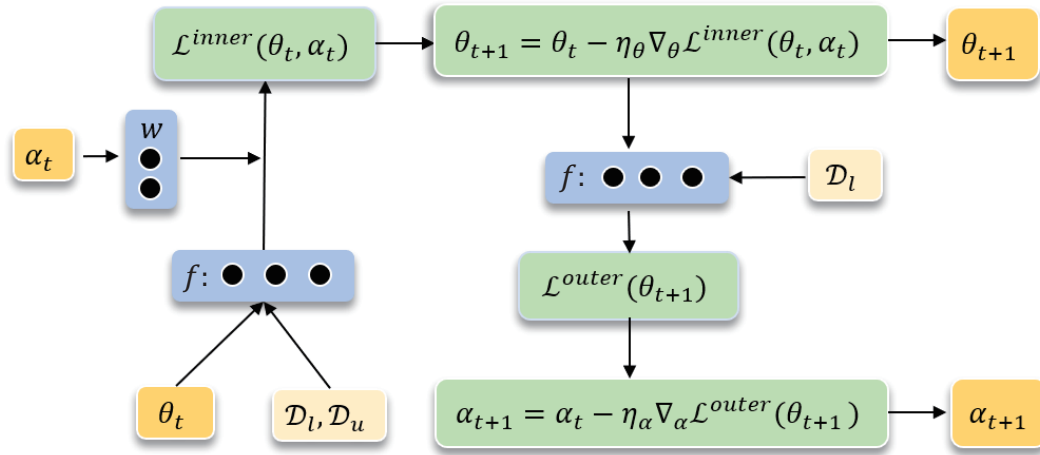


图 2-3: DS3L 的优化步骤。

权重参数 α 的更新。在上一步得到模型参数 θ_{t+1} 之后，我们可以直接计算该模型在标注数据上的损失，并通过如下步骤更新权重参数 α ：

$$\alpha_{t+1} = \alpha_t - \eta_\alpha \nabla_\alpha \mathcal{L}^{outer}(\theta_{t+1}) \quad (2-13)$$

式 2-13 的主要难点在于计算 ∇_α 需要计算双层梯度，利用链式法则，我们推导出双层梯度的计算过程如下所示：

$$\begin{aligned} & \nabla_\alpha \mathcal{L}^{outer}(\theta_{t+1}) \\ &= \nabla_\alpha \mathcal{L}^{outer}(\theta_t - \eta_\theta \nabla_\theta \mathcal{L}^{inner}(\theta_t, \alpha_t)) \\ &= \nabla_\theta \mathcal{L}^{outer}(\theta_t) (-\eta_\theta \nabla_\alpha \nabla_\theta \mathcal{L}^{inner}(\theta_t, \alpha_t)) \end{aligned} \quad (2-14)$$

此外，值得一提的是，上述双层梯度的计算过程在开源的深度学习框架，如 Pytorch^①，Tensorflow^②中已有实现，在实际的训练过程中我们可以利用开源深度学习框架提供的自动求导工具来进行双层梯度 $\nabla_\alpha \mathcal{L}^{outer}(\theta_{t+1})$ 的计算。

图 2-3 展示了算法优化的核心步骤，算法 2.1 总结了优化过程的伪代码。

复杂度分析。相比与普通的单层梯度优化算法，上述双层优化算法需要额外计算学习模型 f 中的前向传播和反向传播过程以及权重模型 w 中的前向传播和反向传播过程。因此，DS3L 算法的时间复杂度约为普通深度半监督学习算法

^①<https://pytorch.org/>

^②www.tensorflow.org

算法 2.1 DS3L 算法优化流程。

输入： 标注数据集 $\mathcal{D}_l = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$, 无标注数据集 $\mathcal{D}_u = \{\mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+m}\}$, 有标注批训练样本个数 B , 无标注批训练样本个数 μB , 优化轮数 T , 模型参数更新步长 η_θ , 权重参数更新步长 η_α 。

输出： 权重参数 α_T 和模型参数 θ_T 。

- 1: 初始化权重参数 α_0 和模型参数 θ_0
 - 2: **for** $t = 0$ to $T - 1$ **do**
 - 3: 从 \mathcal{D}_l 随机采样得到有标注批训练数据 $\{(\mathbf{x}_b^l, \mathbf{y}_b^l) : b \in (1, \dots, B)\}$
 - 4: 从 \mathcal{D}_u 随机采样得到无标注批训练数据 $\{\mathbf{x}_b^u : b \in (1, \dots, \mu B)\}$
 - 5: 计算监督损失: $\mathcal{L}_s = \frac{1}{B} \sum_{b=1}^B \mathcal{L}_s(f(\mathbf{x}_b^l; \theta_t), \mathbf{y}_b^l)$
 - 6: 计算加权无监督损失: $\mathcal{L}_u = \frac{1}{\mu B} \sum_{b=1}^{\mu B} w(\mathbf{x}_b^u, \alpha_t) \mathcal{L}_u(\mathbf{x}_b^u; \theta_t)$
 - 7: 计算内层优化损失: $\mathcal{L}^{inner}(\theta_t, \alpha_t) = \mathcal{L}_s + \mathcal{L}_u$
 - 8: 更新模型参数: $\theta_{t+1} = \theta_t - \eta_\theta \nabla_\theta \mathcal{L}^{inner}(\theta_t, \alpha_t)$
 - 9: 计算外层优化损失: $\mathcal{L}^{outer}(\theta_{t+1}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_s(f(\mathbf{x}_i; \theta_{t+1}), \mathbf{y}_i)$
 - 10: 计算双层梯度: $\nabla_\alpha \mathcal{L}^{outer}(\theta_{t+1})$
 - 11: 更新权重参数: $\alpha_{t+1} = \alpha_t - \eta_\alpha \nabla_\alpha \mathcal{L}^{outer}(\theta_{t+1})$
 - 12: **end for**
 - 13: 返回 θ_T, α_T 。
-

训练复杂度的 3 倍。

此外，我们研究了上述迭代优化算法的收敛性，得到了如下结论：

定理 2-1 (收敛性) 假设损失函数是 L 利普希茨连续的，令模型参数 θ 的优化步长满足 $\eta_\theta \leq \frac{2G}{L}$, G 为大于 0 的常数，那么，基于本文提出的优化算法，模型在标注数据上的监督损失将会随着训练轮数的增加单调递减，即，

$$\mathcal{L}^{outer}(\theta_{t+1}) \leq \mathcal{L}^{outer}(\theta_t) \quad (2-15)$$

进一步，上式等号成立，当且仅当外层优化目标对权重参数 α 的梯度为 0，即，

$$\mathcal{L}^{outer}(\theta_{t+1}) = \mathcal{L}^{outer}(\theta_t) \quad (2-16)$$

当且仅当

$$\nabla_\alpha \mathcal{L}^{outer}(\theta_t) = 0 \quad (2-17)$$

该定理说明了根据我们的优化算法，DS3L 的优化目标将会单调下降，并且在外层优化目标对权重参数 α 的梯度变为 0 时收敛。

进一步，我们分析了该优化算法的收敛率，结论如下所示：

定理 2-2 (收敛率) 假设损失函数是 L 利普希茨连续的, 并且梯度是 ρ 限制的。令模型参数 θ 的优化步长 η_θ 满足 $\eta_\theta = \min\{1, \frac{k}{T}\}$, $k > 0$ 并且 $\frac{k}{T} < 1$; 权重参数 α 的优化步长满足 $\eta_\alpha = \min\{\frac{1}{L}, \frac{c}{\sqrt{T}}\}$, $c > 0$ 并且 $\frac{\sqrt{T}}{c} \geq L$ 。那么, 本文的优化算法能够以 $O(1/\epsilon^2)$ 的阶数实现 $\mathbb{E}[\|\nabla_\alpha \mathcal{L}^{outer}(\theta_t)\|_2^2] \leq \epsilon$, 即:

$$\min_{0 \leq t \leq T} \mathbb{E}[\|\nabla_\alpha \mathcal{L}^{outer}(\theta_t)\|_2^2] \leq O\left(\frac{1}{\sqrt{T}}\right) \quad (2-18)$$

该定理说明了, 通过调整模型参数 θ 和权重参数 α 的优化步长, DS3L 的优化过程以 $O(\frac{1}{\sqrt{T}})$ 的速率收敛到最优解。

定理 2-1 和定理 2-2 的结论证明了本文提出的迭代优化策略针对复杂的双层优化问题能够高效的收敛。此外, 值得一提的是, $O(\frac{1}{\sqrt{T}})$ 的收敛速率已经是目前双层优化问题最优的收敛速率。接下来我们给出上述两定理的具体证明过程。

2.3.4 定理 2-1 证明

在本节, 我们给出定理 2-1 的证明, 首先, 我们给出证明所需的关于利普希茨连续以及梯度限制的定义。

定义 2-1 (利普希茨连续) 函数 $f(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ 满足 L 利普希茨连续, 如果,

$$\|\nabla f(x_1) - \nabla f(x_2)\| \leq L\|x_1 - x_2\|, \quad \forall x_1, x_2 \in \mathbb{R}^d \quad (2-19)$$

定义 2-2 (梯度限制) 函数 $f(x)$ 的梯度是 ρ 限制的, 如果,

$$\|\nabla f(x)\| \leq \rho, \quad \forall x \in \mathbb{R}^d \quad (2-20)$$

接下来, 我们给出具体的证明过程。

证明: 首先, 损失函数满足如下性质:

$$(\nabla_\theta \mathcal{L}^{outer}(\theta))^\top \nabla_\theta \mathcal{L}^{inner}(\theta, \alpha) \geq G \|\nabla_\theta \mathcal{L}^{inner}(\theta, \alpha)\|_2^2 \quad (2-21)$$

其中 $G > 0$ 。

因为在半监督学习的梯度下降过程中，整体半监督损失的梯度方向与监督损失一致，所以上式显然成立。

在优化过程的第 t 步到第 $t + 1$ 步，优化目标的变化如下所示：

$$\begin{aligned}
& \mathcal{L}^{outer}(\theta_{t+1}) - \mathcal{L}^{outer}(\theta_t) & (2-22) \\
& = \mathcal{L}^{outer}(\theta_t - \eta_\theta \nabla_\theta \mathcal{L}^{inner}(\theta_t, \alpha_t)) - \mathcal{L}^{outer}(\theta_t) \\
& \leq -\eta_\theta (\nabla_\theta \mathcal{L}^{outer}(\theta_t))^\top \nabla_\theta \mathcal{L}^{inner}(\theta_t, \alpha_t) + \frac{L}{2} \|\eta_\theta \nabla_\theta \mathcal{L}^{inner}(\theta_t, \alpha_t)\|_2^2 \\
& \leq \left(\frac{L\eta_\theta^2}{2} - \eta_\theta G \right) \|\nabla_\theta \mathcal{L}^{inner}(\theta_t, \alpha_t)\|_2^2 \\
& \leq 0
\end{aligned}$$

其中第一个不等式成立的原因是损失函数满足 L 利普希茨连续性，第二个不等式成立的原因是在半监督学习优化过程中，整体的半监督学习损失与监督学习损失梯度方向一致，第三个不等式成立的原因是 $0 \leq \eta_\theta \leq \frac{2G}{L}$ 。上述结论完成了定理 2-1 的证明。 \square

2.3.5 定理 2-2 证明

在本节，我们给出 DS3L 优化算法收敛率，即，定理 2-2 的证明。

证明： 为了证明过程的简略，我们使用符号 $g(\theta_t, \alpha_t) = \theta_t - \eta_\theta \nabla_\theta \mathcal{L}^{inner}(\theta_t, \alpha_t)$ 代替梯度优化中模型参数更新的过程。根据 DS3L 的优化步骤，我们有，

$$\begin{aligned}
& \mathcal{L}^{outer}(\theta_{t+1}) - \mathcal{L}^{outer}(\theta_t) & (2-23) \\
& = \mathcal{L}^{outer}(g(\theta_t, \alpha_t)) - \mathcal{L}^{outer}(g(\theta_{t-1}, \alpha_{t-1})) \\
& = \{ \mathcal{L}^{outer}(g(\theta_t, \alpha_t)) - \mathcal{L}^{outer}(g(\theta_{t-1}, \alpha_t)) \} \\
& \quad + \{ \mathcal{L}^{outer}(g(\theta_{t-1}, \alpha_t)) - \mathcal{L}^{outer}(g(\theta_{t-1}, \alpha_{t-1})) \}
\end{aligned}$$

针对上式第一项，我们有：

$$\begin{aligned}
& \mathcal{L}^{outer}(g(\theta_t, \alpha_t)) - \mathcal{L}^{outer}(g(\theta_{t-1}, \alpha_t)) & (2-24) \\
& \leq \langle \nabla_\theta \mathcal{L}^{outer}(g(\theta_{t-1}, \alpha_t)), g(\theta_t, \alpha_t) - g(\theta_{t-1}, \alpha_t) \rangle + \frac{L}{2} \|g(\theta_t, \alpha_t) - g(\theta_{t-1}, \alpha_t)\|_2^2
\end{aligned}$$

$$\begin{aligned}
&\leq \eta_\theta \rho^2 + \frac{L}{2} \eta_\theta^2 \rho^2 \\
&= \eta_\theta \rho^2 \left(\frac{\eta_\theta L}{2} + 1 \right)
\end{aligned}$$

其中第一个不等式成立的原因是损失函数满足 L 利普希茨连续性，第二个不等式成立的原因是损失函数的梯度是被 ρ 限制住的。

对第二项，我们根据损失函数 \mathcal{L}^{outer} 对参数 α 的利普希茨连续性，即，

$$\|\nabla_\alpha \mathcal{L}^{outer}(g(\theta, \alpha_t)) - \nabla_\alpha \mathcal{L}^{outer}(g(\theta, \alpha_{t+1}))\| \leq L \|\alpha_t - \alpha_{t+1}\|, \quad \forall t \quad (2-25)$$

可以得到如下结论：

$$\begin{aligned}
&\mathcal{L}^{outer}(g(\theta_{t-1}, \alpha_t)) - \mathcal{L}^{outer}(g(\theta_{t-1}, \alpha_{t-1})) \quad (2-26) \\
&\leq \langle \nabla_\alpha \mathcal{L}^{outer}(g(\theta_{t-1}, \alpha_{t-1})), \alpha_t - \alpha_{t-1} \rangle + \frac{L}{2} \|\alpha_t - \alpha_{t-1}\|_2^2 \\
&= -\left(\eta_\alpha - \frac{L}{2} \eta_\alpha^2 \right) \|\nabla_\alpha \mathcal{L}^{outer}(\theta_t)\|_2^2
\end{aligned}$$

将上述两项合并，我们可以得到：

$$\begin{aligned}
&\mathcal{L}^{outer}(\theta_{t+1}) - \mathcal{L}^{outer}(\theta_t) \quad (2-27) \\
&\leq \eta_\theta \rho^2 \left(\frac{\eta_\theta L}{2} + 1 \right) - \left(\eta_\alpha - \frac{L}{2} \eta_\alpha^2 \right) \|\nabla_\alpha \mathcal{L}^{outer}(\theta_t)\|_2^2
\end{aligned}$$

将上述不等式从 $t = 1$ 到 $t = T$ 相加并化简，我们可以得到：

$$\begin{aligned}
&\sum_{t=1}^T \left(\eta_\alpha - \frac{L}{2} \eta_\alpha^2 \right) \|\nabla_\alpha \mathcal{L}^{outer}(\theta_t)\|_2^2 \quad (2-28) \\
&\leq \mathcal{L}^{outer}(\theta_1) - \mathcal{L}^{outer}(\theta_{T+1}) + \eta_\theta \rho^2 \left(\frac{\eta_\theta L T}{2} + T \right) \\
&\leq \mathcal{L}^{outer}(\theta_1) + \eta_\theta \rho^2 \left(\frac{\eta_\theta L T}{2} + T \right)
\end{aligned}$$

基于上式可以推导得出：

$$\min_t \mathbb{E}[\|\nabla_\alpha \mathcal{L}^{outer}(\theta_t)\|_2^2] \quad (2-29)$$

$$\begin{aligned}
&\leq \frac{\sum_{t=1}^T (\eta_\alpha - \frac{L}{2}\eta_\alpha^2) \|\nabla_\alpha \mathcal{L}^{outer}(\theta_t)\|_2^2}{\sum_{t=1}^T (\eta_\alpha - \frac{L}{2}\eta_\alpha^2)} \\
&\leq \frac{1}{T(2\eta_\alpha - L\eta_\alpha^2)} [2\mathcal{L}^{outer}(\theta_1) + \eta_\theta \rho^2(2T + \eta_\theta LT)] \\
&\leq \frac{1}{T\eta_\alpha} [2\mathcal{L}^{outer}(\theta_1) + \eta_\theta \rho^2(2T + \eta_\theta LT)] \\
&\leq \frac{2\mathcal{L}^{outer}(\theta_1)}{T} \frac{1}{\eta_\alpha} + \frac{\eta_\theta \rho^2(2 + L)}{\eta_\alpha} \\
&= \frac{2\mathcal{L}^{outer}(\theta_1)}{T} \max\left\{L, \frac{\sqrt{T}}{C}\right\} \\
&\quad + \min\left\{1, \frac{k}{T}\right\} \max\left\{L, \frac{\sqrt{T}}{C}\right\} \rho^2(2 + L) \\
&\leq \frac{2\mathcal{L}^{outer}(\theta_1)}{C\sqrt{T}} + \frac{k\rho^2(2 + L)}{C\sqrt{T}} \\
&= O\left(\frac{1}{\sqrt{T}}\right)
\end{aligned}$$

其中第三个不等式成立的原因是 $\eta_\alpha \leq \frac{1}{L}$ ，第四个不等式成立的原因是 $\eta_\theta \leq 1$ 。上述结论完成了定理 2-2 的证明。 \square

2.4 理论分析

在本节，我们给出 DS3L 安全性的理论结果，并展示结果的证明过程。在给出具体的理论证明之前，我们首先给出直观上 DS3L 相比与监督学习以及现有深度半监督学习方法的优越性。

与监督学习方法相比。 监督学习算法直接根据标注数据优化模型参数 θ ，但当标注数据规模较小而模型参数维度较高时，监督学习算法并不能学到泛化性很好的模型。而 DS3L 方法利用标注数据优化权重参数 α ，相比与参数 θ ， α 的参数维度非常小，只需要少量的数据即可学到很好的近似结果。在 DS3L 中，模型参数 θ 是利用标注数据和无标注数据一起训练得到的，无标注数据的引入可以帮助学到相比监督学习更好的参数 θ 。

与现有深度半监督学习方法相比。 以往的深度半监督学习方法认为所有的无标注样本都是同等作用的，在开放环境下，当无标注数据中存在分布外样本时，由于这些样本对学习任务没有帮助，同等的利用所有无标注数据会对模型造成负面影响，导致模型出现性能退化的问题。而 DS3L 通过对无标注样本进行

加权，选择性的利用无标注数据，相比现有的深度半监督学习算法可以减少分布外样本对模型性能的伤害。

接下来，我们从理论上给出 DS3L 性能的具体分析，包括经验风险和泛化风险两个方面。

2.4.1 经验风险分析

本节分析 DS3L 的经验风险，证明了 DS3L 在经验风险上不会比只利用有标注数据的监督学习更差。

定理 2-3 (经验风险分析) 令 θ^{SL} 表示只利用有标注数据训练得到的监督学习模型参数，即， $\theta^{SL} = \arg \min_{\theta \in \Theta} \sum_{i=1}^n \mathcal{L}_s(f(\mathbf{x}_i; \theta), \mathbf{y}_i)$ 。经验风险的定义如下：

$$\hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n [\mathcal{L}_s(f(\mathbf{x}_i; \theta), \mathbf{y}_i)] \quad (2-30)$$

DS3L 训练得到的模型参数 $\hat{\theta}$ 性能永远不会比 θ^{SL} 差，即 $\hat{R}(\hat{\theta}) \leq \hat{R}(\theta^{SL})$ 。

证明：采用反证法进行证明。首先假设 $\hat{R}(\hat{\theta}) > \hat{R}(\theta^{SL})$ ，显然，我们可以通过将无标注样本权重全部设置为 0，得到与监督学习相同的模型，则有 $\hat{R}(\hat{\theta}) = \hat{R}(\theta^{SL})$ ，与假设矛盾。因此， $\hat{\theta}$ 在经验风险上永远不会比监督学习模型 θ^{SL} 更差。 □

定理 2-3 表明，DS3L 在经验风险上可以实现安全性，即，性能不会比只利用标注数据的监督学习模型差，这是以往的深度半监督学习算法均未有的性质。

2.4.2 泛化风险分析

本节分析 DS3L 的泛化风险，证明了 DS3L 具有比依赖大模型的深度监督学习更优的泛化能力。

定理 2-4 (泛化风险分析) 令 $\alpha \in \mathbb{B}^{d_w}$ 表示权重参数 α 的取值空间为 d_w 维的单位球，假设损失函数对参数 α 是利普希茨连续的。泛化风险的定义如下：

$$R(\theta) = \mathbb{E}_{(X,Y)}[\mathcal{L}_s(f(X; \theta), Y)] \quad (2-31)$$

令 $\alpha^* = \arg \max_{\alpha \in \mathbb{B}^{d_w}} R(\hat{\theta}(\alpha))$ 表示单位球中最优的权重参数, $\hat{\alpha} = \arg \max_{\alpha \in \mathcal{A}} \hat{R}(\hat{\theta}(\alpha))$ 表示在候选空间 \mathcal{A} 学到的经验风险最小的权重参数, 以不小于 $1 - \delta$ 的概率, $R(\hat{\theta}(\hat{\alpha})) - R(\hat{\theta}(\alpha^*))$ 被下式限制住:

$$\frac{\left(3\lambda + \sqrt{4d_w \ln(n) + 8 \ln(2/\delta)}\right)}{\sqrt{n}} \quad (2-32)$$

证明: 在证明定理 2-4 之前, 我们先给出证明所需的重要定义。

定义 2-3 (布尔不等式, union bound) 对于任意集合 $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n$, 有:

$$P\left(\bigcup_{i=1}^n \mathcal{A}_i\right) \leq \sum_{i=1}^n P(\mathcal{A}_i) \quad (2-33)$$

定义 2-4 (Hoeffding 不等式) 对于 n 个独立的随机变量 Z_1, \dots, Z_n , 其中 $Z_i \in [0, 1], \forall i$ 。对于所有的 $t \geq 0$, 有:

$$P\left(\frac{1}{n} \sum_{i=1}^n (Z_i - \mathbb{E}(Z_i)) \geq t\right) \leq \exp(-2nt^2) \quad (2-34)$$

以及,

$$P\left(\frac{1}{n} \sum_{i=1}^n (Z_i - \mathbb{E}(Z_i)) \leq -t\right) \leq \exp(-2nt^2) \quad (2-35)$$

定义 2-5 (ϵ 覆盖) 集合 \mathcal{A} 是集合 \mathcal{B} 的一个 ϵ 覆盖, 如果 $\forall \alpha \in \mathcal{B}, \exists \alpha' \in \mathcal{A}$ 满足 $\|\alpha - \alpha'\| \leq \epsilon$ 。

接下来我们给出具体的证明过程, 首先, 令

$$\epsilon = \frac{3}{\sqrt{n}}, \quad \Delta = \frac{\sqrt{2d_w \ln(3/\epsilon) + 2 \ln(2/\delta)}}{\sqrt{n}} \quad (2-36)$$

对于任意 α , 根据 Hoeffding 不等式, 我们有,

$$\begin{aligned} P\{|\hat{R}(\hat{\theta}(\alpha)) - R(\hat{\theta}(\alpha))| > \Delta\} &\leq 2 \exp\left(-\frac{n\Delta^2}{2}\right) \\ &= \frac{\delta}{(3/\epsilon)^{d_w}} \end{aligned} \quad (2-37)$$

令 \mathcal{A} 为 \mathbb{B}^{d_w} 的一个 ϵ 覆盖，则我们有：

$$|\mathcal{A}| \leq (1 + 2/\epsilon)^{d_w} \leq (3/\epsilon)^{d_w}. \quad (2-38)$$

根据布尔不等式，对于 \mathcal{A} 中的所有元素以不小于 $1 - \delta$ 的概率有：

$$\forall \alpha \in \mathcal{A} : |\hat{R}(\hat{\theta}(\alpha)) - R(\hat{\theta}(\alpha))| \leq \sqrt{\frac{2d_w \ln(3/\epsilon) + 2 \ln(2/\delta)}{n}} \quad (2-39)$$

对于 $\forall \alpha' \in \mathcal{A}$ ，我们可以得到：

$$\begin{aligned} R(\hat{\theta}(\hat{\alpha})) &\leq \hat{R}(\hat{\theta}(\hat{\alpha})) + \sqrt{\frac{2d_w \ln(3/\epsilon) + 2 \ln(2/\delta)}{n}} \\ &\leq \hat{R}(\hat{\theta}(\alpha')) + \sqrt{\frac{2d_w \ln(3/\epsilon) + 2 \ln(2/\delta)}{n}} \\ &\leq R(\hat{\theta}(\alpha')) + 2\sqrt{\frac{2d_w \ln(3/\epsilon) + 2 \ln(2/\delta)}{n}} \end{aligned} \quad (2-40)$$

其中第二个不等式成立是因为 $\hat{\alpha} = \arg \min_{\alpha \in \mathcal{A}} \hat{R}(\hat{\theta}(\alpha))$ 。

因为 \mathcal{A} 是 \mathbb{B}^{d_w} 一个 ϵ 覆盖，根据损失函数对参数 α 的利普希茨连续性，对于所有的 $\forall \alpha \in \mathbb{B}^{d_w}$ 我们有：

$$\begin{aligned} R(\hat{\theta}(\hat{\alpha})) &\leq R(\hat{\theta}(\alpha)) + \lambda\epsilon + 2\sqrt{\frac{2d_w \ln(3/\epsilon) + 2 \ln(2/\delta)}{n}} \\ &\leq R(\hat{\theta}(\alpha)) + \frac{(3\lambda + \sqrt{4d_w \ln(n) + 8 \ln(2/\delta)})}{\sqrt{n}} \end{aligned} \quad (2-41)$$

化简上式即可完成定理 2-4 的证明。 \square

定理 2-4 表明 DS3L 算法学到的样本权重能够以 $O(\sqrt{d_w \ln(n)/n})$ 的速率接近最优样本权重，其中 d_w 表示权重模型 w 的参数维度，其取值通常较小。而在标注数据上训练监督学习模型的收敛率为 $O(\sqrt{d_\theta \ln(d_\theta) \ln(n)/n})$ [120]，其中 d_θ 表示模型参数 θ 的维度。值得一提的是，对于神经网络而言， d_θ 往往是百万级的，而在 DS3L 算法中 d_w 的取值通常不超过 100，因此 d_w 远小于 d_θ 。以上结果表明 DS3L 算法能够比依赖大模型的深度监督学习算法以更快的速率收敛。

2.5 实验验证

在本节中，我们在深度半监督学习基准数据集 MNIST 和 CIFAR-10 上进行实验，以进一步评估本文所提出的 DS3L 算法的有效性。

2.5.1 数据信息

我们在深度半监督学习的基准数据集 MNIST 和 CIFAR-10 上进行实验，这两个数据集在深度半监督学习的研究中被广泛采用 [109]，可以有效评估深度半监督学习算法的性能。

MNIST 是用于手写数字识别任务的基准数据集，其中包括 60,000 个训练样本和 10,000 个测试样本，每个样本为 28×28 的手写数字图像。MNIST 数据集一共包含 10 个类别，即数字“1”到数字“10”。在我们的实验中，我们基于 MNIST 构建一个数字 1 到数字 6 的六分类任务数据集，具体而言，我们从类别 1 到类别 6 中分别选取 10 个样本作为标注数据，总计 60 个标注样本，同时，为了构造包含分布外样本的无标注数据集，我们从类别 1 到 10 中选取 30,000 样本作为无标注数据，通过变化来自类别 1-6 的样本的比例来模拟标注数据和无标注数据分布不匹配的程度。比如，当数据分布不匹配程度为 0% 时，意味着所有的无标注样本均来自类别 1-6，和标注数据集一致；当数据分布不匹配程度为 50% 时，意味着有一半的无标注样本来自类别 1-6，其余样本来自类别 7-10。此外，测试数据中只包含来自类别 1-6 的样本。

CIFAR-10 数据集是用于图像分类任务的基准数据集，其中包括了 60,000 个训练样本和 10,000 个测试样本，每个样本为 32×32 的图像。CIFAR-10 数据集一共包含 10 个类别，分别为：“飞机”、“汽车”、“船”、“卡车”、“鸟”、“猫”、“鹿”、“狗”、“青蛙”、“马”，其中四个类别属于交通工具类，其余六个类别属于动物类。在我们的实验中，我们基于 CIFAR-10 构造了一个六类动物图像分类的任务，具体而言，我们在每个动物类别中分别选取 400 个样本作为标注数据集，总计 2,400 个标注样本，同时，为了构造包含分布外样本的无标注数据集，我们从所有的 10 个类中选择 20,000 个样本作为无标注数据。与 MNIST 数据集的构造方式相似，我们通过变化来自动物类别和交通工具类别的无标注样本的比例来模拟数据分布不匹配程度的变化。

2.5.2 对比算法

为了展示我们所提出的 DS3L 算法的有效性，我们将 DS3L 算法与下列主流的深度半监督学习算法进行对比：

- Pseudo-Labeling 算法 [84]: Pseudo-Labeling 算法在神经网络的训练过程中利用网络的预测结果为无标注样本赋予伪标注，如果伪标注的预测置信度大于预先设定的阈值，则将伪标注样本添加至标注数据集中进行训练。
- Π -Model [118]: Π -Model 算法采用一致性正则作为无监督损失函数，最小化模型在原始样本和随机数据增广后的样本上预测概率之间的均方误差。
- Temporal Ensembling [80]: 在 Π -Model 算法的基础上，Temporal Ensembling 算法进一步对模型输出的预测概率做集成，即，计算多轮模型预测结果的均值，然后计算模型在原始样本和数据增广后样本上的集成结果之间的均方误差，通过引入集成的操作进一步提升了 Π -Model 算法的稳定性。
- Mean Teacher [128]: 相比于 Temporal Ensembling 算法对模型的预测值做集成，Mean Teacher 算法对模型参数进行集成，通过引入 EMA 操作 (Exponential Moving Average)，得到多轮训练产生的模型参数的集成，然后对于输入样本计算当前轮的模型和集成模型预测概率之间的均方误差。
- Virtual Adversarial Training (VAT) [105]: VAT 算法不再对数据进行随机的增广操作，而是求解一个在保证样本类别不变的情况下，使样本输出分布变化最大的特征扰动，然后计算模型在原始输入图像和添加该扰动后的图像上预测概率之间的均方误差。

此外，我们也将 DS3L 算法和上述深度半监督学习方法与只利用有标注数据的深度监督学习方法进行对比。

2.5.3 实验设置

对于 MNIST 数据集，我们采用了一个两层的神经网络作为分类器模型，网络包含两个卷积层，其尺寸分别为 $1 \times 16 \times 3$ 和 $16 \times 32 \times 3$ ，两个池化层，其尺寸为 3，步长为 2，填充值为 1。我们采用 ReLu 函数作为网络的激活函数。监督学习损失为标注数据上的交叉熵损失，无监督数据的损失为模型在原始输入图片及添加高斯噪声扰动的图片上预测结果的均方误差 (Mean Squared Error)。该

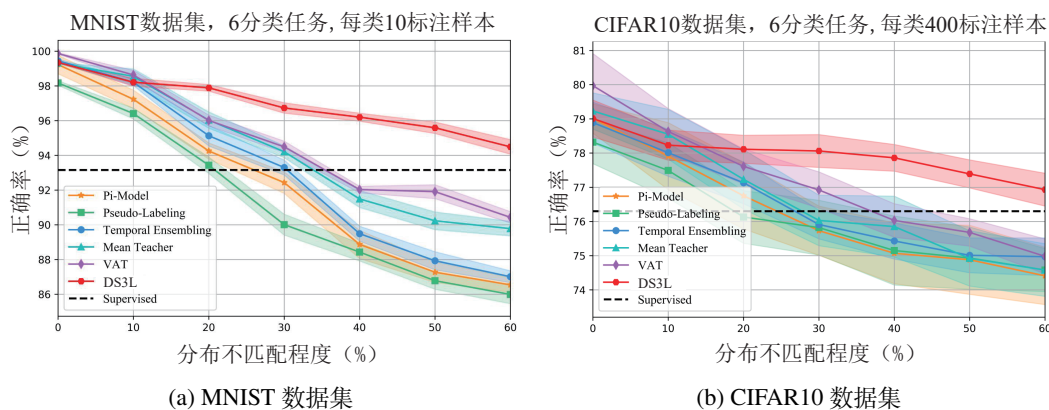


图 2-4: 在 MNIST 和 CIFAR10 数据集中, 随类别分布不匹配程度的变化, DS3L 与对比算法的性能变化。

网络采用随机梯度下降算法 (SGD) 进行优化, 学习率为 1×10^{-3} , 训练轮数为 200,000, 批处理数据的规模 (Batch Size) 为 100。

对于 CIFAR10 数据集, 我们采用 Wide ResNet-28-2 [148] 作为分类器模型。我们对输入的图片做了全局的对比正则化 (Global Contrast Normalization), 以及 ZCA 正则化, 这两者均为 CIFAR-10 数据集上常用的预处理操作, 此外, 我们采用了随机水平翻转、随机像素平移的数据增广操作。类似于 MNIST 数据集, 监督学习损失为标注数据上的交叉熵损失, 无监督数据损失为模型在原始输入图片及添加高斯噪声扰动的图片上预测结果的均方误差。该网络采用随机梯度下降的变种算法 Adam 进行优化, 学习率为 3×10^{-4} , 权重衰减率为 0.2, 训练轮数为 400,000, 批处理数据的规模为 100。

2.5.4 实验结果

MNIST 数据集实验结果。 在 MNIST 数据集上对比方法在测试集上的分类正确率 (Accuracy) 随类别分布不匹配程度变化的实验结果如图 2-4(a)所示。从图 2-4(a)我们可以发现, 当标注数据集和无标注数据集数据分布相同时, 所有的深度半监督学习方法都明显优于基线的监督学习方法, 这是因为半监督学习引入了更多的无标注数据帮助提升性能。然而, 随着数据分布不匹配程度的增加, 现有深度半监督学习算法的性能严重下降, 例如, 当有 40% 的无标注样本为分布外样本时, 许多深度半监督学习算法的性能就已经不如简单监督学习方法的性能, 而我们所提出的 DS3L 算法在分布外无标注样本的比例超过 60% 时仍然能够保持明显的性能提升。

CIFAR10 数据集实验结果。在 CIFAR10 数据集上的实验结果如图 2-4(b)所示，我们可以观察到与 MNIST 数据集上的实验相似的结果。当其它深度半监督学习算法随着数据分布不匹配程度的增加而性能显著下降时，DS3L 仍然能够保持安全的性能提升。以上的实验结果充分验证了在标注数据和无标注数据分布不匹配时，DS3L 方法相比现有深度半监督学习方法的优越性。

DS3L 通用性验证。在上述 MNIS 和 CIFAR10 数据集的实验中，DS3L 采用了最简单的深度半监督学习方法，即，无监督损失为模型在原始输入图像和添加了高斯扰动的图像中输出结果之间的均方误差。值得一提的是，DS3L 是一个通用的半监督学习框架，可以与现有的任意深度半监督学习算法相结合。因此，我们通过进一步将 DS3L 与四种主流半监督学习算法（Pi-Model、Temporal Ensembling、Mean Teacher 和 VAT）相结合，展示了 DS3L 的通用性。实验结果如图 2-5 所示。从图 2-5 我们可以发现，DS3L 与任意半监督学习算法相结合都可以保证性能不退化，这证明了 DS3L 算法的通用性。此外，DS3L 与性能表现更好的半监督学习算法结合可以实现更好的性能。

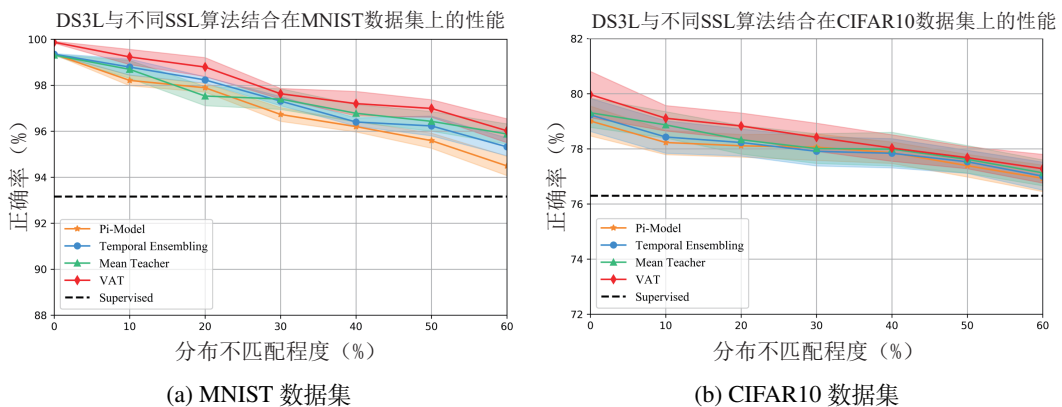


图 2-5: DS3L 与四种不同的深度半监督学习算法相结合在 MNIST 和 CIFAR10 数据集上的实验性能。

分布外样本检测能力验证。虽然 DS3L 算法并非特地为了解决分布外样本检测问题，但是其学得的权重依然具有分布外样本检测的能力，因为分布外样本往往具有比分布内样本更小的权重。为了验证 DS3L 的分布外样本检测能力，我们在 MNIST 和 CIFAR10 数据集上进行实验。我们将 DS3L 与基线的分布外样本检测算法概率筛选法 [63] 进行比较。概率筛选法利用模型对样本的预测概率进行分布外样本检测，当预测概率的最大值低于某个阈值时就认为该样本是分布外样本，其背后的思想是模型在分布内样本上会产生置信度更高的预测结

表 2-1: 分布外样本检测性能。汇报指标为排序误差，数值越低代表检测能力越强，加粗的数字代表表现更好的方法。

数据集	数据分布不匹配程度	概率筛选法	DS3L
MNIST	10%	4.33 ± 0.29	1.67 ± 0.04
	20%	4.78 ± 0.41	0.53 ± 0.23
	30%	4.57 ± 0.33	1.19 ± 0.19
	40%	4.73 ± 0.35	1.50 ± 0.20
	50%	5.67 ± 0.43	2.31 ± 0.13
	60%	7.32 ± 0.51	3.57 ± 0.32
CIFAR-10	10%	7.69 ± 0.67	4.37 ± 0.48
	20%	7.99 ± 0.63	5.34 ± 0.41
	30%	7.67 ± 0.72	5.33 ± 0.43
	40%	8.37 ± 0.75	5.19 ± 0.47
	50%	9.77 ± 0.88	6.51 ± 0.39
	60%	15.03 ± 1.03	10.47 ± 0.78

果。分布外样本检测的评价指标 (1-AUC) 如表 2-1 所示，我们可以发现，所有数据分布不匹配的场景下，DS3L 算法相比于概率筛选法都可以实现更好的性能表现，这也验证了我们所提出的 DS3L 算法可以有效的检测出分布外样本。

2.6 小结

在本章中，我们研究了深度半监督学习标注数据和无标注数据分布不匹配的问题，提出了一种新颖的深度半监督学习方法 DS3L，给出了其经验风险和泛化风险的理论分析，并在两个基准图像分类数据集上进行实验验证了其有效性。DS3L 采用基于双层优化的样本赋权机制，通过对无标注样本进行赋权，降低分布外样本对学习性能的影响。理论上，DS3L 的经验风险不会比只利用少量标注数据的监督学习模型更高，泛化风险比依赖大模型的深度监督学习更优。实验上，目前主流的深度半监督学习方法在数据分布不匹配程度超过 40% 时性能就已经不如简单的监督学习方法，而 DS3L 方法在数据分布不匹配程度超过 60% 时依然可以取得性能提升。此外，DS3L 是通用的学习框架，可以与任意深度半监督学习算法结合以提升其在分布失配情况下的学习性能。

在本章工作的基础上，未来仍有若干方向值得进一步研究。首先，目前的深

度半监督学习研究主要关注图像分类任务，相比与图像这种结构化的数据，现实世界中存在更多复杂的非结构化数据类型，比如欺诈检测、推荐系统等任务中常见的表格数据 [121]，在表格数据上现有深度半监督学习方法常用的数据增广策略不再适用，因此如何在表格数据上训练深度半监督学习模型是仍然是一个具有挑战的任务。其次，本章研究的深度半监督学习方法主要针对的是深度神经网络模型，在深度神经网络之外还存在一些其它的新型深度学习模型，比如基于决策树集成的深度森林模型 [158]，如何面向更广泛的深度学习模型研究分布失配下的学习性能也是一个开放且重要的问题。

本章的主要工作已经成文发表，包括：

- Lan-Zhe Guo, Zhen-Yu Zhang, Yuan Jiang, Yu-Feng Li, Zhi-Hua Zhou. Safe Deep Semi-Supervised Learning for Unseen-Class Unlabeled Data. In: **Proceedings of the 37th International Conference on Machine Learning (ICML'20)**, Virtual Event, pp.3897-3906, 2020. (中国计算机学会 A 类会议，第一作者)
- Zhi Zhou, Lan-Zhe Guo (co-first author), Zhan-Zhan Cheng, Yu-Feng Li, Shi-Liang Pu. STEP: Out-of-Distribution Detection in the Presence of Limited In-Distribution Labeled Data. In: **Advances in Neural Information Processing Systems (NeurIPS'21)**, Virtual Event, pp.29168-29180, 2021. (中国计算机学会 A 类会议，共同一作)

第三章 适于数据动态流式的 稳健半监督学习

3.1 引言

半监督学习的目标是当标注数据不足时利用无标注数据提升机器学习模型的性能，传统的半监督学习方法通常是在离线环境下运行的，即模型在训练时能够获取完整的标注样本和无标注样本集合，并且依赖两个假设：一是数据分布是固定不变的，二是数据的存储资源是不受限制的。然而，在很多现实任务中，这两种假设并不能实现。例如，在汽车自动驾驶任务中 [130] 中，一些有标注的路况信息可以在起始阶段给定的，在汽车运行阶段，会接收到海量无标注的路况数据，这些数据难以被全部存储下来，并且数据的分布显然会随着路况和环境的变化而逐渐变化；在社交媒体情感分析任务中 [9]，初始阶段可以标注少量的情感数据，而网络上每天都会持续产生大量的无标注数据，系统无法将海量数据全部存储，并且相应的数据内容也会随着时间的变化而不断发生变化；在网约车评价任务中 [57]，机器学习模型用少量标注的网约车评价数据进行初始化，然而每天都有大量的无标注数据随着新网约车订单的出现而产生，无法全部存储用于模型训练，其分布也会随时间、区域等因素产生变化。类似的情形也存在于其他在线应用程序中，例如垃圾邮件检测与商品推荐系统等 [39, 79, 130, 155, 170]。

现实开放环境中半监督学习任务的数据往往具有如下特性：首先，无标注数据是连续产生的，且数据分布随时间不断变化，其次，存储资源往往是受限的，不能完整地存储所有的数据用于模型训练。此问题与以往的半监督学习研究有明显的不同，我们将这种新颖的问题设置称为资源受限的流式半监督学习 (Resource Constrained Streaming Semi-Supervised Learning)。据我们所知，该问题还鲜有研究。图 3-1 展示了该问题的设置，该问题的正式设定如下所述：

资源受限的流式半监督学习。在训练初始阶段，即 $t = 0$ 时，给定标注数据集 $\mathcal{L} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$ ，无标注数据以数据流的形式不断收集得到：

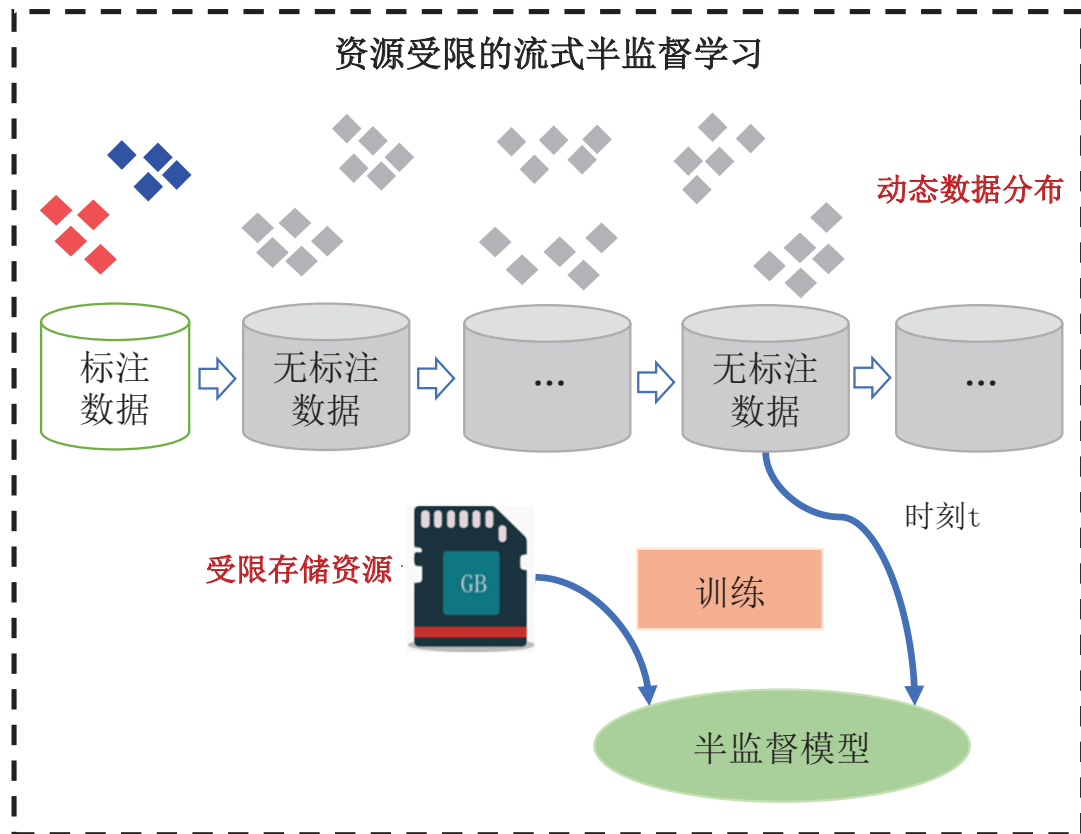


图 3-1: 资源受限的流式半监督学习。

$\mathcal{U}^1, \dots, \mathcal{U}^t, \dots$, 其中 $\mathcal{U}^t = \{\mathbf{x}_1, \dots, \mathbf{x}_{m(t)}\}$, 并且数据分布 $p(\mathbf{x})$ 随环境动态变化。定义 B 为内存预算, 即, 最多仅有 B 个样本可以被存储。学习目标是在每一个时刻 t , 都可以利用存储的 B 个样本与当前到达的无标注数据 \mathcal{U}^t 训练半监督学习模型, 在当前数据分布上表现良好。

以往的半监督学习研究显然无法很好地解决本文提出的问题设置, 因为它们只能在数据分布恒定的封闭静态环境中运行, 无法很好地处理开放环境下分布变化的流式数据。有一些半监督学习研究 [51, 37, 52, 130, 165] 尝试处理流式到来的训练数据, 但没有考虑到数据分布会逐渐发生变化的问题, 与开放世界的真实场景不符; [39] 考虑数据流式到来且数据分布变化的问题, 但是没有考虑到现实环境中流式半监督学习会面临存储资源受限的问题。资源受限半监督学习 [160] 考虑了存储资源受限的条件, 然而仅适用于静态场景, 无法解决流式数据中分布不断变化的问题。

为此, 在本章中, 我们针对上述问题提出了一个系统性的半监督学习方案 Record。Record 方法提出对无标注样本进行子集选择, 根据内存资源的约束选择

在分布变化的场景下对模型训练最有帮助的样本进行存储。特别地，Record 采用一种基于影响力机制的样本选择技术，首先通过模型预测置信度进行分布外样本检测，区分旧分布和新分布下的样本，然后通过计算旧分布下的样本在新的数据分布上的影响力，选择对学习任务最有帮助的同时与新数据分布关联性最强的样本。Record 是通用的半监督学习框架，可以嵌入任意的半监督学习算法，大量的实验结果验证了本文提出的 Record 方法在存储资源受限的流式半监督学习问题中的有效性，例如，与当前最先进的流式半监督学习方法相比，Record 在多个流式数据场景上平均分类正确率提升 20% 以上。

3.2 相关工作

当前绝大多数半监督学习算法是为离线场景设计的，也有一些方法 [51, 37, 52, 39, 68, 130, 165] 尝试解决流式数据的半监督学习问题。例如，[51] 提出了一种在线的流形正则化算法，该算法可以通过求解核空间内的凸规划问题对无标注数据进行学习。[52] 通过结合半监督似然函数和蒙特卡罗框架，提出了一种在线主动半监督学习方法。[68] 基于局部一致性的标注传播算法，提出了一种基于图的在线半监督学习方法。然而，这些方法没有考虑流式数据中伴随的数据分布变化的问题。[88, 37] 考虑了具有分布变化的流式半监督学习问题，然而，这些工作均假设标注数据在每个时刻都可以获取，在现实任务中，由于数据标注的困难性，该假设往往不能满足。[59, 165] 考虑数据类别变化的流式数据，假设数据流中会不断出现模型未见过的类别，其任务目标是识别新类样本并对其进行分类，与本文研究目标有所区别。

资源受限的半监督学习研究也非常有限。[160] 是较早研究资源受限半监督学习的工作，通过提出一种近似算法，使得基于图的半监督学习方法能适应给定的内存资源约束。[40] 采用基于一种密度的度量来进行无标注样本的选择，通过降低样本复杂度减少半监督学习的计算消耗，使其适应计算资源受限的场景。[82] 提出了一种资源受限的半监督支持向量机，该算法利用频谱图中携带的相邻信息和相离信息来更有效地使用内存资源。然而，这些方法主要是在静态离线环境中运行，无法应用于数据分布动态变化的流数据学习环境。

和本文最相关两个工作是 [39] 和 [130]，这两个工作同样考虑标注数据只在

初始阶段给出，并且无标注数据以数据流的形式连续收集到，其数据分布随环境不断变化。具体而言，[39]提出了一种基于几何算法进行无标注样本选择的半监督学习方法 COMPOSE，然而，他们没有考虑存储资源的限制，且所提出的几何算法只在数据特征低维的场景下有效。[130]针对图半监督学习提出了一种具有内存约束的标记传播算法 TLP，然而该方法没有考虑到数据分布变化的问题，并且只能应用到图结构数据中。

3.3 本文工作

在本节，我们提出了一个系统性的解决方案 Record 来处理资源受限的流式半监督学习问题。具体而言，该问题两个主要挑战为：1) 数据分布不断发生变化，如何让模型适应分布变化，避免性能下降；2) 无标注样本无法完整存储，在满足存储资源约束的条件下，如何高效利用无标注样本，最大化半监督学习性能增益。Record 方法针对上述挑战，提供了一个完整通用的解决方案，包括两个核心的技术：基于置信度的分布偏移样本检测和基于影响力机制的样本选择。接下来，我们首先介绍 Record 的总体框架，包括符号和基本设置，然后依次介绍具体的技术细节，最后给出复杂度分析。

3.3.1 总体框架

考虑从输入空间 $\mathcal{X} \in \mathbb{R}^d$ 到输出空间 $\mathcal{Y} \in \mathbb{R}^C$ 的预测问题，其中， d 是特征空间维度， C 是类别空间维度。在训练初始阶段，即 $t = 0$ 时刻，给定包含 n 个标注样本的标注数据集 $\mathcal{L} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ ，其中 $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{X} \times \mathcal{Y}$ 。在后续其它时刻，我们只能接收到无标注数据集 $\mathcal{U}^t = \{\mathbf{x}_i\}_{i=1}^{m^t}$ ，每个时刻无标注样本的数量 m^t 可能存在差异。目标是学习预测模型 $f: \mathcal{X} \rightarrow \mathcal{Y}$ ，能够在每个时刻 t 的数据分布 $p^t(\mathbf{x}, \mathbf{y})$ 都表现出良好的性能。

在本文中，我们考虑时刻 t 的数据分布 $p^t(\mathbf{x})$ 是从时刻 $t-1$ 的数据分布逐渐变化得到的，并且数据分布 $p^t(\mathbf{x})$ 和 $p^{t-1}(\mathbf{x})$ 存在一定的重合。这种假设在实际任务中是合理的，因为完全随机的数据分布变化是不可学习的 [39]。图 3-2 展现了数据分布变化的例子。

为了进一步满足内存资源的约束，Record 方法在每个时刻 t 对无标注样本

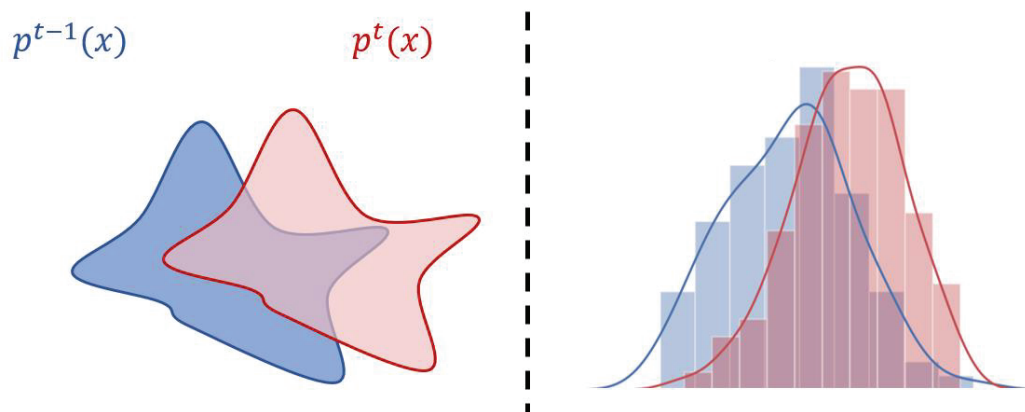


图 3-2: 时刻 $t-1$ 的数据分布 $p^{t-1}(\mathbf{x})$ 和时刻 t 的数据分布 $p^t(\mathbf{x})$ 示例。

进行子集选择，选取对后续时刻最具帮助的样本进行存储，即，维护一个最多包含 B 个样本的数据集 \mathcal{D}_s ，并持续更新。在每个时刻 t ，Record 框架利用存储的数据集 \mathcal{D}_s 和当前时刻收集到的无标注数据 \mathcal{U}^t 训练半监督学习模型。因此，Record 的核心问题是如何选择无标注样本并更新数据集 \mathcal{D}_s ，使其能够在满足内存资源约束的条件下，适应连续的数据分布变化。具体而言，Record 在每个时刻 t 试图解决如下问题：

$$\begin{aligned}
 & \max_{\mathcal{D}_s^t} \text{performance}(f) & (3-1) \\
 & \text{s.t. } f = \mathcal{A}(\mathcal{D}_s^t, \mathcal{U}^t) \\
 & \mathcal{D}_s^t \subseteq \mathcal{D}_s^{t-1} \cup \mathcal{U}^{t-1} \\
 & \text{size}(\mathcal{D}_s^t) \leq B
 \end{aligned}$$

其中 $\text{performance}(\cdot)$ 表示模型的评价指标，如分类正确率， \mathcal{A} 表示某种半监督学习算法， B 表示内存资源约束。

Record 整体框架如图 3-3 所示。在时刻 $t=0$ ，给定少量标注样本，其中不同颜色的图形代表不同的类别，在后续时刻只能接收到无标注样本，并且其分布会随时间逐渐变化。Record 方法的目标是维护数据集 \mathcal{D}_s 辅助进行半监督模型训练，并且在每个时刻 t 的数据分布上都取得良好的性能。接下来，我们介绍 Record 方法的两个核心技术：基于置信度的分布偏移样本检测和基于影响力机制的样本选择。

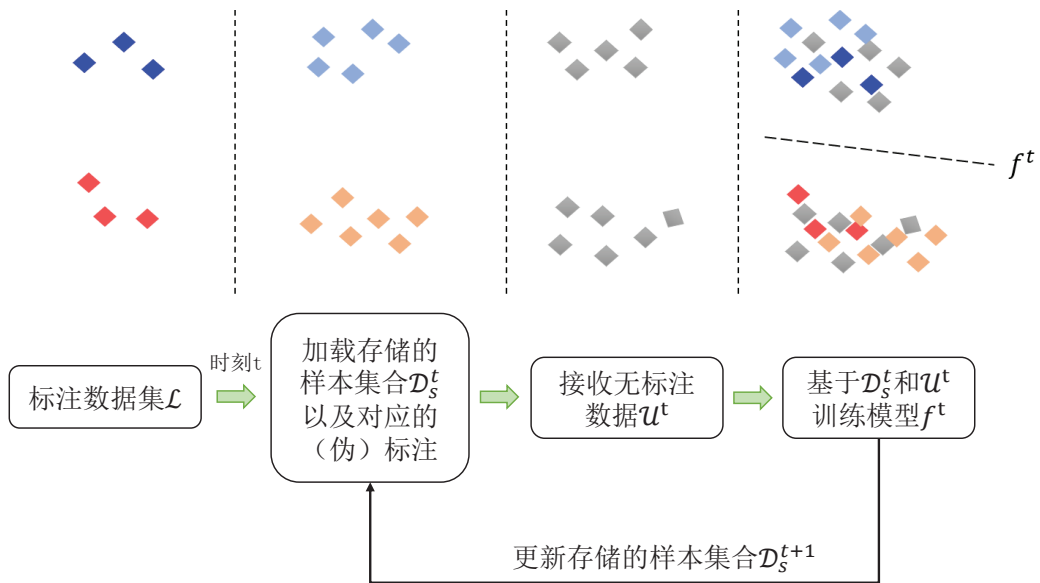


图 3-3: Record 学习框架示意图。

3.3.2 基于置信度的分布偏移样本检测

Record 方法需要解决的一个核心问题是如何根据内存资源的约束选择样本子集进行存储。一种直接的样本选择方法是在每个时刻都随机选择 B 个样本存储，显然，这种随机采样的方法不能取得良好的性能，因为它没有考虑到无标注数据的分布信息，无法保证选择的样本是对学习任务最有帮助的。此外，一种启发式的样本选择方法是根据模型预测结果，选择具有高预测概率的样本，因为预测概率揭示了模型对该样本预测结果的置信度，置信度高的样本往往是对学习任务更有帮助的样本，类似的思想已在分布外样本检测任务中成功应用 [63]。这种方法虽然能够选出对于当前任务中更有价值的样本，但是并不能很好的解决数据分布变化的问题。由于 \mathcal{U}^t 的数据分布 $p^t(\mathbf{x})$ 会随着时间持续变化，我们不仅需要选择对当前分布 $p^t(\mathbf{x})$ 有价值的样本，还需要找到能够反映数据分布变化趋势的样本，使其能够在未来时刻发挥作用。

为了实现上述目标，我们首先需要识别分布偏移的样本，即， $p^{t-1}(\mathbf{x})$ 和 $p^t(\mathbf{x})$ 非重叠区域的数据样本。具体而言，我们使用模型输出的预测概率分布作为指标，因为处于当前分布中的样本往往比处于偏移后数据分布的样本具有更高的预测概率 [64]。将模型的预测概率记为 $f(\mathbf{x}) \in \mathbb{R}^{1 \times C}$ ，其中 $f(\mathbf{x})_c$ 表示后验概率

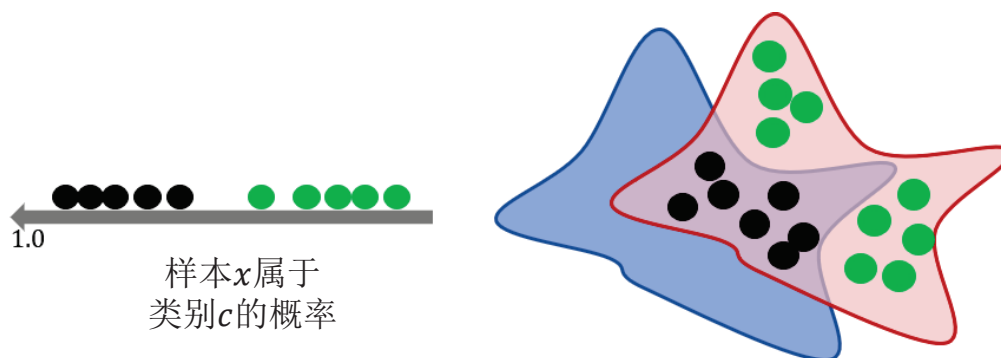


图 3-4: \mathcal{U}_c^{in} 和 \mathcal{U}_c^{out} 的划分过程实例。集合 \mathcal{U}_c^{in} 中的样本用黑色圆形表示，集合 \mathcal{U}_c^{out} 中的样本用绿色圆形表示。 \mathcal{U}_c^{in} 的样本具有更高的预测置信度，因此更有可能落在当前分布。 \mathcal{U}_c^{out} 中的样本预测置信度相对较低，因此更有可能落在偏移后的数据分布。

$p(c|\mathbf{x}, f)$ ，我们可以根据该预测概率为无标注样本赋予伪标注，即：

$$\hat{\mathbf{y}} = \underset{c \in \{1, \dots, C\}}{\operatorname{argmax}} f(\mathbf{x})_c \quad (3-2)$$

定义 \mathcal{U}_c 为 t 时刻的无标注样本集合 \mathcal{U}^t 中伪标注被预测为类别 c 的所有样本的集合，基于模型的预测概率，可以将 \mathcal{U}_c 分为大小相等的两个子集 \mathcal{U}_c^{in} 和 \mathcal{U}_c^{out} 。具体来说， \mathcal{U}_c^{in} 包含具有更高预测置信度的样本，因为这些样本有更大的概率属于当前数据分布，而 \mathcal{U}_c^{out} 包含剩余的样本，这些样本更容易落在偏移后的数据分布中（即相邻时刻数据分布的非重叠区域）。图 3-4 展示了集合 \mathcal{U}_c^{in} 和 \mathcal{U}_c^{out} 的划分过程。

3.3.3 基于影响力机制的样本选择

在上一步得到集合 \mathcal{U}_c^{in} 和 \mathcal{U}_c^{out} 之后，我们进一步考虑如何选择对后续时刻的数据分布最有帮助的样本。基本的思路是采用影响力机制，通过计算 \mathcal{U}_c^{in} 中的样本在集合 \mathcal{U}_c^{out} 上影响力，选取影响力最大的样本进行存储。因为 \mathcal{U}_c^{in} 中的样本具有更高的预测置信度，表示这些样本对学习任务具有较大的帮助，而同时在 \mathcal{U}_c^{out} 上有更高的影响力，表示这些样本能够更好的反映数据分布变化的趋势。

具体而言，我们提出利用影响函数（Influence Function）[78] 来评估 \mathcal{U}_c^{in} 中的样本在集合 \mathcal{U}_c^{out} 上的影响力。影响函数是一种帮助寻找对潜在数据分布最有

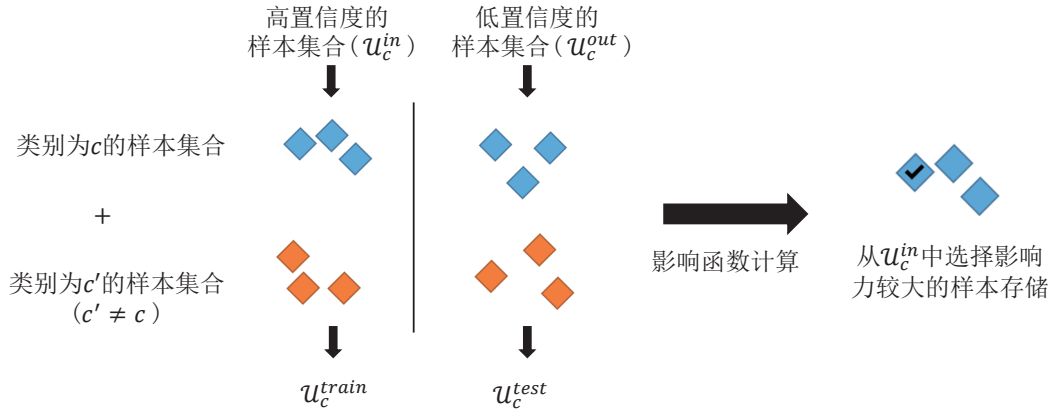


图 3-5: 对于任意类别 c , 基于影响力机制的样本选择过程。

影响的训练样本的有效方法。为了计算影响函数, 首先需要对于每一个类别 c 构造一个训练集和测试集, 我们从其它任意类别 c' , $c' \neq c$ 的样本集合 $\mathcal{U}_{c'}^{in}$ 和 $\mathcal{U}_{c'}^{out}$ 随机采样了 $|\mathcal{U}_c^{in}|$ 和 $|\mathcal{U}_c^{out}|$ 个样本作为负样本, (即, 标注 $\tilde{y} = -1$), 并将集合 \mathcal{U}_c^{in} 和 \mathcal{U}_c^{out} 中的样本作为正样本 (即, 标注 $\tilde{y} = 1$) 来构造一个二分类任务的训练集 \mathcal{U}_c^{train} 和测试集 \mathcal{U}_c^{test} , 如下所示:

$$\mathcal{U}_c^{train} = \mathcal{U}_c^{in} \cup \text{Random}_{|\mathcal{U}_c^{in}|}(\mathcal{U}_{c'}^{in}) \quad (3-3)$$

$$\mathcal{U}_c^{test} = \mathcal{U}_c^{out} \cup \text{Random}_{|\mathcal{U}_c^{out}|}(\mathcal{U}_{c'}^{out}) \quad (3-4)$$

其中 $\text{Random}_n(\mathcal{D})$ 表示从集合 \mathcal{D} 随机采样得到 n 个样本。

然后我们基于训练集 \mathcal{U}_c^{train} 构建一个二分类的逻辑回归 (Logistic Regression) 模型, 并且计算集合 \mathcal{U}_c^{in} 中的样本 \mathbf{x} 在单个测试样本 $\mathbf{z}_{test} = (\mathbf{x}_{test}, \tilde{y})$ 上的影响函数值 $I(\mathbf{x}, \mathbf{z}_{test})$ 。

具体而言, 对于集合 \mathcal{U}_c^{in} 中的样本 \mathbf{x} , 逻辑回归的损失函数为:

$$L(\mathbf{x}, \theta) = \log(1 + \exp(-\theta^\top \mathbf{x})) \quad (3-5)$$

其中 θ 为逻辑回归模型的参数。

基于 [78], 影响函数值 $I(\mathbf{x}, \mathbf{z}_{test})$ 可以由如下公式计算得到:

$$\tilde{y} \sigma(-\tilde{y} \theta^\top \mathbf{x}_{test}) \cdot \sigma(-\theta^\top \mathbf{x}) \mathbf{x}_{test}^\top \mathbf{H}_\theta^{-1} \mathbf{x} \quad (3-6)$$

算法 3.1 Record 算法流程。

输入: 任意半监督学习算法 \mathcal{A} , 标注数据集 $\mathcal{L} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, 存储样本集 $\mathcal{D}_s^0 = \mathcal{L}$ 。

- 1: **for** $t = 0, 1, \dots$ **do**
- 2: 获取当前时刻无标注数据集 $\mathcal{U}^t = \{\mathbf{x}_i\}_{i=1}^{m^t}$
- 3: 根据 \mathcal{D}_s^t 和 \mathcal{U}^t 训练半监督学习模型: $f^t: \mathcal{X} \rightarrow \mathcal{Y}$
- 4: 获得 \mathcal{U}^t 中样本的预测概率: $f^t(\mathbf{x}), \forall \mathbf{x} \in \mathcal{U}^t$
- 5: 根据预测概率指派伪标注: $\hat{\mathbf{y}} = \arg \max_{c \in \{1, \dots, C\}} f^t(\mathbf{x})_c$
- 6: 令 $\mathcal{R} = \emptyset$
- 7: **for** $c = 1, \dots, C$ **do**
- 8: $\mathcal{U}_c = \{\mathbf{x}_i \in \mathcal{U}^t, \text{其中 } \mathbf{x}_i \text{ 对应的伪标注 } \hat{\mathbf{y}}_i = c\}$
- 9: 将 \mathcal{U}_c 划分为集合 \mathcal{U}_c^{in} 和集合 \mathcal{U}_c^{out}
- 10: 从类别 c' 随机采样 $|\mathcal{U}_c^{in}|$ 个样本, 并入集合 \mathcal{U}_c^{in} , 作为训练集 \mathcal{U}_c^{train}
- 11: 从类别 c' 随机采样 $|\mathcal{U}_c^{out}|$ 个样本, 并入集合 \mathcal{U}_c^{out} , 作为测试集 \mathcal{U}_c^{test}
- 12: 运行算法 3.2, 根据样本影响力得到集合 \mathcal{R}_c
- 13: $\mathcal{R} = \mathcal{R} \cup \mathcal{R}_c$
- 14: **end for**
- 15: 根据公式 3-8 更新集合 \mathcal{D}_s^{t+1}
- 16: **end for**

其中 $\sigma(t) = \frac{1}{1+\exp(-t)}$, \mathbf{H}_θ 表示逻辑回归函数的 Hessian 矩阵。

得到样本 \mathbf{x} 在单个测试样本 \mathbf{z}_{test} 上的影响函数值 $I(\mathbf{x}, \mathbf{z}_{test})$ 之后, 样本 \mathbf{x} 的影响力可以由如下公式计算得到:

$$\text{IF}(\mathbf{x}) = \frac{1}{|\mathcal{U}_c^{test}|} \sum_{\mathbf{z}_{test} \in \mathcal{U}_c^{test}} I(\mathbf{x}, \mathbf{z}_{test}) \quad (3-7)$$

即, 在整个测试集 \mathcal{U}_c^{test} 上的平均影响函数值。

最终, 对于任意类别 c , 令 \mathcal{R}_c 表示集合 \mathcal{U}_c^{in} 中影响力为正的样本, 即 $\text{IF}(\mathbf{x}) \geq 0$ 的样本 \mathbf{x} , 如果样本数量 $|\mathcal{R}_c|$ 大于每一类的存储限制 $\frac{B}{C}$, 则选择影响力最大的 $\frac{B}{C}$ 个样本进行存储, 否则, 用 $|\mathcal{R}_c|$ 中的样本替换 \mathcal{D}_s^t 最旧的样本, 即距离当前时刻最远的样本, 如下所示:

$$\mathcal{D}_s^{t+1} = \bigcup_{c=1}^C \left\{ \begin{array}{ll} \text{Top-IF}_{\frac{B}{C}}(\mathcal{R}_c) & |\mathcal{R}_c| \geq \frac{B}{C} \\ \text{Latest}_{\frac{B}{C}-|\mathcal{R}_c|}(\mathcal{D}_s^t) \cup \mathcal{R}_c & |\mathcal{R}_c| < \frac{B}{C} \end{array} \right\} \quad (3-8)$$

存储的样本及其伪标注将用于后续时刻半监督学习模型的训练。图 3-5 展示了对于任意类别 c 进行样本选择的示例, Record 框架的具体步骤总结在算法 3.1 中。

算法 3.2 影响函数计算。

输入: \mathcal{U}_c^{train} , \mathcal{U}_c^{test} , \mathcal{U}_c^{in} 。

- 1: 根据训练集 \mathcal{U}_c^{train} 训练逻辑回归模型
 - 2: **for** \mathbf{x} in \mathcal{U}_c^{in} **do**
 - 3: $\text{IF}(\mathbf{x}) = 0$
 - 4: **for** \mathbf{z}_{test} in \mathcal{U}_c^{test} **do**
 - 5: 利用公式3-6计算 $I(\mathbf{x}, \mathbf{z}_{test})$
 - 6: $\text{IF}(\mathbf{x}) = \text{IF}(\mathbf{x}) + I(\mathbf{x}, \mathbf{z}_{test})$
 - 7: **end for**
 - 8: **end for**
 - 9: 对集合 \mathcal{U}_c^{in} 中的样本根据 $\text{IF}(\mathbf{x})$ 进行排序
 - 10: 返回 \mathcal{U}_c^{in} 中 $\text{IF}(\mathbf{x}) > 0$ 的样本。
-

3.3.4 复杂度分析

进一步, 我们给出了 Record 框架的时间复杂度和空间复杂度分析。具体而言, 在每个时刻 t , Record 首先需要将每个类别 c 中的样本划分为集合 \mathcal{U}_c^{in} 和集合 \mathcal{U}_c^{out} , 由于划分过程需要将样本按照预测概率进行排序, 因此这一步的时间复杂度为排序复杂度, 即, $O(m^t \log(m^t))$ 。然后, Record 需要计算每个样本在集合 \mathcal{U}_c^{out} 上的影响力, 根据 [78] 可知, 计算每个样本影响函数的时间复杂度为 $O\left(\frac{(m^t)^2}{4}d\right)$, 其中 d 是样本特征空间的维度。综上, Record 方法在每个时刻 t 的时间复杂度为 $O\left(\frac{(m^t)^2}{4}d + m^t \log(m^t)\right)$ 。对于空间复杂度, 每一步 Record 只能利用存储的 B 个样本以及当前到达的无标注数据集 \mathcal{U}^t , 因此 Record 的空间复杂度为 $O(B + m^t)$ 。

3.4 实验验证

在本节中, 我们针对多种流式数据基准数据集和主流方法进行了大量的对比实验, 以评估 Record 方法在资源受限流式半监督学习问题中的有效性。

3.4.1 实验设定

我们一共在八种数据集上进行了实验以评估算法的性能表现, 其中包括四种常用的分类数据集: *Optdigits*、*Satimage*、*Twonorm*、*Spam* 和四种数据分布持续变化的流数据学习基准数据集: *1CHT*、*2CDT*、*UG_2C_2D*、*UG_2C_5D* [126]。

数据集所包含的样本数量最小为 5,620，最大为 200,000，特征空间维度从 2 到 500 不等。这些数据集的统计信息总结在表 3-1 中，其中包括样本数量、特征空间维度、类别空间维度，以及每个类别所选取的标注样本的数量。

表 3-1: 数据集统计信息

数据集	样本数量	特征空间维度	类别空间维度	每类标注数量
Optdigits	5,620	64	10	10
Satimage	6,435	36	7	10
Twonorm	7,400	20	2	10
Spam	9,324	500	2	50
1CHT	16,000	2	2	1
2CDT	16,000	2	2	5
UG_2C_2D	100,000	2	2	1
UG_2C_5D	200,000	5	2	5

为了验证 Record 框架与各种半监督学习算法结合的通用性，我们考虑了如下三种主流的半监督学习方法：

- Mean Teacher [128]: Mean Teacher 算法是代表性的深度半监督学习算法。具体来说，Mean Teacher 采用一致性正则损失作为无监督损失函数，通过 EMA 操作对历史模型进行集成，并最小化集成模型与当前模型在无标注样本上输出的预测概率之间的均方误差。目前，基于一致性正则的半监督学习方法在多种半监督学习任务中取得了最先进的结果 [109]。
- Label Propagation [164]: 标注传播算法是一种经典的基于图的半监督学习算法，可以将标注信息沿着图中的边由标注数据传播至无标注数据。标注传播算法的内在假设是平滑假设，即，相似的样本应该具有相似的标注。
- S³VM [70]: 半监督支持向量机 (S³VM) 也是经典的半监督算法之一，其采用低密度分割假设，认为决策边界应该穿越数据密度较低的区域。

数据集预处理。对于四个常用的分类数据集，我们通过对样本进行重新分组来手动模拟分布变化。对于 *Optdigits*、*Twonorm* 和 *Satimage* 数据集，每个时间步有 200 个样本到达，其中 160 个作为无标注训练样本，40 个作为测试样本；对于 *Spam* 数据集，每个时间步有 400 个样本到达，其中 280 个作为无标注训练样本，120 个作为测试样本。对于四个流数据基准数据集，我们采用 [126] 中提供的默

认预处理方式。在每个时刻，*1CHT* 和 *2CDT* 中有 400 个样本到达，*UG_2C_2D* 中有 1000 个样本到达，*UG_2C_5D* 中有 2000 个样本到达，其中 70% 作为无标注训练样本，30% 作为测试样本。对于 S^3VM 算法，因为其处理多分类任务的效率低下，我们基于 *Optdigits* 和 *Satimage* 数据集构建了一个类别 1 对类别 7 的二分类任务用于实验。

参数详情。对于 S^3VM 和 Label Propagation 算法，我们采用开源机器学习库 `sklearn`^① 中的实现，其中 S^3VM 算法的核函数设置为 RBF 核，Spam 数据集上的标签传播算法采用 9NN 算法，其它所有超参数均设置为默认值。对于 Mean Teacher 算法，我们采用官方开源代码实现^②，优化轮数设置为 1,000 次。实验中存储资源预算 *B* 设置为 100，即最多只有 100 个无标注样本可以被存储。

3.4.2 与基线方法对比结果

由于现有半监督研究中没有与本文研究的问题设置相同的工作，我们首先与多种基线学习方法进行比较，以验证 Record 方法在数据分布变化的流数据上的性能表现，具体采用的对比方法如下：

- **Supervised:** 监督学习方法，忽略所有的无标注数据，只利用初始阶段给定的标注数据集训练监督学习模型。监督学习方法可以看做是性能的下界。
- **Proba:** 概率选择法，利用模型对样本输出的预测概率进行子集选择，即，具有更高预测概率的样本优先被选择。此类方法在多种任务中被证明是有效的，例如分布外样本检测 [64] 等。
- **Random:** 随机选择法，在每个时刻，根据存储资源的约束随机选取部分样本进行存储。

此外，我们还比较了利用所有数据标注进行训练的 Oracle 方法：

- **Oracle:** 假设所有样本的真实标注都可获取，在每个时刻利用当前到来的所有标注样本训练监督学习模型。由于标注获取的困难性，Oracle 在实践中是无法实现的，可以视为性能的上界。

图 3-6、图 3-7 和图 3-8 分别展示了采用 Mean Teacher、Label Propagation 和 S^3VM 作为半监督学习算法时，Record 方法在 8 个数据集上运行 10 次实验得到

^①<https://scikit-learn.org/>

^②<https://github.com/CuriousAI/mean-teacher>

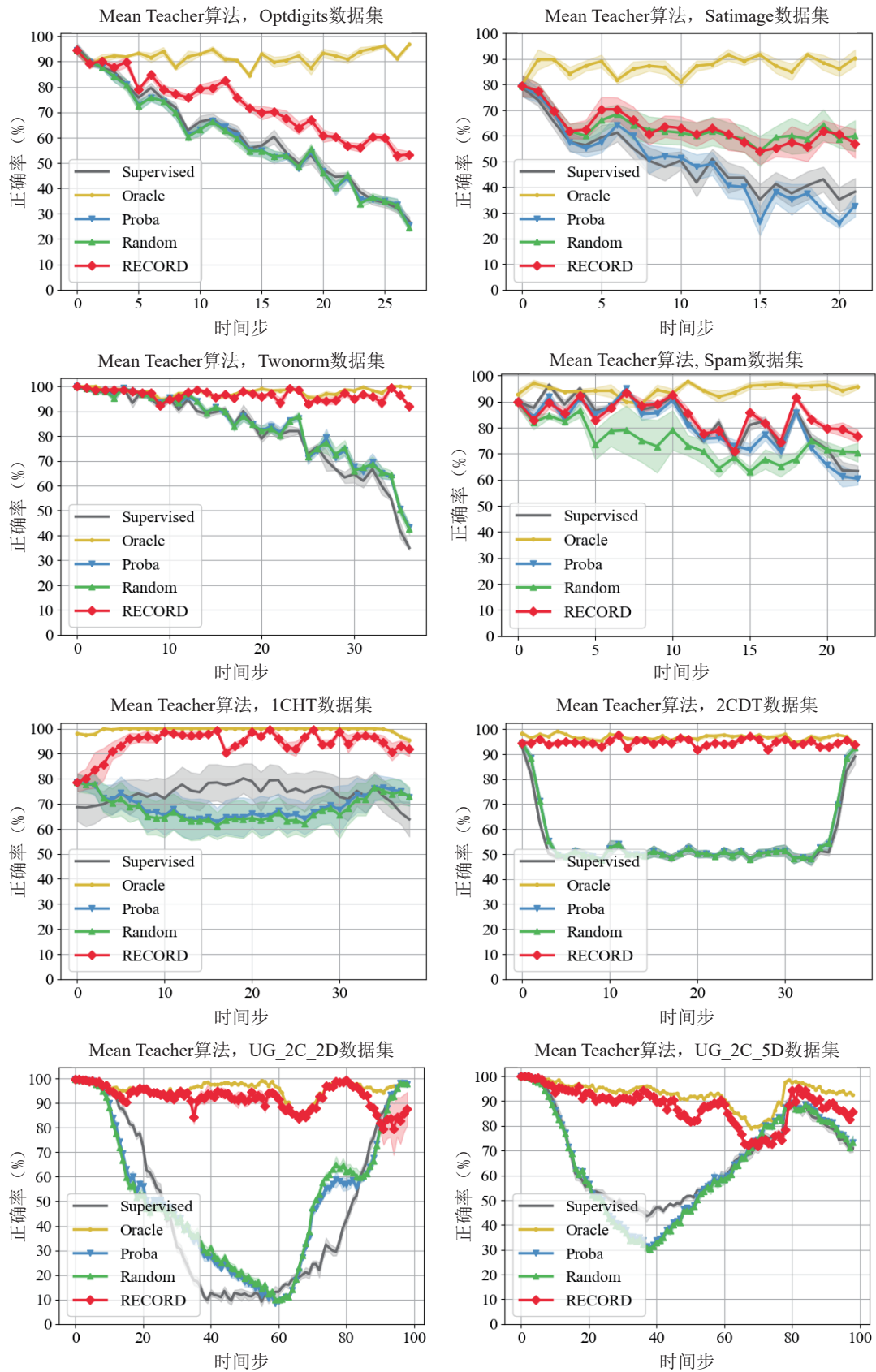


图 3-6: 以 *Mean Teacher* 作为半监督学习算法时, Record 框架在 8 种数据集上的分类正确率。曲线阴影部分表示重复 10 次实验性能的标准差。

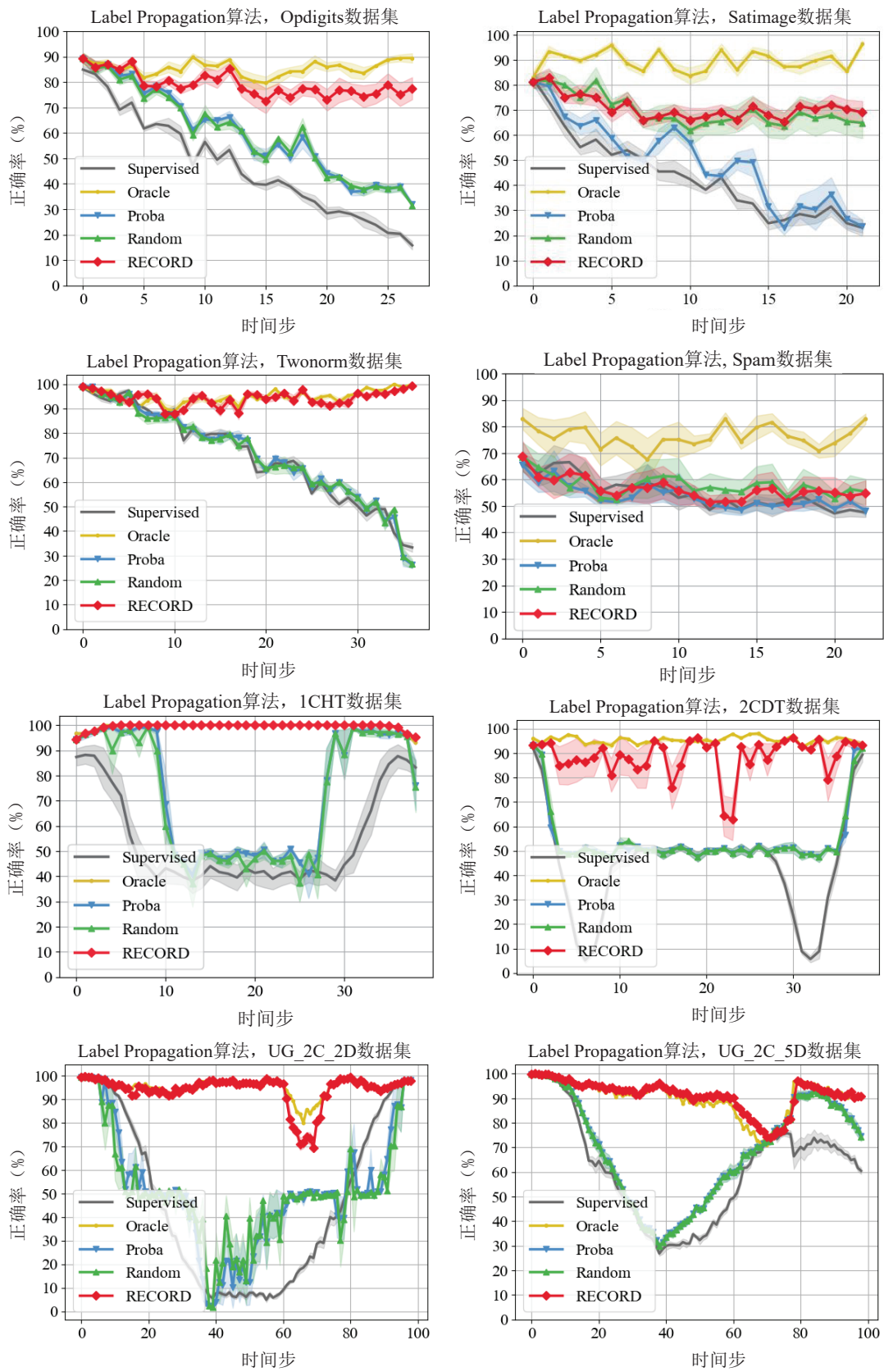


图 3-7: 以 *Label Propagation* 作为半监督学习算法时, Record 框架在 8 种数据集上的分类正确率。曲线阴影部分表示重复 10 次实验性能的标准差。

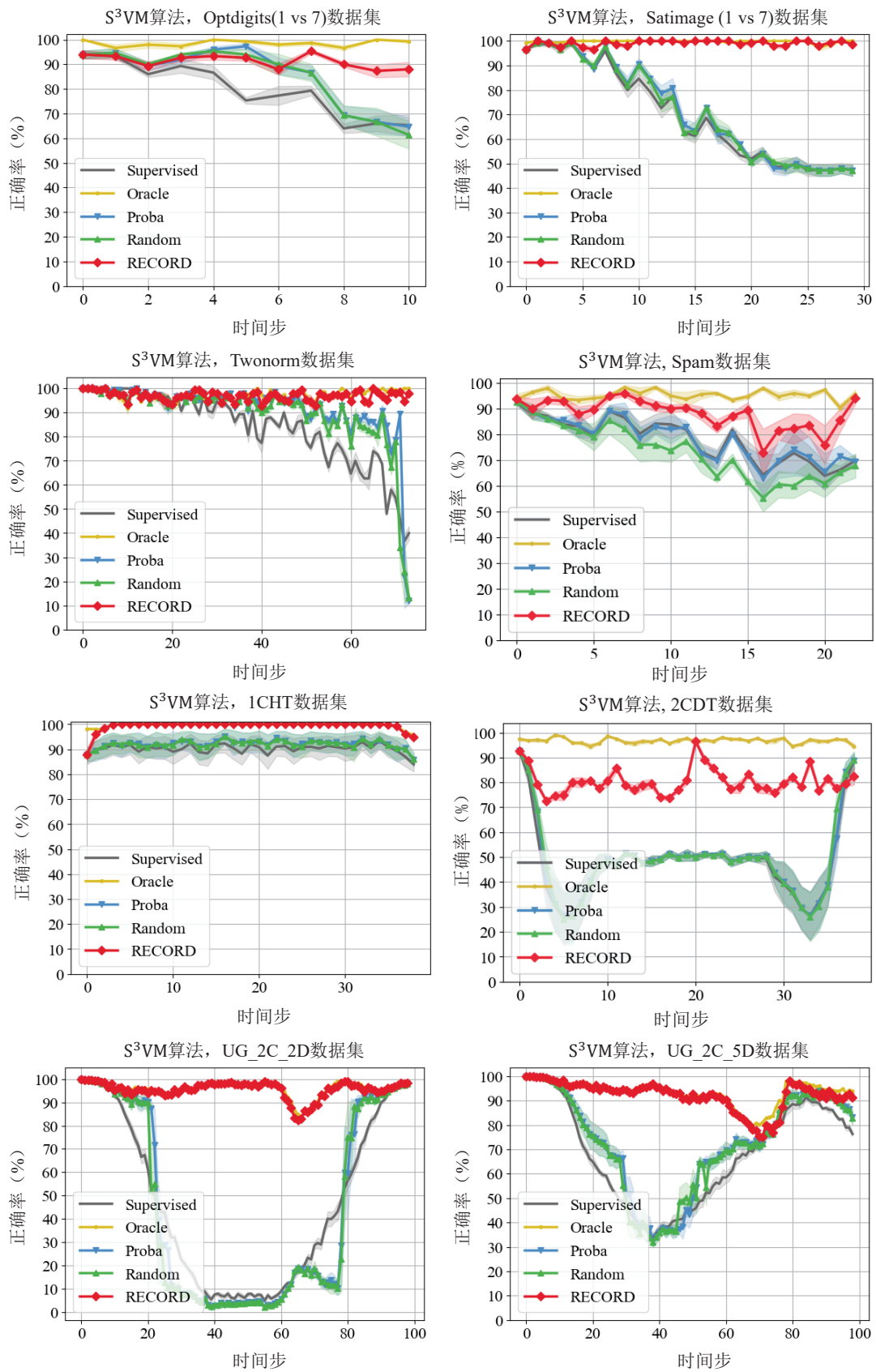


图 3-8: 以 S^3VM 作为半监督学习算法时, Record 框架在 8 种数据集上的分类正确率。曲线阴影部分表示重复 10 次实验性能的标准差。

的分类正确率的均值与标准差。从实验结果我们可以观察到, 无论使用何种半监督学习算法, 在任意数据集上, Record 都在所有的比较方法中获得了最优的分类正确率。与 Supervised 方法相比, 简单的 Proba 和 Random 方法都因为分布变化的现象出现了性能退化的现象, 而 Record 始终保持性能的提升。此外, Record 方法非常接近 Oracle 方法的性能, 在某些情况下甚至优于 Oracle, 主要原因是相比 Oracle 利用当前时刻到来的所有数据及标注训练模型, Record 存储了以往时刻的样本, 这种现象验证了 Record 可以保留对后续时刻的数据分布有帮助的样本。以上结果证明了 Record 方法的有效性和通用性。

3.4.3 与最先进方法对比结果

进一步, 我们将本文提出的 Record 方法, 与目前在该问题上最先进的两种方法 TLP [130] 和 COMPOSE [39] 进行对比。

TLP 算法针对流式数据的图半监督学习问题, 提出通过维护数据流的一个有效子集, 实现新样本到来时模型快速更新。TLP 的有效性已在各种实际任务中得到证明, 例如心电图分析、自动驾驶等 [130]。

本文提出的 Record 方法与 TLP 方法的对比结果如图 3-9 所示, 其中 Record 方法采用 Label Propagation 作为半监督学习算法, TLP 算法的超参数均设置为论文推荐的默认值。从实验结果我们可以发现, TLP 方法会随着数据分布的变化而出现性能退化, 而本文提出的 Record 方法在分布变化时表现更稳定, 不仅没有出现性能退化的问题, 而且与 TLP 方法相比取得了显著的性能提升。

此外, 我们还汇报了 Record 方法和 TLP 方法在整个数据流上各个时刻的平均性能, 如表 3-2 所示。我们可以看到 TLP 方法在 *Twonorm*、*Spam* 和 *UG_2C_5D* 数据集上的性能表现甚至比基线的监督学习方法更差, 而本文提出的 Record 方法则实现了更优的性能, 在某些情况下甚至可以接近性能上界的 Oracle 方法。在 8 个数据集上的平均分类正确率 Record 相比 TLP 提升近 25%。以上结果证明了相比于目前最先进的 TLP 方法, Record 能更好的处理资源受限的流式半监督学习问题, 在分布变化的数据流中具有更加稳健的性能表现。

COMPOSE 算法采用基于样本集几何原型的子集选择方法选择样本进行存储, 但是这种基于几何的方法只能应用于特征低维的数据集, 所以我们仅在 *1CHT* 和 *2CDT* 数据集上与 COMPOSE 方法进行对比, 其中 COMPOSE 方法两

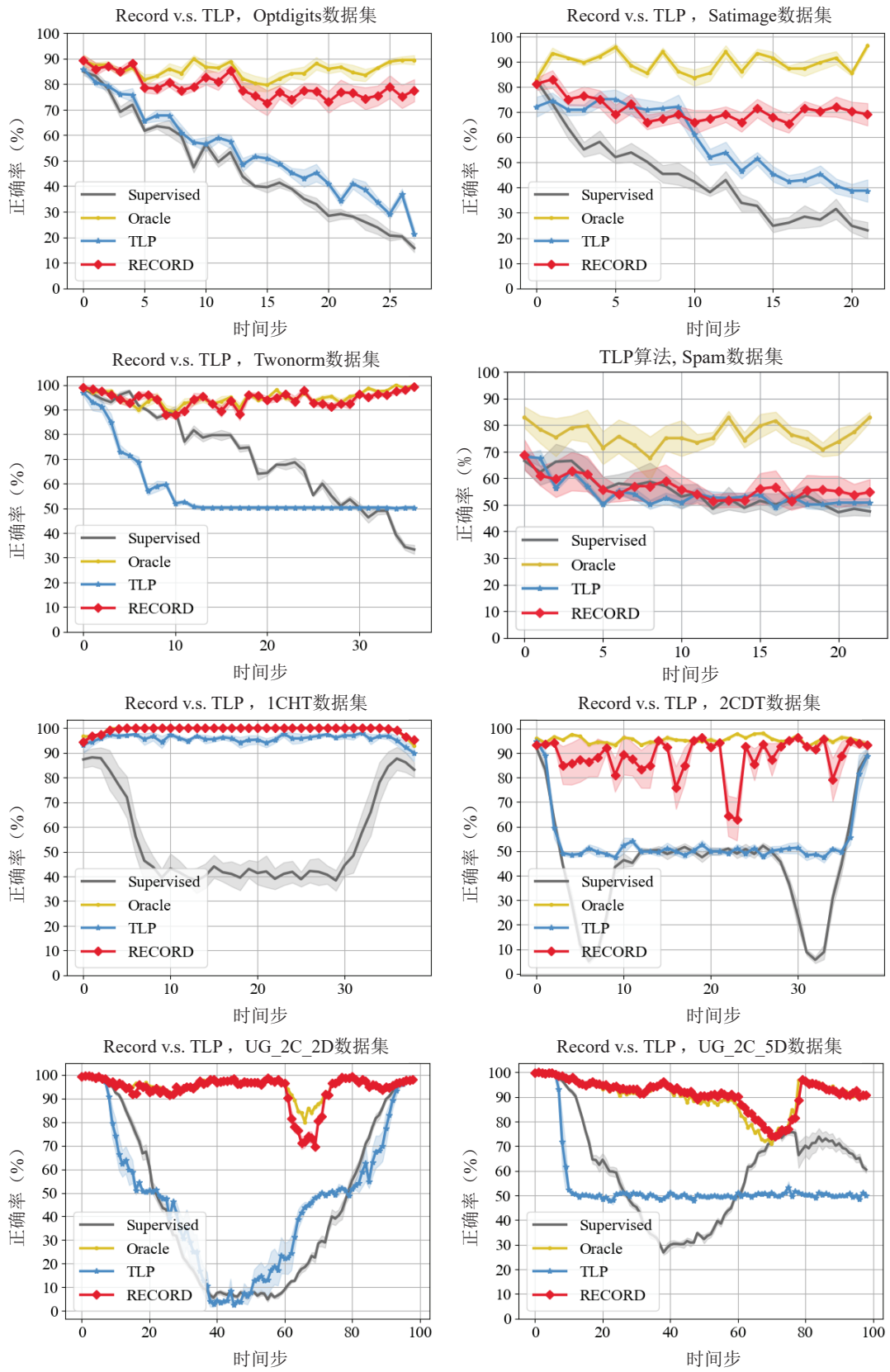


图 3-9: Record 方法和 TLP 方法在 8 种数据集上的分类正确率。曲线阴影部分表示重复 10 次实验性能的标准差。

表 3-2: Record 方法和 TLP 方法在所有时刻分类正确率的均值和标准差。其中, 粗体部分表示最优方法, 下划线部分表示比基线的监督学习方法更差。

数据集	Supervised	TLP	Record	Oracle
Optdigits	46.64 ± 1.04	53.56 ± 1.42	79.20 ± 2.01	85.61 ± 0.25
Twonorm	70.84 ± 0.36	<u>57.24 ± 0.86</u>	94.31 ± 0.38	94.99 ± 0.13
Satimage	43.47 ± 2.61	58.43 ± 2.54	71.02 ± 2.93	89.59 ± 0.70
Spam	55.12 ± 3.48	<u>54.18 ± 1.39</u>	56.55 ± 4.72	76.35 ± 2.41
1CHT	54.79 ± 4.93	95.78 ± 1.27	99.44 ± 0.05	99.44 ± 0.01
2CDT	44.47 ± 0.46	54.28 ± 0.48	88.68 ± 1.57	95.36 ± 0.13
UG_2C_2D	46.40 ± 0.23	47.56 ± 1.12	94.11 ± 0.12	95.35 ± 0.04
UG_2C_5D	61.00 ± 1.01	<u>54.40 ± 0.17</u>	91.39 ± 0.08	90.41 ± 0.05
平均性能	52.84 ± 1.77	59.43 ± 1.16	84.34 ± 1.48	90.89 ± 0.47

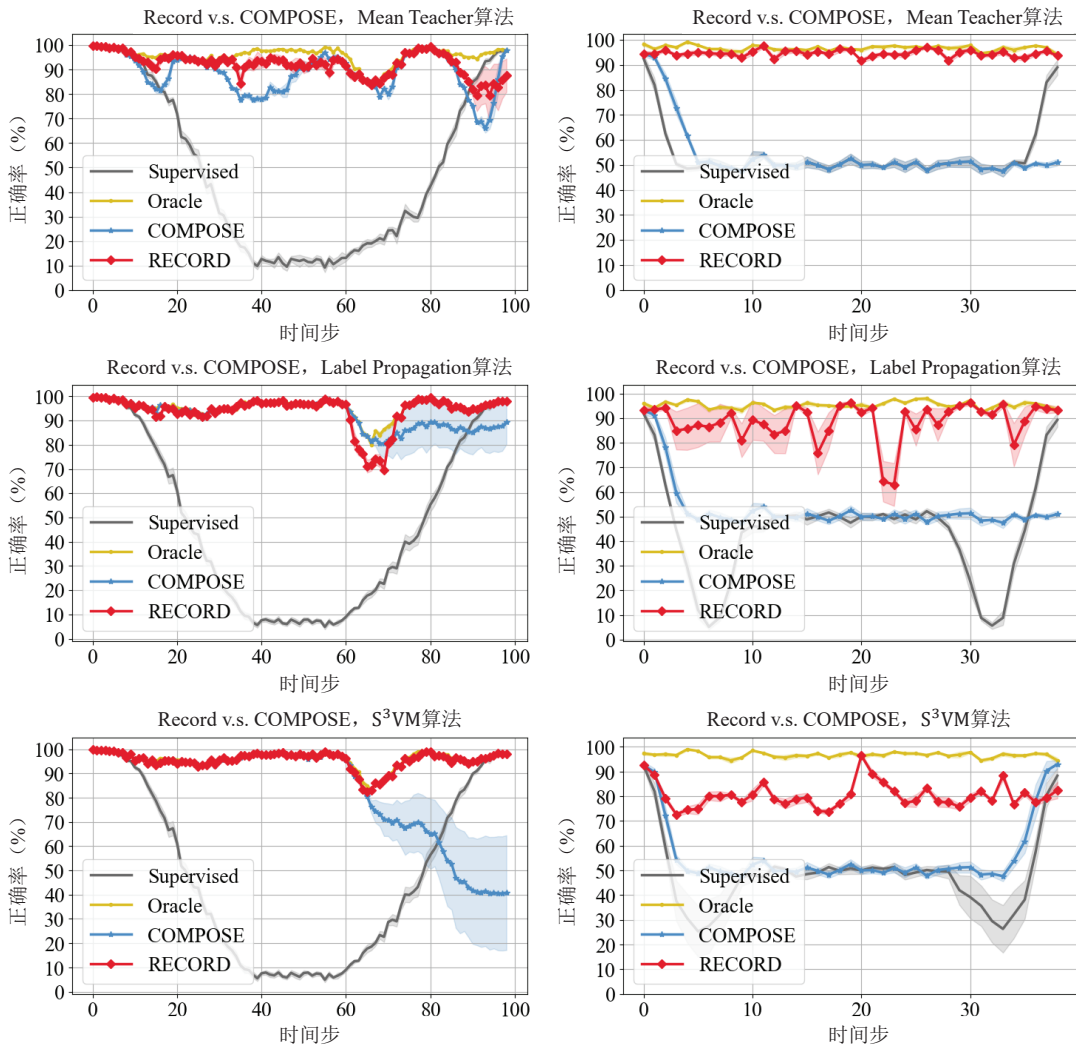


图 3-10: 分别采用 Mean Teacher、Label Propagation 和 S³VM 作为半监督学习算法时, Record 和 COMPOSE 方法在数据集 UG_2C_2D (左列) 和数据集 2CDT (右列) 中的分类正确率。曲线阴影部分表示重复 10 次实验性能的标准差。

表 3-3: 采用 3 种不同的半监督学习算法时, Record 和 COMPSE 在数据集 2CDT 和 UG_2C_2D 上的正确率均值和标准差。其中, 粗体部分表示最优方法, 下划线部分表示比基线的监督学习方法更差。

2CDT				
半监督算法	Supervised	COMPOSE	Record	Oracle
Mean Teacher	54.40 ± 0.36	<u>54.09 ± 0.24</u>	94.60 ± 0.13	96.72 ± 0.18
Label Propagation	44.47 ± 0.46	53.19 ± 0.20	88.70 ± 1.57	95.36 ± 0.13
S ³ VM	48.40 ± 2.29	56.05 ± 0.70	80.41 ± 0.19	96.71 ± 0.16
平均性能	49.04 ± 1.04	54.44 ± 0.38	87.90 ± 0.63	96.26 ± 0.47
UG_2C_2D				
Mean Teacher	47.82 ± 0.65	88.76 ± 0.35	92.34 ± 0.66	95.94 ± 0.02
Label Propagation	46.40 ± 0.23	92.63 ± 2.89	94.11 ± 0.12	95.35 ± 0.04
S ³ VM	46.49 ± 0.19	83.05 ± 5.09	95.55 ± 0.05	96.02 ± 0.04
平均性能	46.90 ± 0.36	88.15 ± 2.78	94.0 ± 0.28	95.77 ± 0.03

个重要的超参数 α 和 cp 分别设置为 0.4 和 0.7, 与论文推荐设置相同。图 3-10 展示了分别以 Mean Teacher、Label Propagation 和 S³VM 算法作为半监督学习算法时, Record 和 COMPOSE 方法的分类正确率。从实验结果我们可以发现, 在 UG_2C_2D 数据集中, Record 和 COMPOSE 方法在前 60 个时间步都可以实现优异的性能, 然而, 在 60 步之后, COMPOSE 方法面临严重的性能下降问题, 而 Record 方法始终保持接近 Oracle 方法的性能。在 2CDT 数据集上, COMPOSE 在前 5 个时间步就出现了性能下降的问题, 而 Record 方法相比 COMPOSE 方法始终有显著的性能提升。

我们也汇报了在整个数据流所有时刻上的平均性能, 如表 3-3 所示。从实验结果我们可以观察到类似的现象, 即, COMPOSE 在某些情况下性能弱于基线的监督学习方法, 出现了性能退化的现象, 而本文提出的 Record 方法始终优于 COMPOSE 方法, 并接近性能上界 Oracle 方法。在 2CDT 数据集中, 三种实验设置下的平均分类正确率 Record 相比 COMPOSE 提升 30% 以上, 在 UG_2C_2D 数据集中, 平均分类正确率 Record 相比 COMPOSE 提升 5% 以上。上述结果验证了在资源受限且数据分布变化的流式半监督学习场景中, 本文提出的 Record 方法能够实现最先进的性能。

3.4.4 Record 方法对存储资源的稳健性

在数据分布逐渐变化的流式半监督学习场景中，可以存储的样本越多，模型就能更好的适应数据分布的变化，因此，研究存储资源的限制对 Record 方法的影响是一个有意义的问题。为此，我们在四个基准数据集上研究了当可存储样本的数量从 80 下降到 10 时，Record 方法性能的变化情况，结果如图 3-11 所示，其中 S^3VM 算法作为 Record 框架中的半监督学习算法。从实验结果我们可以发现，只要可以存储的样本数量超过 20 个，Record 就不会出现性能显著下降的问题，这验证了 Record 方法对存储资源变化的稳健性。

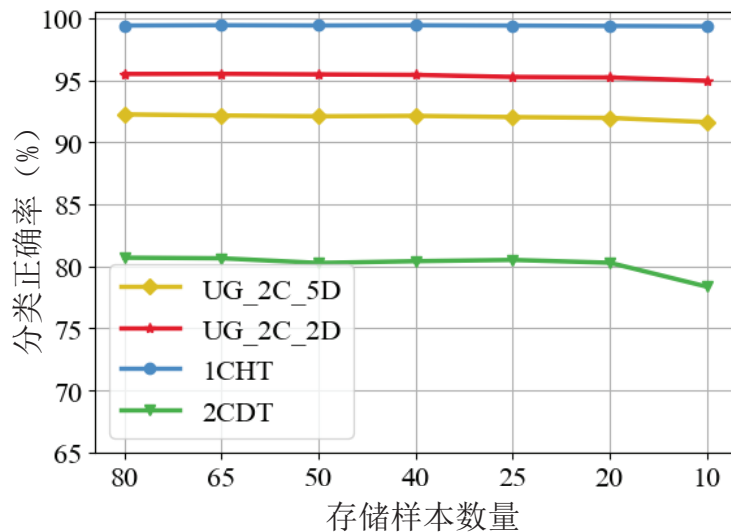


图 3-11: 随存储资源变化 Record 性能的变化情况。

3.4.5 可视化分析

为了进一步说明 Record 方法的有效性，我们对 Record 所选择的样本进行可视化分析。在 UG_2C_2D 数据集中，以 S^3VM 作为 Record 框架中的半监督学习算法时，Record 从数据流的时刻 5 到时刻 14 所选择的样本如图 3-12 所示，其中蓝色和红色圆形分别表示两个类别的样本，浅色圆形表示在当前时刻到达的所有无标注样本，深色圆形表示 Record 方法保存的样本。从实验结果我们可以看到 Record 算法可以保留处于数据分布 $p^t(\mathbf{x})$ 和 $p^{t+1}(\mathbf{x})$ 之间重叠区域的样本，这有助于当前分布的学习并跟踪未来分布偏移的趋势，与我们的目标一致。可视化结果验证了 Record 方法采用的基于影响力机制的样本选择技术可以在数据分

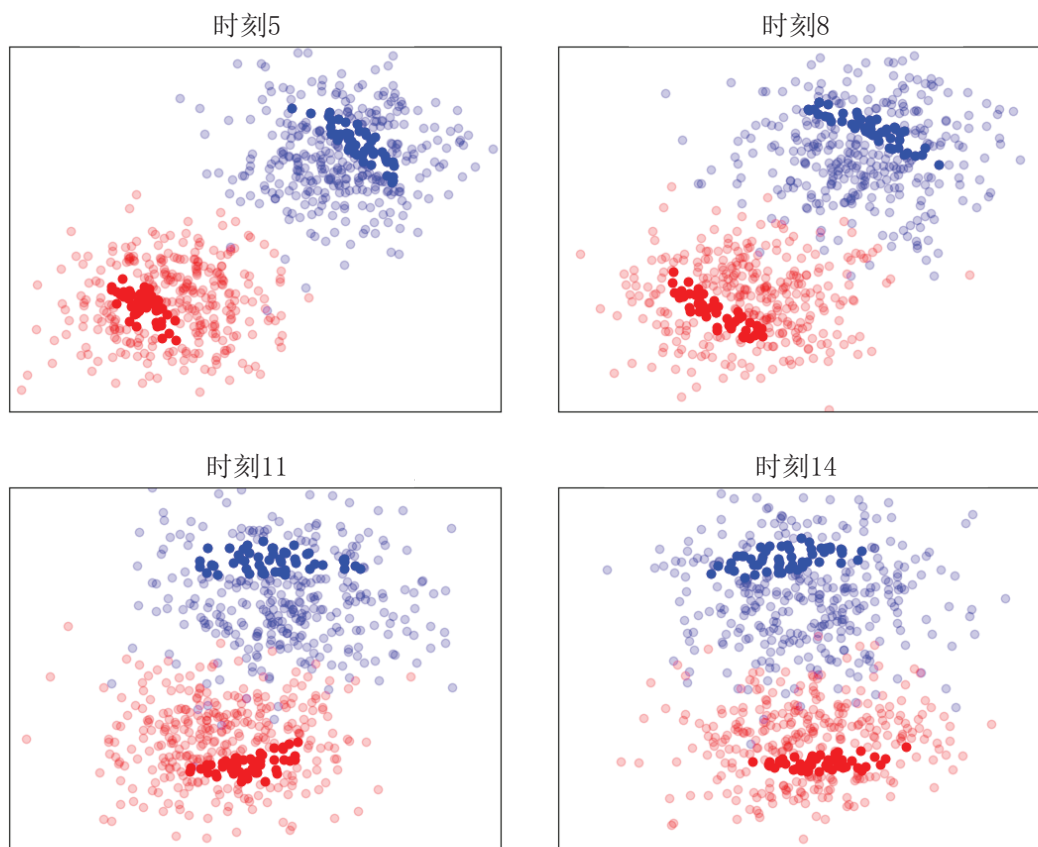


图 3-12: 可视化展示从时刻 5 到时刻 14, *UG_2C_2D* 数据集的分布变化情况及 Record 所选择的样本。其中蓝色和红色圆形分别表示两个类别的样本, 加深的颜色表示 Record 在当前时刻存储的样本。

布变化的条件下有效进行样本选择。

3.5 小结

在本章中, 我们考虑开放环境下数据流式到来, 数据分布随环境不断动态变化且存储资源受限的场景, 提出了一种新的半监督学习问题设置: 资源受限的流式半监督学习。在该场景中, 标注数据只在初始阶段给出, 无标注数据以流的形式不断收集得到的, 且其数据分布会逐渐发生变化, 此外, 由于存储资源的约束, 无法将所有时刻的无标注样本全部存储下来用于模型训练。这种新颖的问题设置目前还没有被很好的研究, 我们针对该问题提出了一种系统性的解决方案 Record。Record 的基本思想是根据内存资源的约束选取样本子集进行存储, 具体而言, Record 采用一种先进的基于影响力机制的样本选择策略, 计算旧数据分布下的样本在新数据分布上的影响力, 根据存储资源的约束选择与新

数据分布最具关联的样本子集，使得模型能够快速有效地适应变化的数据分布。在大量数据集上的实验验证了现有半监督学习方法在资源受限的流式半监督学习场景下性能会随分布的变化出现严重退化的情况，而 Record 方法在整个数据流中性能稳健，显著优于现有的方法并且能够接近性能上界的 Oracle 方法。此外，Record 是通用性的学习框架，可以与任意半监督学习算法相结合。

目前我们只考虑了流数据中分布会产生变化的问题，后续我们将进一步考虑数据的类别、属性都有可能发生变化的场景。此类问题在现实任务中十分常见，例如，对于自动驾驶任务，可能存在某些传感器失灵导致某个时间段部分属性缺失；对于在线学习任务，可能会不断产生新的未见过的类别。如何保证在更复杂的流数据环境中半监督学习的稳健性是未来一个重要的研究方向。此外，在数据不断产生的过程中，虽然我们无法在每个时间点都获取到所有数据的准确标注，但是仍然有可能与环境进行交互获取一些带噪的反馈，如何在流数据学习中利用带噪反馈提升性能，降低对高质量监督信息的依赖，也是一个值得研究的方向。

本章的主要工作已经成文发表，包括：

- [Lan-Zhe Guo, Zhi Zhou, Yu-Feng Li](#). RECORD: Resource Constrained Semi-Supervised Learning under Distribution Shift. In: **Proceedings of the 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD'20)**, San Diego, CA, USA, pp.1636-1644, 2020. (中国计算机学会 A 类会议，第一作者)

第四章 适于先验知识不足的 稳健半监督学习

4.1 引言

半监督学习的研究目标是当标注数据不足时通过引入大量易于获取的无标注数据提升模型性能，而要有效利用无标注数据，必然要做一定的假设将无标注样本中所包含的数据分布信息与数据标注信息进行关联。常见的半监督假设如：聚类假设 (Cluster Assumption)，假设数据分布在不同的簇结构中，位于同一个簇中的样本具有相同的数据标注；流形假设 (Manifold Assumption)，假设数据分布在一个流形中，位置相近的样本应该具有相似的类别标注；低密度假设 (Low-Density Assumption)，假设分类器的决策边界应该穿过数据密度比较低的区域等。采取不同的假设将产生不同的半监督学习模型，如半监督支持向量机、图半监督学习等。

一个自然的问题是，当给定少量标注数据和大量无标注数据，该如何选择半监督学习模型，使其能够在当前任务上表现良好。不同于监督学习可以利用大量标注数据通过交叉验证等方式进行选择，半监督学习中标注数据通常不足以提供可靠的模型选择，因此传统封闭环境下的半监督学习依赖充分的先验知识选取最适合当前任务的学习模型，如数据的真实分布等知识。然而，在现实世界的开放环境中，往往无法获得充分的先验知识，而一旦模型选择错误，不符合真实数据中数据分布信息与数据标注信息之间的关联，就可能会导致半监督学习出现不安全的问题，即，性能甚至弱于只利用有标注数据的简单监督学习。例如，对于生成式半监督学习，Cozman 等人 [29] 指出，当生成方法中采取的数据分布假设不正确时，生成式半监督学习会出现性能退化的不安全现象；对于半监督支持向量机，Li 和 Zhou [92] 指出，基于给定的少量标注数据和大量无标注数据训练，可以获得多个满足低密度假设的决策边界，而错误的模型选择会导致半监督支持向量机性能退化等。图 4-1 展示了这一现象。

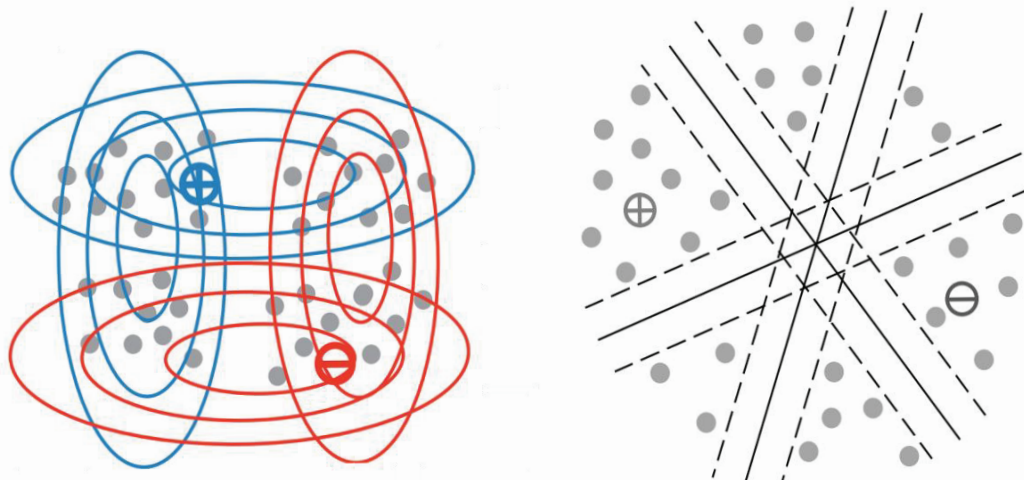


图 4-1: 给定训练数据, 可以获得多个半监督学习模型, 由于先验知识不充分, 无法进行可靠的模型选择。

半监督学习不安全的现象无疑违背了半监督学习利用无标注数据提升模型性能的出发点, 阻碍了半监督学习在更多现实开放环境任务中的应用。本章针对该问题展开研究, 提出了一种基于模型集成的安全半监督学习框架 SafeW。给定多个半监督学习模型, 由于先验知识不充分, 无法进行可靠的模型选择, SafeW 通过对多个模型的预测结果进行集成以获取最终的预测结果。具体而言, 我们提出了一种基于最大最小优化的集成学习框架, 优化模型在最坏情况下相比于基线监督学习模型的性能增益, 保证模型在最坏情况下依然能取得安全的性能提升。SafeW 框架具有如下优势: 1) 该框架对于多种半监督学习常用的损失函数均可以在理论上实现安全的性能, 如用于回归任务的均方损失 (Mean Square Loss), 用于分类任务的铰链损失 (Hinge Loss) 等; 2) 该框架对若干个不确定的模型进行加权集成, 同时可以灵活嵌入与模型性能相关的先验知识; 3) 该框架可以转化为简单的二次优化问题或线性优化问题, 从而高效地获得全局最优解; 4) 该框架是通用的学习框架, 除半监督学习之外, 容易扩展应用到其它弱监督学习场景, 如标注噪声学习、领域自适应学习、多示例学习等。

4.2 相关工作

本章的工作主要与半监督学习的安全性相关, 同时也可以扩展以解决其它弱监督学习问题, 如领域自适应学习、多示例学习、标注噪声学习。

4.2.1 安全半监督学习

能够利用少量标注数据和大量无标注数据构建机器学习模型的半监督学习方法在机器学习领域重要的研究方向,已经有大量的半监督学习方法被提出,包括生成式半监督学习 [102],图半监督学习 [164],半监督支持向量机 [70],基于分歧的半监督学习 [10] 等。这些半监督模型均依赖一定的数据假设,如流形假设、聚类假设、大间隔假设等,当先验知识不充分时,无法判断数据符合何种假设,从而无法选取最可靠的半监督模型,而错误的模型选择将导致模型性能退化,甚至不如简单的监督学习性能。因此,安全半监督学习得到了一定的关注,例如, Li 和 Zhou [92] 指出在半监督支持向量机中,可能产生多个大间隔的分类器,错误的模型选择将导致模型性能退化,因此,他们提出在给定多个候选大间隔分类器的情况下,可以通过优化最坏情况下的性能增益来构建安全的半监督支持向量机,并且证明当真实标注指派可以由某个决策边界实现时,该方法可以实现安全性。Balsubramani 等人 [4] 指出,在给定多个候选半监督分类器时,如果数据的真实标注在某个特定的候选集内,可以通过模型集成的方式实现安全的半监督学习。然而,目前关于安全半监督学习的研究均局限于特定的任务,如分类任务,或者特定的算法,如半监督支持向量机,对于通用的半监督学习算法和任务,如何构建安全的半监督学习模型仍然是一个开放的问题。

4.2.2 领域自适应学习

领域自适应学习和半监督学习同属于弱监督学习的一种,处理的是监督信息不完整的弱监督数据。领域自适应学习和半监督学习均假设当前任务中只能获取少量的标注数据,但区别在于领域自适应学习假设可以从与当前域数据分布不同的源域中获取大量的标注样本,如何利用分布不同的源域样本帮助提升目标域上的模型泛化性能是领域自适应学习的目标。目前已经有大量领域自适应学习的方法被提出,例如,基于样本迁移的方法 [31],基于特征表示迁移的方法 [114],基于参数迁移的方法 [11],基于关联知识迁移的方法 [101] 等。在领域自适应学习中,也存在不安全的现象,即利用了更多源域信息反而导致目标域泛化性能下降的负迁移 (Negative Transfer) 的现象 [111]。尽管该问题被认为是领域自适应学习中一个关键的问题,但目前的解决方案还比较少。Rosenstein

等人 [117] 通过实验分析, 如果源域和目标域数据分布差异较大时, 迁移可能会损害目标任务的性能, 不如只利用少量的目标域数据进行模型训练。Bakker 和 Heskes [3] 提出一种可以用于多个任务先验分布估计的贝叶斯方法, 用于判断不同任务之间是否适合迁移。Argyriou 等人 [2] 利用数据表示将任务进行分组, 认为同一组内的任务是更容易进行迁移的。Ge 等人 [50] 提出可以为来自不同源域的数据赋予不同的权重, 通过对数据进行加权减弱负迁移的影响。

4.2.3 标注噪声学习

标注噪声学习也是机器学习领域的一个热门研究方向, 其处理的是监督信息不正确的弱监督数据, 即, 在训练数据中, 存在部分数据其给定的标注不是真实对应的标注。目前, 已经有大量标注噪声学习的方法被提出, 如基于数据清洗的方法 [108], 基于样本赋权的方法 [116], 基于损失校正的方法 [64] 等。但是, 有研究表明, 标注噪声学习算法并不安全, 有些情况下, 利用了更多的不可靠数据, 反而不如只利用少量正确的标注数据训练得到的模型性能 [44, 47, 65, 42]。在理论方面, Manwani 和 Satry [99] 研究了在经验风险最小化框架下损失函数的稳健性, 指出 0-1 损失函数对标注噪声更具稳健性, 而其它损失函数理论上并不具备该性质。在算法方面, 集成学习方法, 如 Bagging 和 Boosting 对标注噪声稳健性更强 [44], 并且 Bagging 通常具有比 Boosting 更好的性能表现 [35]。

4.2.4 多示例学习

多示例学习是弱监督学习的一种, 其处理的是监督信息不具体的数据。例如, 对于二分类问题, 给定若干个样本作为一个包 (Bag), 我们只知道包内是否包含正类样本, 但并不知道具体哪一个样本是正类。多示例学习领域也已经有大量算法被提出, 如基于密度的方法 [154], 基于 k 近邻的方法 [131], 基于支持向量机的方法 [1], 基于集成学习的方法 [139, 161, 18] 等。然而, 多示例学习也存在不安全的现象, 在利用了更多标注不具体的样本后, 性能可能不如只利用少量标注样本的监督学习模型。例如, Ray 和 Craven [115] 比较了多示例学习方法和对应的监督学习模型的性能, 他们发现, 在很多情况下, 简单监督学习模型可以实现更优的性能。Carbonneau 等人 [18] 研究了多种多示例学习方法识别

正类样本的能力，他们发现，多示例学习方法严重依赖数据的性质，当数据与算法相匹配时，多示例学习算法表现良好，反之，则出现性能退化的不安全现象。

综上所述，半监督学习、领域自适应学习、标注噪声学习、多示例学习等弱监督学习场景均存在不安全的现象，即，利用了更多的弱监督数据后反而不如只利用少量标注数据的监督学习性能。目前，还没有方法能够解决通用弱监督学习场景下安全性的问题。

4.3 本文工作

在本节，我们介绍本文提出的安全半监督学习框架，首先，我们介绍相应的背景知识，接着，介绍本文提出的框架 SafeW，并给出相应的理论分析，然后，我们展示了 SafeW 框架在分类任务和回归任务上的优化方法，最后，我们通过大量场景上的实验验证该学习框架的通用性和安全性。

4.3.1 背景介绍

在半监督学习中，由于缺少大量精确的标注信息，通过集成的方式综合考虑多个不确定的模型是提升半监督学习稳健性的一种常用方案 [156]。传统的集成学习方法只考虑如何利用多个基学习器构建性能提升最大化的模型，而忽略了最坏情况下的安全性，但现实任务要求我们不仅要实现性能的提升，更要保证在最坏情况下模型性能不会下降，防止模型出现重大的错误。

具体而言，给定 b 个模型在无标注数据上的预测结果 $\{\mathbf{f}_1, \dots, \mathbf{f}_b\}$ ，其中 $\mathbf{f}_i \in \mathbb{H}^u$ ， $i = 1, \dots, b$ ， u 表示无标注样本的个数。有大量策略可以生成多个不确定的半监督学习模型，例如，通过不同的学习算法、不同的数据采样、不同的模型参数等 [156]。令 $\mathbf{f}_0 \in \mathbb{H}^u$ 表示基线模型的预测结果，即，只利用少量标注数据进行训练的监督学习模型。本文的目标是利用多个基学习器构建获得安全的预测结果 $\mathbf{f} = g(\{\mathbf{f}_1, \dots, \mathbf{f}_b\}, \mathbf{f}_0)$ ，使其能够在通常情况下优于基线模型 \mathbf{f}_0 ，同时在最坏情况下也不会比 \mathbf{f}_0 更差。

由于 SafeW 是面向通用半监督学习的框架，我们同时考虑分类任务和回归任务，对于分类任务， $\mathbb{H} = \{+1, -1\}$ ，对于回归任务， $\mathbb{H} = \mathbb{R}$ ，本章中采用的符号及含义如表 4-1。

表 4-1: 本章采用的符号及含义总结。

符号	含义
u	无标注样本个数
b	模型个数
\mathbb{H}	标注空间, 对于分类任务 $\mathbb{H} = \{+1, -1\}$, 对于回归任务 $\mathbb{H} = \mathbb{R}$
f_1, \dots, f_b	b 个基学习模型
$\mathbf{f}_1, \dots, \mathbf{f}_b \in \mathbb{H}^u$	模型在无标注数据上的预测结果
$\mathbf{f}_0 \in \mathbb{H}^u$	基线的监督学习模型预测结果
$\mathbf{f}^* \in \mathbb{H}^u$	无标注数据的真实标注 (未知)
$\hat{\mathbf{f}} \in \mathbb{H}^u$	SafeW 预测结果
$\ell(\cdot, \cdot)$	损失函数
α	各模型权重
\mathcal{M}	权重 α 的取值范围
\mathbf{C}^{clf}	分类任务中 b 个基学习器的协方差矩阵
\mathbf{C}^{reg}	回归任务中 b 个基学习器的协方差矩阵

4.3.2 安全的半监督学习框架 SafeW

我们首先考虑一个简单的情况, 即, 所有无标注样本的真实标注都是已知的。令 \mathbf{f}^* 表示无标注样本的真实标注, 我们的学习目标是得到模型在无标注数据上的预测结果 \mathbf{f} , 使其相比于基线模型 \mathbf{f}_0 取得最大的性能提升, 该目标可以通过优化如下目标实现:

$$\max_{\mathbf{f} \in \mathbb{H}^u} \ell(\mathbf{f}_0, \mathbf{f}^*) - \ell(\mathbf{f}, \mathbf{f}^*) \quad (4-1)$$

其中 $\ell(\cdot, \cdot)$ 表示损失函数, 损失函数的值越小, 表示模型性能越好。表 4-2 总结了常用的用于分类和回归任务的损失函数。

然而, 在现实任务中显然真实标注 \mathbf{f}^* 是未知的, 为了解决该问题, 我们提出采用多个半监督学习模型的集成结果近似真实标注。具体而言, 令 $\mathbf{f}^* = \sum_{i=1}^b \alpha_i \mathbf{f}_i$, 其中 $\alpha = [\alpha_1; \alpha_2; \dots; \alpha_b] \geq \mathbf{0}$ 表示多个基学习器的模型权重, 并且满足 $\sum_{i=1}^b \alpha_i = 1$, 然后我们可以将式 4-1 中的 \mathbf{f}^* 用集成模型的预测结果替代, 得到如下目标:

$$\max_{\mathbf{f} \in \mathbb{H}^u} \ell\left(\mathbf{f}_0, \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) - \ell\left(\mathbf{f}, \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) \quad (4-2)$$

表 4-2: 分类任务及回归任务常用损失函数 $\ell(\mathbf{p}, \mathbf{q})$ 。 $\mathbf{q} = [q_1; \dots; q_u] \in \mathbb{R}^u$ 为模型预测结果, $\mathbf{p} = [p_1; \dots; p_u] \in \mathbb{H}^u$ 为样本真实标注。对于分类任务, $\mathbb{H}^u = \{+1, -1\}^u$, 对于回归任务, $\mathbb{H}^u = \mathbb{R}^u$ 。 η 表示损失函数的利普希茨常数, 对于回归任务, $M = \max\{|a|, |b|\}$, $[a, b]$ 表示标注的取值范围。

损失函数	$\ell(\mathbf{p}, \mathbf{q})$ 定义	任务	η
铰链损失	$\frac{1}{u} \sum_{i=1}^u \max\{1 - p_i q_i, 0\}$	分类	1
交叉熵损失	$\frac{1}{u} \sum_{i=1}^u -p_i \ln(q_i) - (1 - p_i) \ln(1 - q_i)$	分类	1
均方损失	$\frac{1}{u} \sum_{i=1}^u (p_i - q_i)^2 = \frac{1}{u} (1 - \mathbf{p}\mathbf{q})^2$	分类	4
均方损失	$\frac{1}{u} \sum_{i=1}^u (p_i - q_i)^2 = \frac{1}{u} \ \mathbf{p} - \mathbf{q}\ _2^2$	回归	$2 + M$
平均绝对损失	$\frac{1}{u} \sum_{i=1}^u p_i - q_i = \frac{1}{u} \ \mathbf{p} - \mathbf{q}\ _1$	回归	1
ϵ 不敏感损失	$\frac{1}{u} \sum_{i=1}^u \max\{ p_i - q_i - \epsilon, 0\}$	回归	1

在实践中, 我们无法得知每个基学习器准确的权重, 因此, 我们假设权重 α 来自一个凸集 \mathcal{M} , 其中 \mathcal{M} 包含了关于每个模型可靠性的先验知识, 在后续内容中我们会进一步针对分类任务和回归任务讨论集合 \mathcal{M} 的设置。

由于获取的半监督学习模型具有不确定性, 为了保证最终预测结果的安全性, 我们提出优化最坏情况下的性能增益。直观来讲, 如果在最坏的情况下, 模型相比与基线的监督学习模型仍然能取得性能提升, 则可以认为该模型是安全的, 基于此, 我们可以得到如下最大最小的优化目标:

$$\max_{\mathbf{f} \in \mathbb{H}^u} \min_{\alpha \in \mathcal{M}} \ell\left(\mathbf{f}_0, \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) - \ell\left(\mathbf{f}, \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) \quad (4-3)$$

该目标可以应用于任意半监督、弱监督学习模型的预测结果, 也可以同时应用于分类任务及回归任务, 因此这是一个通用的安全半(弱)监督学习框架。

关于基学习器权重 α 的取值集合如何设置的问题, 我们可以简单地将 \mathcal{M} 设置为如下单纯形, 即,

$$\mathcal{M} = \left\{ \alpha \mid \sum_{i=1}^b \alpha_i = 1, \alpha \geq 0 \right\} \quad (4-4)$$

类似的设置方法在其它集成学习工作, 如 [89, 91] 中已有应用。但是该策略的局限性是没有考虑任何关于基学习器性能的先验知识, 过于保守。显然, \mathcal{M} 的设置可以很容易嵌入各种先验知识, 例如, 如果已知基学习器 \mathbf{f}_i 相比 \mathbf{f}_j 更为可靠, 并且所有满足此类关系的模型索引集合为 \mathcal{S} , 那么 \mathcal{M} 可以设置为 $\{\alpha \mid \alpha_i - \alpha_j \geq 0, (i, j) \in \mathcal{S}; \alpha^\top \mathbf{1} = 1; \alpha \geq \mathbf{0}\}$, 其中 $\mathbf{1}(\mathbf{0})$ 分别表示全为 1 或者全为 0 的

向量；如果已知基学习器的重要性程度，记为 $\{r_1, \dots, r_b\}$ ，那么 \mathbf{M} 可以设置为 $\{\boldsymbol{\alpha} | -\gamma \leq \alpha_i - r_i \leq \gamma, \forall i = 1, \dots, b; \boldsymbol{\alpha}^\top \mathbf{1} = 1; \boldsymbol{\alpha} \geq \mathbf{0}\}$ ，其中 γ 为一个取值很小的常数。以上的权重设置方法依赖基学习器的先验知识，然而，现实任务中依然存在无法获取类似知识的情况，这种情况下，一个简单的方法是通过交叉验证的方法对多个基学习器的性能进行评估。然而，交叉验证方法非常耗时，并且在半监督学习中，由于标注样本的不足，往往无法获得可靠的模型验证结果。因此，我们在本文中进一步提出一种从数据中自动学习模型权重的方法。

首先，对于回归任务，令 $\mathbf{C}^{reg} \in \mathbb{R}^{b \times b}$ 表示 b 个基学习器 $\{f_1, \dots, f_b\}$ 在整个数据分布 (X, Y) 上的协方差矩阵，其元素为：

$$C_{ij}^{reg} = \mathbb{E}_X[(f_i(X) - \mu_i)^\top (f_j(X) - \mu_j)] \quad (4-5)$$

其中 μ_i 表示预测结果的均值，即 $\mu_i = \mathbb{E}_X[f_i(X)]$ 。

令 $\boldsymbol{\rho}^{reg} = [\rho_1^{reg}; \dots; \rho_b^{reg}]$ 表示基学习器预测结果与真实标注指派 f^* 之间的协方差，即，

$$\rho_i^{reg} = \mathbb{E}_X[(f^*(X) - \mu^*)^\top (f_i(X) - \mu_i)] \quad (4-6)$$

其中 $\mu^* = \mathbb{E}_X[f^*(X)]$ 。

我们通过最小化加权集成的预测结果与真实预测结果之间的误差来优化权重 $\boldsymbol{\alpha}$ ，其目标如下所示：

$$\boldsymbol{\alpha}^* = \arg \min_{\boldsymbol{\alpha}} \mathbb{E}_X \left[\text{MSE} \left(\sum_{i=1}^b \alpha_i f_i(X), f^*(X) \right) \right] \quad (4-7)$$

其中 MSE 表示均方误差 (Mean Squared Error)。

关于上式的闭式解，有如下定理：

定理 4-1 (Bates and Granger, 1969 [6]) 最优的 $\boldsymbol{\alpha}^*$ 满足下列条件：

$$\boldsymbol{\rho}^{reg} = \mathbf{C}^{reg} \boldsymbol{\alpha}^* \quad (4-8)$$

基于以上定理我们可知求解最优的权重 $\boldsymbol{\alpha}^*$ 需要估计 \mathbf{C}^{reg} 和 $\boldsymbol{\rho}^{reg}$ 的值。对

于 \mathbf{C}^{reg} ，由于完整的数据分布未知，我们可以基于训练数据来获得 \mathbf{C}^{reg} 的无偏估计 $\hat{\mathbf{C}}^{reg}$ ，即，

$$\hat{\mathbf{C}}_{ij}^{reg} = (\mathbf{f}_i - \hat{\boldsymbol{\mu}}_i)^\top (\mathbf{f}_j - \hat{\boldsymbol{\mu}}_j) \quad (4-9)$$

其中 $\hat{\boldsymbol{\mu}}_i = \frac{1}{b} \sum_{i=1}^b \mathbf{f}_i$ 。

对于 $\boldsymbol{\rho}^{reg}$ ，以下结论表明其取值与基学习器的性能密切相关。

定理 4-2 假设模型预测结果 $\{f_i(X)\}_{i=1}^b$ 的均值满足 $\mu_i = 0, \forall i = 1, \dots, b$ ，方差为 1，当损失函数为均方损失时， ρ_i^{reg} 越大，模型 f_i 的损失越小。

证明: 对于 $\boldsymbol{\rho}^{reg}$ ，我们有：

$$\rho_i^{reg} = \mathbb{E}_X[(f^*(X) - \mu^*)^\top (f_i(X) - \mu_i)] = \mathbb{E}_X[(f^*(X))^\top f_i(X)] \quad (4-10)$$

对于均方误差，我们有：

$$\begin{aligned} \text{MSE}(f^*(X), f_i(X)) & \quad (4-11) \\ &= \mathbb{E}_X[(f^*(X) - f_i(X))^2] \\ &= \mathbb{E}_X[\|f^*(X)\|^2 + \|f_i(X)\|^2 - 2f^*(X)^\top f_i(X)] \\ &= 2 - 2\mathbb{E}_X[f^*(X)^\top f_i(X)] \\ &= 2 - 2\rho_i^{reg} \end{aligned}$$

由以上两式可知， ρ_i^{reg} 的值越大，模型 f_i 对应的损失越小。 \square

根据上述结论，我们可以将 \mathcal{M} 设置为 $\{\boldsymbol{\alpha} | \hat{\mathbf{C}}^{reg} \boldsymbol{\alpha} \geq \mathbf{1}\delta, \boldsymbol{\alpha}^\top \mathbf{1} = 1, \boldsymbol{\alpha} \geq \mathbf{0}\}$ ，其中 δ 为一个常数，表示基学习器的性能下界，例如，优于随机猜测的模型。非常容易验证集合 \mathcal{M} 为一个凸集合。

对于分类任务，类似地，令 $\mathbf{C}^{clf} \in \mathbb{R}^{b \times b}$ 表示 b 个基学习器预测结果的协方差，其元素为：

$$\mathbf{C}_{ij}^{clf} = \mathbb{E}_X[f_i(X)^\top f_j(X)] \quad (4-12)$$

令 $\boldsymbol{\rho}^{clf} = [\rho_1^{clf}; \rho_2^{clf}; \dots; \rho_b^{clf}]$ 表示模型预测结果与真实标注之间的协方差，

即,

$$\rho_i^{clf} = \mathbb{E}_X [f^*(X)^\top f_i(X)] \quad (4-13)$$

当以分类正确率 (Accuracy) 作为评价指标时, 我们可以得到如下结论:

定理 4-3 分类任务中的最优权重 α^* 满足:

$$\rho^{clf} = \mathbf{C}^{clf} \alpha^* \quad (4-14)$$

类似地, 我们可以基于训练数据获得 \mathbf{C}^{clf} 的无偏估计 $\hat{\mathbf{C}}^{clf}$, 其元素为:

$$\hat{C}_{ij}^{clf} = \mathbf{f}_i^\top \mathbf{f}_j \quad (4-15)$$

根据上述结论, 与回归任务类似, 我们可以将 \mathcal{M} 设置为 $\{\alpha | \hat{\mathbf{C}}^{clf} \alpha \geq \mathbf{1}\delta, \alpha^\top \mathbf{1} = 1, \alpha \geq \mathbf{0}\}$, δ 为一个常数, 表示基学习器的性能下界。同样可以验证 \mathcal{M} 是一个凸集合。

综上所述, 一方面如果可以获取到与基学习器性能有关的先验知识, SafeW 框架可以直接嵌入这些先验知识, 另一方面, 当无法获取先验知识时, 该框架仍然能够有效利用回归任务与分类任务的协方差矩阵分析获得的权重估计。这进一步说明了 SafeW 框架的灵活性与通用性。

接下来, 我们讨论如何获得 SafeW 的优化目标式 4-3 的最优解。原始目标函数是由两个损失函数相减得到的, 通常来说, 该目标式是非凸的, 而对于非凸优化问题, 求解其全局最优解是非常困难的 [13]。幸运的是, 在本文中, 我们发现对于一类常用的损失函数, 目标式 4-3 可以等价转换为一个凸优化问题, 从而高效地获得全局最优解, 基于此, 我们分别给出 SafeW 对于回归问题和分类问题常用损失函数的优化算法。

首先, 对于回归问题, 我们有如下定理:

定理 4-4 对于回归任务, 假设 $\ell(\cdot, \sum_{i=1}^b \alpha_i \mathbf{f}_i)$ 对 α 是凸的, 并且存在 $\mathbf{f} \in \mathbb{R}^u$ 满足 $\ell(\mathbf{f}, \sum_{i=1}^b \alpha_i \mathbf{f}_i) = 0$, 此时, 目标式 4-3 为一个凸优化问题。

在证明定理 4-4 之前，我们首先给出下列引理，

引理 4-5 在定理 4-4 的条件下，最优解 $\hat{\mathbf{f}}$ 和 $\hat{\alpha}$ 满足下列关系：

$$\ell\left(\hat{\mathbf{f}}, \sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i\right) = 0 \quad (4-16)$$

证明：我们利用反证法进行上述引理的证明。假设 $\ell(\hat{\mathbf{f}}, \sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i) \neq 0$ ，根据条件，存在 $\tilde{\mathbf{f}}$ 满足 $\ell(\tilde{\mathbf{f}}, \sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i) = 0$ ，显然， $0 = \ell(\tilde{\mathbf{f}}, \sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i) < \ell(\hat{\mathbf{f}}, \sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i)$ 。因此， $\hat{\mathbf{f}}$ 不是最优解，出现矛盾。 \square

接下来，我们证明定理 4-4。

证明：根据引理 4-5，对于分类任务，目标式 4-3 可以写为：

$$\min_{\alpha \in \mathcal{M}} \ell\left(\mathbf{f}_0, \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) \quad (4-17)$$

由于 $\ell(\cdot, \sum_{i=1}^b \alpha_i \mathbf{f}_i)$ 关于 α 为凸，因此，式 4-3 为一个凸优化问题。 \square

评注 4-1 定理 4-4 假设损失函数是关于权重 α 的凸函数，该假设非常容易满足。大多数回归任务的损失函数，如均方损失、 ϵ 不敏感损失，以及平均绝对损失 [133] 均满足该条件。

基于引理 4-5 和定理 4-4 可知，对于回归任务，目标式 4-3 可以高效地获得全局最优解。接下来，我们以回归任务中最常用的损失函数均方损失为例展示 SafeW 的具体优化方法。

当以均方损失作为损失函数时，SafeW 的目标式 4-3 可以写作如下只与 α 相关的等价形式：

$$\min_{\alpha \in \mathcal{M}} \|\sum_{i=1}^b \alpha_i \mathbf{f}_i - \mathbf{f}_0\|^2 \quad (4-18)$$

进一步，将公式 4-18 中的二次项进行展开，可以得到：

$$\min_{\alpha \in \mathcal{M}} \alpha^\top \mathbf{F} \alpha - \mathbf{v}^\top \alpha \quad (4-19)$$

算法 4.1 SafeW 框架在回归问题上的优化流程。

输入： 多个基学习器预测结果 $\{\mathbf{f}_i\}_{i=1}^b$ ，监督学习器预测结果 \mathbf{f}_0

输出： SafeW 的预测结果 $\hat{\mathbf{f}}$

- 1: 构建线性核矩阵 \mathbf{F} ，其中， $F_{ij} = \mathbf{f}_i^\top \mathbf{f}_j, \forall 1 \leq i, j \leq b$
- 2: 推导出 $\mathbf{v} = [2\mathbf{f}_1^\top \mathbf{f}_0; \dots; 2\mathbf{f}_b^\top \mathbf{f}_0]$
- 3: 求解凸二次优化问题公式 4-19，得到最优权重 $\boldsymbol{\alpha}^* = [\alpha_1^*, \dots, \alpha_b^*]$
- 4: 返回 $\hat{\mathbf{f}} = \sum_{i=1}^b \alpha_i^* \mathbf{f}_i$ 。

其中 $\mathbf{F} \in \mathbb{R}^{b \times b}$ 是预测结果 \mathbf{f}_i 的线性核矩阵，即， $F_{ij} = \mathbf{f}_i^\top \mathbf{f}_j$ ， $\mathbf{v} = [2\mathbf{f}_1^\top \mathbf{f}_0; \dots; 2\mathbf{f}_b^\top \mathbf{f}_0]$ 。因为 \mathbf{F} 是半正定矩阵，所以式 4-19 是一个凸二次优化问题 [13]，可以高效地获得全局最优解。我们可以借助开源的凸优化工具包，例如 MOSEK^①包，对该问题进行求解。

在求解得到最优的权重 $\boldsymbol{\alpha}^*$ 之后，可以得到 SafeW 算法的最终预测结果为：

$$\hat{\mathbf{f}} = \sum_{i=1}^b \alpha_i^* \mathbf{f}_i \quad (4-20)$$

算法 4.1总结了对于回归任务 SafeW 的优化流程。

值得一提的是，公式 4-18可以看作一个几何投影问题，具体而言，令 $\Omega = \{\mathbf{f} | \sum_{i=1}^b \alpha_i \mathbf{f}_i, \boldsymbol{\alpha} \in \mathcal{M}\}$ ，公式 4-18可以写作：

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f} \in \Omega} \|\mathbf{f} - \mathbf{f}_0\|^2 \quad (4-21)$$

该目标是一个学习 \mathbf{f}_0 到集合 Ω 上的投影的过程。

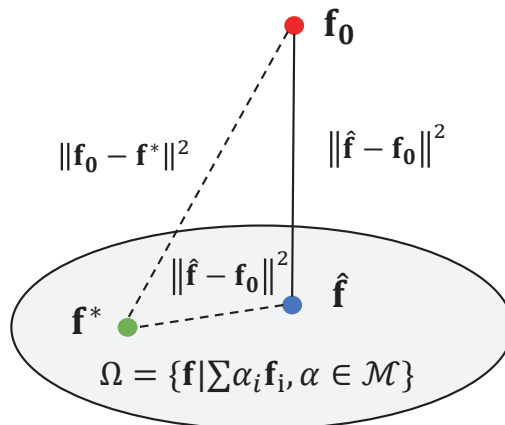


图 4-2: SafeW 优化算法的投影解释。直观上讲，该算法学习一个 \mathbf{f}_0 到凸集合 Ω 的投影。

^①<https://www.mosek.com/resources/downloads>

图 4-2 从投影的角度展示了回归任务中 SafeW 的优化过程。根据 Pythagorean 定理 [19], 如果 $\mathbf{f}^* \in \Omega$, 则 $\|\hat{\mathbf{f}} - \mathbf{f}^*\|$ 之间的距离应当小于 $\|\mathbf{f}_0 - \mathbf{f}^*\|$, 该结论证明了当真实标注指派可以由多个基学习器预测结果组合构造得到时, SafeW 可以实现安全性, 这从几何投影角度给了我们关于安全半监督学习的新理解。

对于分类问题, 由于模型输出空间的不连续性, 不能直接将回归任务中的引理 4-5 应用到分类任务中。我们进一步对分类任务中的优化问题展开分析, 针对分类任务中的常用损失, 铰链损失, 证明了可以将目标式 4-3 转化为线性优化问题从而高效获得最优解, 对于另一种深度学习中流行的分类损失, 交叉熵损失 (Cross-Entropy Loss), 也可以通过一个简单的凸松弛技术转化为凸优化问题, 从而高效求得最优解。

首先, 我们提出如下引理:

引理 4-6 对于分类任务, 最优的 $\hat{\mathbf{f}}$ 和 $\hat{\alpha}$ 满足如下关系:

$$\hat{\mathbf{f}} = \text{sign}\left(\sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i\right) \quad (4-22)$$

如果 $s \geq 0$, $\text{sign}(s) = 1$, 反之, $\text{sign}(s) = 0$ 。

证明: 我们采用反证法证明该引理。假设 $\hat{\mathbf{f}} \neq \text{sign}\left(\sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i\right)$, 根据条件, 我们可知存在 $\tilde{\mathbf{f}}$ 满足 $\tilde{\mathbf{f}} = \text{sign}\left(\sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i\right)$, 那么, 显然, $\ell(\tilde{\mathbf{f}}, \sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i) < \ell(\hat{\mathbf{f}}, \sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i)$ 。因此, $\hat{\mathbf{f}}$ 不是最优解, 存在矛盾, 因此假设不正确。 \square

根据引理 4-6, 我们可以得到下述定理,

定理 4-7 如果 $\mathbf{f}_i \in \{+1, -1\}^u, \forall i = 1, \dots, b$, 损失函数 $\ell(\cdot, \cdot)$ 为铰链损失, 则 SafeW 的优化目标 4-3 是一个线性规划问题。

证明: 根据引理 4-6, 式 4-3 可以写为

$$\min_{\alpha \in \mathcal{M}} \ell\left(\mathbf{f}_0, \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) - \ell\left(\text{sign}\left(\sum_{i=1}^b \alpha_i \mathbf{f}_i\right), \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) \quad (4-23)$$

因为 $\mathbf{f}_i \in \{+1, -1\}^u, \forall i = 1, \dots, b$, 并且 $\ell(\cdot, \sum_{i=1}^b \alpha_i \mathbf{f}_i)$ 对预测结果是线性的,

那么,

$$\ell\left(\text{sign}\left(\sum_{i=1}^b \alpha_i \mathbf{f}_i\right), \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) \quad (4-24)$$

可以等价转换为

$$\ell\left(\left\|\sum_{i=1}^b \alpha_i \mathbf{f}_i\right\|_1\right) \quad (4-25)$$

因此, 式 4-23 等价于

$$\min_{\alpha \in \mathcal{M}} \ell\left(\mathbf{f}_0, \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) + \ell\left(\left\|\sum_{i=1}^b \alpha_i \mathbf{f}_i\right\|_1\right) \quad (4-26)$$

令 $\tilde{\mathbf{f}} = \sum_{i=1}^b \alpha_i \mathbf{f}_i$, 式 4-26 可以写为:

$$\begin{aligned} \min_{\alpha \in \mathcal{M}} \quad & \ell(\mathbf{f}_0, \tilde{\mathbf{f}}) + \ell(\|\tilde{\mathbf{f}}\|_1) \\ \text{s.t.} \quad & \tilde{\mathbf{f}} = \sum_{i=1}^b \alpha_i \mathbf{f}_i \end{aligned} \quad (4-27)$$

通过引入两个辅助变量:

$$\mathbf{z} = \frac{|\tilde{\mathbf{f}}| + \tilde{\mathbf{f}}}{2}, \quad \mathbf{w} = \frac{|\tilde{\mathbf{f}}| - \tilde{\mathbf{f}}}{2} \quad (4-28)$$

式 4-27 可以等价转换为:

$$\begin{aligned} \min_{\alpha \in \mathcal{M}, \mathbf{z}, \mathbf{w}} \quad & \ell(\mathbf{f}_0, \tilde{\mathbf{f}}) + \ell(\mathbf{1}^\top (\mathbf{z} + \mathbf{w})) \\ \text{s.t.} \quad & \tilde{\mathbf{f}} = \sum_{i=1}^b \alpha_i \mathbf{f}_i \\ & \tilde{\mathbf{f}} + \mathbf{z} - \mathbf{w} = \mathbf{0} \\ & \mathbf{z} \geq \mathbf{0}, \mathbf{w} \geq \mathbf{0} \end{aligned} \quad (4-29)$$

因为损失函数 $\ell(\cdot, \tilde{\mathbf{f}})$ 是关于 $\tilde{\mathbf{f}}$ 是线性函数, 所以上式优化目标和约束条件对于变量 α 、 \mathbf{z} 、 \mathbf{w} 都是线性的。由此可知, 式 4-29 是一个线性规划问题。□

式 4-29 同样可以通过 MOSEK 等优化工具高效地获得全局最优解, 在得到最优权重 α^* 之后, 最优的预测结果 $\hat{\mathbf{f}} = \sum_{i=1}^b \alpha_i^* \mathbf{f}_i$ 。算法 4.2 总结了对于分类任务

算法 4.2 SafeW 框架在分类任务上的优化流程。

输入: 多个基学习器预测结果 $\{\mathbf{f}_i\}_{i=1}^b$, 监督学习器预测结果 \mathbf{f}_0

输出: SafeW 的预测结果 $\hat{\mathbf{f}}$

- 1: 构建 $\tilde{\mathbf{f}} = \sum_{i=1}^b \alpha_i \mathbf{f}_i$
 - 2: 引入 $\mathbf{z} = (|\tilde{\mathbf{f}}| + \tilde{\mathbf{f}})/2$, $\mathbf{w} = (|\tilde{\mathbf{f}}| - \tilde{\mathbf{f}})/2$
 - 3: 求解线性优化问题公式 4-29, 得到最优解 $\alpha^* = [\alpha_1^*, \dots, \alpha_b^*]$
 - 4: 返回 $\hat{\mathbf{f}} = \sum_{i=1}^b \alpha_i^* \mathbf{f}_i$
-

SafeW 的优化流程。

接下来我们证明, 上述凸问题的性质对于深度学习中最常用的分类损失交叉熵损失 [53] 依然成立。

首先, 令

$$\hat{\ell}(p) = \begin{cases} \ln(p) & 0.5 \leq p \leq 1 \\ \ln(1-p) & 0 \leq p < 0.5 \end{cases} \quad (4-30)$$

可以发现, 当损失函数 $\ell(\cdot, \cdot)$ 为交叉熵损失时下式成立,

$$-\ell\left(\text{sign}\left(\sum_{i=1}^b \alpha_i \mathbf{f}_i\right), \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) = \sum_{j=1}^u \hat{\ell}\left(\left(\sum_{i=1}^b \alpha_i \mathbf{f}_i\right)_j\right) \quad (4-31)$$

其中 $(\sum_{i=1}^b \alpha_i \mathbf{f}_i)_j$ 表示 $\sum_{i=1}^b \alpha_i \mathbf{f}_i$ 中的第 j 个元素。

令

$$g(p) = \begin{cases} (2 \ln 2)p - 2 \ln 2 & 0.5 \leq p \leq 1 \\ -(2 \ln 2)p & 0 \leq p < 0.5 \end{cases} \quad (4-32)$$

可以发现, $g(p)$ 实现了 $\hat{\ell}(p)$ 的一个凸松弛, 且 $g(p)$ 是 $\hat{\ell}(p)$ 的一个凸包。

基于上述观察, 我们可以得到如下定理:

定理 4-8 令 $\tilde{\mathbf{f}} = \sum_{i=1}^b \alpha_i \mathbf{f}_i$, 当损失函数为交叉熵损失时, 优化问题

$$\min_{\alpha} \ell(\mathbf{f}_0, \tilde{\mathbf{f}}) + \sum_{j=1}^u g(\tilde{f}_j) \quad (4-33)$$

是目标式 4-3 的一个凸松弛。

证明: 根据引理 4-6, 最优的 \mathbf{f} 为 $\text{sign}(\sum_{i=1}^b \alpha_i \mathbf{f}_i)$, 因此, 目标式 4-3 可以等价写

为

$$\min_{\alpha} \ell(\mathbf{f}_0, \tilde{\mathbf{f}}) + \sum_{j=1}^u \hat{\ell}(\tilde{\mathbf{f}}_j) \quad (4-34)$$

因为 $\ell(\mathbf{f}_0, \tilde{\mathbf{f}})$ 是凸函数，并且 $g(p)$ 是 $\hat{\ell}(p)$ 的一个凸包，所以，式 4-33 是凸的，并且是式 4-3 的一个凸松弛。□

同样地，得到最优的权重 α^* 之后，最优的预测结果 $\hat{\mathbf{f}} = \sum_{i=1}^b \alpha_i^* \mathbf{f}_i$ 。此外，值得一提的是，以上的优化技巧并不单适用于交叉熵损失，其它相似的凸损失函数也可以使用。

综上所述，对于多种回归任务和分类任务中常用的损失函数，SafeW 都可以将其转换为凸优化问题从而高效获得全局最优解，这表明了 SafeW 方法的通用性和有效性。

4.4 理论分析

在本节，我们从理论上对 SafeW 框架的安全性进行分析，证明了 SafeW 的优化目标式 4-3 对于多种常用的凸损失函数，如表 4-2 中所示，可以实现安全性。

首先，我们给出如下结论：

定理 4-9 (安全性) 假设无标注样本的真实标注 \mathbf{f}^* 可以通过多个基学习器预测结果的组合构造得出，即， $\mathbf{f}^* \in \{\mathbf{f} | \sum_{i=1}^b \alpha_i \mathbf{f}_i, \alpha \in \mathcal{M}\}$ 。令 $\hat{\mathbf{f}}$ 和 $\hat{\alpha}$ 表示目标式 4-3 的最优解，我们有：

$$\ell(\hat{\mathbf{f}}, \mathbf{f}^*) \leq \ell(\mathbf{f}_0, \mathbf{f}^*) \quad (4-35)$$

并且， $\hat{\mathbf{f}}$ 相比于 \mathbf{f}_0 已经实现了最优的性能提升。

证明：首先，我们定义

$$L(\mathbf{f}, \alpha) = \ell\left(\mathbf{f}_0, \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) - \ell\left(\mathbf{f}, \sum_{i=1}^b \alpha_i \mathbf{f}_i\right) \quad (4-36)$$

因为目标式 4-3 式一个最大最小优化目标，所以下列不等式对于任意的 \mathbf{f} 与 α 成立：

$$L(\mathbf{f}, \hat{\alpha}) \leq L(\hat{\mathbf{f}}, \hat{\alpha}) \leq L(\hat{\mathbf{f}}, \alpha) \quad (4-37)$$

如果 α^* 的取值使得 $\mathbf{f}^* = \sum_{i=1}^b \alpha_i^* \mathbf{f}_i$ 成立, 通过将 \mathbf{f} 和 α 设置为 \mathbf{f}_0 和 α^* 的取值, 我们有:

$$\ell\left(\mathbf{f}_0, \sum_{i=1}^b \hat{\alpha}_i \mathbf{f}_i\right) - \ell\left(\mathbf{f}_0, \sum_{i=1}^b \alpha_i^* \mathbf{f}_i\right) \leq \ell\left(\mathbf{f}_0, \sum_{i=1}^b \alpha_i^* \mathbf{f}_i\right) - \ell\left(\hat{\mathbf{f}}, \sum_{i=1}^b \alpha_i^* \mathbf{f}_i\right) \quad (4-38)$$

因此, 下式成立:

$$\ell(\hat{\mathbf{f}}, \mathbf{f}^*) \leq \ell(\mathbf{f}_0, \mathbf{f}^*) \quad (4-39)$$

此外, 由于我们已经在最坏的情况下实现了性能增益的最大化, $\hat{\mathbf{f}}$ 已经实现了相比于 \mathbf{f}_0 的最大性能增益。□

评注 4-2 定理 4-9 证明了本文提出的优化目标 4-3 是合理的优化目标, 即, 通过优化该目标得出的最优解 $\hat{\mathbf{f}}$ 可以实现性能永远不会弱于简单的监督学习模型 \mathbf{f}_0 。同时, 相比于以往的安全半监督学习工作, 如 [4, 91, 92], 该优化目标具有如下优势: 相比于 [92] 假设无标注样本的真实标注可以由某个基学习器预测得到, 我们只假设真实标注由多个基学习器的预测结果组合得到, 该假设更容易实现。相比于 [4], 我们在本文中显式的考虑最大化相比于基线模型的性能增益, 可以实现性能增益最大化; 相比于 [91] 只能应用于回归任务, 我们的方法更为通用, 在回归任务以及分类任务中均可适用。

进一步, 假设模型的损失函数 $\ell(\cdot, \cdot)$ 满足利普希茨连续性, 即,

$$\|\ell(\mathbf{f}_1, \mathbf{f}_2) - \ell(\mathbf{f}_1, \mathbf{f}_3)\| \leq \eta \|\mathbf{f}_2 - \mathbf{f}_3\|, \quad \forall \mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3 \in [-1, 1] \quad (4-40)$$

令 $\beta^* = [\beta_1^*, \dots, \beta_b^*] \in \mathcal{M}$ 表示优化目标 4-3 的最优解, 即,

$$\beta^* = \arg \min_{\beta \in \mathcal{M}} \ell\left(\sum_{i=1}^b \beta_i \mathbf{f}_i, \mathbf{f}^*\right) \quad (4-41)$$

ϵ 表示与真实标注之间的偏差,

$$\epsilon = \mathbf{f}^* - \sum_{i=1}^b \beta_i^* \mathbf{f}_i \quad (4-42)$$

我们可以得到如下结论：

定理 4-10 SafeW 框架的预测结果 $\hat{\mathbf{f}}$ 相比于基线模型的预测结果 \mathbf{f}_0 的性能增益，即， $\ell(\mathbf{f}_0, \mathbf{f}^*) - \ell(\hat{\mathbf{f}}, \mathbf{f}^*)$ ，有性能下界 $-2\eta\|\epsilon\|_1$ 。

证明：注意到， $\sum_{i=1}^b \beta_i^* \mathbf{f}_i \in \{\mathbf{f} \mid \sum_{i=1}^b \alpha_i \mathbf{f}_i, \alpha \in \mathcal{M}\}$ ，根据定理 4-9，我们有，

$$\ell\left(\mathbf{f}_0, \sum_{i=1}^b \beta_i^* \mathbf{f}_i\right) - \ell\left(\hat{\mathbf{f}}, \sum_{i=1}^b \beta_i^* \mathbf{f}_i\right) \geq 0 \quad (4-43)$$

因为 $\mathbf{f}^* = \sum_{i=1}^b \beta_i^* \mathbf{f}_i + \epsilon$ ，可以得到：

$$\left| \ell(\hat{\mathbf{f}}, \mathbf{f}^*) - \ell\left(\hat{\mathbf{f}}, \sum_{i=1}^b \beta_i^* \mathbf{f}_i\right) \right| \leq \eta\|\epsilon\|_1 \quad (4-44)$$

该不等式成立是因为损失函数满足 η 利普希茨连续性。

相似地，我们有：

$$\left| \ell(\mathbf{f}_0, \mathbf{f}^*) - \ell\left(\mathbf{f}_0, \sum_{i=1}^b \beta_i^* \mathbf{f}_i\right) \right| \leq \eta\|\epsilon\|_1 \quad (4-45)$$

该不等式意味着下列两个不等式成立：

$$-\eta\|\epsilon\|_1 \leq \ell(\hat{\mathbf{f}}, \mathbf{f}^*) - \ell\left(\hat{\mathbf{f}}, \sum_{i=1}^b \beta_i^* \mathbf{f}_i\right) \leq \eta\|\epsilon\|_1 \quad (4-46)$$

$$-\eta\|\epsilon\|_1 \leq \ell(\mathbf{f}_0, \mathbf{f}^*) - \ell\left(\mathbf{f}_0, \sum_{i=1}^b \beta_i^* \mathbf{f}_i\right) \leq \eta\|\epsilon\|_1 \quad (4-47)$$

将上述两个不等式合并化简，我们可以得到：

$$\begin{aligned} & \ell(\mathbf{f}_0, \mathbf{f}^*) - \ell(\hat{\mathbf{f}}, \mathbf{f}^*) \\ & \geq \left(\ell\left(\mathbf{f}_0, \sum_{i=1}^b \beta_i^* \mathbf{f}_i\right) - \eta\|\epsilon\|_1 \right) - \left(\ell\left(\hat{\mathbf{f}}, \sum_{i=1}^b \beta_i^* \mathbf{f}_i\right) + \eta\|\epsilon\|_1 \right) \\ & \geq -2\eta\|\epsilon\|_1 \end{aligned} \quad (4-48)$$

其中第二个不等式成立是因为 $\ell(\mathbf{f}_0, \sum_{i=1}^b \beta_i^* \mathbf{f}_i) - \ell(\hat{\mathbf{f}}, \sum_{i=1}^b \beta_i^* \mathbf{f}_i) \geq 0$ 。 \square

评注 4-3 上述结论假设损失函数满足利普希茨连续性假设，事实上，机器学习中大多数常用的损失函数均满足该假设，我们在表 4-2中总结了常用损失函数的利普希茨连续性常数值 η 。

评注 4-4 定理 4-10揭示了最坏情况下的模型性能表现只与基学习器的质量有关，而与基学习器的数量无关。此外，值得一提的是，定理 4-9只给出了安全性的充分条件，而不是必要条件，同样，定理 4-10只给出了性能的下限，而不是具体的性能表现。换言之，即使定理 4-9中的条件不成立，我们的方法依然可以有可能实现安全的性能，大量实验结果证明了该结论。

4.5 实验验证

在本节，我们首先在半监督学习场景进行实验，验证 SafeW 框架的安全性，然后，我们进一步在其它弱监督学习场景进行综合实验，包括领域自适应学习、多示例学习、标注噪声学习，大量的实验结果充分证明了本文提出的 SafeW 框架的通用性和有效性。

4.5.1 半监督学习实验

数据介绍。对于半监督学习，我们在回归任务中采用大量数据集进行实验，包括 *abalone*、*bodyfat*、*cadata*、*cpusmall*、*eunite2001*、*housing*、*mg*、*mpg*、*pyrim*、*space_ga* 共 10 种数据集^①，其中，数据集涉及的领域包括物理测量 (*abalone*)、医疗健康 (*bodyfat*)、经济 (*cadata*)、活动识别 (*mpg*) 等多种重要场景，数据集的范围从大小 100 (*pyrim*) 到超过 20,000 (*cadata*)。

对比方法。我们将本文提出的 SafeW 方法同基线的监督学习算法以及三种先进的半监督回归算法进行对比，

- *k*-NN 算法：基线的监督学习算法，只利用有标注数据训练 *k*-近邻分类器。
- COREG 算法 [159]：一种基于协同训练的半监督回归方法，是半监督回归任务中的代表算法。该算法使用两个采用不同距离度量的 *k*-近邻回归器，每个回归器通过无标注样本的伪标注对标注样本的影响来估计标注置信度，并选

^①<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

择有把握的伪标注数据提供给另一个回归器。

- **Self- k NN 算法**: 基于自训练 (Self-Training) [146] 将监督的 k -近邻算法扩展到半监督回归场景中。首先训练一个仅基于有标注样本的监督 k -近邻模型, 然后利用该模型预测无标注样本的标注, 通过将预测出的标注添加到无标注样本上作为伪标注来不断扩充标注数据集的规模, 并继续训练监督 k -近邻模型, 不断重复此过程, 直到达到最大的迭代轮数, 或者无标注样本上的预测结果不再发生变化。
- **Self-LS 算法**: 基于自训练将监督最小二乘回归算法 [60] 扩展到半监督场景中, 与 Self- k NN 算法类似, 只不过监督算法由 k -近邻算法转变为了最小二乘回归算法。

此外, 我们还与基于投票的集成学习方法 (Voting), 以及基于最优权重的集成学习方法 (OpW, Optimal Weighting) 进行对比。其中, Voting 方法是将多个基学习器的预测结果进行平均, 在多种集成学习任务中都可以取得良好的效果 [156], OpW 方法是利用所有数据的真实标注为基学习器学得最优的权重, 在真实任务中是不可行的, 可以看作是模型集成性能的上限。

参数设置。对于基线的 k -NN 方法, 采用欧氏距离 (Euclidean Distance) 计算样本之间的距离, k 设置为 1; 对于 Self- k NN 方法, 同样采用欧氏距离, k 设置为 3, 最大迭代次数设置为 5, 并且实验发现进一步增加迭代轮数并不会提高性能; 对于 Self-LS 方法, 与标注样本和无标注样本重要性相关的参数分别设置为 1 和 0.1; 对于 COREG 方法, 参数设置为开源代码中推荐的参数, 采用欧氏距离和马氏距离 (Mahalanobis Distance) 作为距离度量; 对于 Voting, OpW 和本文提出的 SafeW 方法, 使用 3 个基学习器的结果做集成, 其中一个来自 Self-LS 方法, 另外两个分别来自采用欧氏距离和余弦距离 (Cosine Distance) 的 Self- k NN 方法。对于本文提出的 SafeW 方法, 参数 δ 在 $[0.5u, 0.7u]$ 范围内通过 5 折交叉验证设置。在我们的实验中, 所有的特征和标注都被归一化到区间 $[0, 1]$ 。对于每个数据集, 随机选择 5 个和 10 个有标注样本作为标注数据集, 其余作为无标注数据。所有实验均重复 30 次, 并报告模型性能的均值与标准差。

实验结果。表 4-3 和表 4-4 汇报了本章提出的 SafeW 方法和对比方法在标注样本分别为 5 个和 10 个时均方损失的均值与标准差, 从实验结果我们可以观察到如下结论:

表 4-3: 标注样本个数为 5 时, 均方损失的均值与标准差。如果方法性能显著优于/弱于基线的 INN 方法, 对应的数字会被标粗/框选。胜出/打平/打败表示相比基线监督学习方法性能提高、不变和退化的次数, 其中表现最优的方法被标粗。

5 个标注样本									
数据集	INN	Self-kNN	Self-LS	COREG	Voting	OpW	SafeW		
abalone	.017 ± .007	.014 ± .003	.013 ± .004	.013 ± .003	.012 ± .003	.005 ± .001	.013 ± .003		
bodyfat	.024 ± .008	.025 ± .009	<u>.054 ± .016</u>	.026 ± .008	<u>.031 ± .011</u>	.018 ± .003	.025 ± .009		
cadata	.090 ± .031	.073 ± .023	.067 ± .022	.069 ± .028	.069 ± .022	.039 ± .014	.070 ± .023		
cpusmall	.027 ± .012	<u>.031 ± .008</u>	<u>.050 ± .021</u>	<u>.031 ± .009</u>	.024 ± .006	.014 ± .003	.028 ± .009		
eumite2001	.052 ± .017	.037 ± .015	.024 ± .012	.037 ± .011	.031 ± .013	.018 ± .005	.032 ± .010		
housing	.042 ± .007	.043 ± .009	<u>.048 ± .012</u>	.041 ± .008	.042 ± .009	.024 ± .002	.041 ± .009		
mg	.071 ± .035	.057 ± .015	.053 ± .011	.054 ± .019	.054 ± .013	.028 ± .009	.053 ± .013		
mpg	.029 ± .012	.030 ± .012	<u>.040 ± .014</u>	.031 ± .012	.031 ± .012	.016 ± .002	.030 ± .012		
pyrim	.032 ± .009	.027 ± .005	<u>.063 ± .012</u>	.029 ± .011	.025 ± .007	.013 ± .002	.025 ± .005		
space_ga	.005 ± .002	.005 ± .003	<u>.030 ± .005</u>	.004 ± .002	<u>.008 ± .002</u>	.001 ± .000	.004 ± .002		
平均性能	.039	.034	.044	.033	.033	.020	.032		
胜出/打平/打败		5/4/1	4/0/6	5/4/1	5/3/2	9/0/0	6/4/0		

表 4-4: 标注样本个数为 10 时, 均方损失的均值与标准差。如果方法性能显著优于/弱于基线的 INN 方法, 对应的数字会被标粗/框选。胜出/打平/打败表示相比基线监督学习方法性能提高、不变和退化的次数, 其中表现最优的方法被标粗。

数据集	10 个标注样本						
	INN	Self-kNN	Self-LS	COREG	Voting	OpW	SafeW
abalone	.020 ± .010	.014 ± .005	.013 ± .004	.012 ± .003	.012 ± .003	.004 ± .001	.013 ± .005
bodyfat	.019 ± .005	.019 ± .007	.041 ± .013	.020 ± .006	.023 ± .009	.010 ± .002	.018 ± .007
cadata	.083 ± .029	.063 ± .012	.056 ± .007	.054 ± .010	.057 ± .009	.033 ± .011	.060 ± .013
cpusmall	.024 ± .012	.027 ± .008	.042 ± .004	.028 ± .008	.020 ± .005	.012 ± .003	.025 ± .008
eumite2001	.044 ± .014	.037 ± .013	.020 ± .006	.031 ± .009	.029 ± .009	.017 ± .002	.029 ± .007
housing	.039 ± .010	.036 ± .009	.036 ± .009	.035 ± .005	.034 ± .008	.021 ± .003	.035 ± .009
mg	.062 ± .019	.046 ± .015	.048 ± .011	.045 ± .015	.043 ± .014	.024 ± .004	.045 ± .014
mpg	.022 ± .007	.020 ± .006	.030 ± .014	.021 ± .007	.021 ± .008	.011 ± .001	.020 ± .006
pyrim	.023 ± .006	.021 ± .005	.052 ± .014	.022 ± .006	.020 ± .007	.009 ± .001	.020 ± .006
space_ga	.004 ± .001	.003 ± .001	.028 ± .002	.003 ± .001	.006 ± .001	.000 ± .000	.003 ± .001
平均性能	.034	.029	.037	.027	.026	.016	.027
胜出/打平/打败		6/3/1	4/1/5	6/3/1	7/1/2	9/0/0	7/3/0

1) Self- k NN 算法通常来说可以提升监督学习模型的性能,但是,在两种情况下出现了性能严重退化的问题。

2) Self-LS 算法效果不明显,一个可能的原因是我们采用的实验数据集中监督的最小二乘回归算法性能不如 k NN 算法的性能。

3) COREG 算法作为半监督回归任务中的代表算法实现了良好的性能,但是在某些情况下也会出现性能退化的问题。

4) 基于投票的集成学习算法 Voting 提高了 Self- k NN 和 Self-LS 的平均性能,但是在 6 种情况下也出现了显著性能退化的问题。

5) 本文提出的方法 SafeW 在多种情况下都显著提升了性能,同时在平均性能上也实现了最优的结果。最重要的是, SafeW 方法并没有出现性能下降的问题,实现了安全性。

6) 即使利用了最优权重的 OpW 方法也不能达到误差为 0,这意味着定理 4-9 中的条件并没有满足,但是我们的方法 SafeW 依然实现了安全的性能,这从侧面证明了 SafeW 方法对理论分析中假设条件变动的稳健性。

总体而言,本文提出的 SafeW 方法显著提高了半监督学习的安全性,同时与最先进的方法相比也获得了相当的性能。

4.5.2 领域自适应学习实验

我们进一步在领域自适应学习进行实验探究本文提出的 SafeW 框架的有效性。与半监督学习类似,领域自适应学习也是假设当前任务中只有少量的标注数据,但不同之处在于,领域自适应学习有大量来自其它数据分布的监督信息可以辅助模型训练。领域自适应学习中也存在不安全的现象,即,利用了更多来自其它分布的数据之后反而导致模型性能下降。

数据介绍。在领域自适应学习的场景中,我们采用两个基准数据集: 20News-grous 和 Landmine 数据集^①进行实验。20Newsgroups 数据集 [81] 包括 19,997 个文档,分别属于 20 个不同的新闻组,参考 [31, 87] 等文中的数据划分方式,我们利用其文档的层次结构生成六个来自不同域的数据集。具体来说,学习任务为对最顶层的新闻类别进行分类,对于每个数据集选择两个顶层类别,一个作为正类,另一个作为负类,然后我们分别选择正类和负类下的一些子类构成一

^①<http://www.cse.ust.hk/TL/>

表 4-5: SafeW 方法和对比方法在 20newsgroup 和 Landmine 数据集中的分类正确率 (均值与标准差)。如果方法性能显著优于/弱于基线方法, 对应的数字会被标粗/框选。胜出/打平/打败表示相比基线监督学习方法性能提高、不变和退化的次数, 其中表现最优的方法被标粗。

20newsgroup 数据集										
数据集	LR	Original	MIDA	TCA	TrAdaBoost	Voting	OpW	SafeW		
Comp vs Rec	.703±.009	.749±.014	.796±.020	.794±.016	.808±.016	.796±.014	.889±.010	.796±.017		
Comp vs Sci	.823±.066	.799±.019	.895±.019	.826±.017	.858±.020	.855±.024	.924±.019	.893±.021		
Comp vs Talk	.842±.069	.802±.018	.823±.016	.843±.011	.825±.014	.823±.017	.893±.015	.845±.016		
Sci vs Talk	.729±.105	.710±.012	.746±.016	.702±.009	.717±.021	.729±.043	.824±.010	.747±.015		
Rec vs Sci	.801±.076	.775±.016	.803±.015	.844±.012	.802±.015	.814±.024	.901±.015	.844±.016		
Rec vs Talk	.828±.045	.828±.012	.857±.011	.858±.013	.842±.011	.857±.012	.913±.012	.858±.011		
平均性能	.787	.777	.820	.811	.808	.807	.891	.831		
胜出/打平/打败		1/2/3	4/1/1	3/2/1	3/2/1	3/2/1	6/0/0	5/1/0		
Landmine 数据集										
Domain-20	.922±.017	.924±.003	.927±.004	.926±.005	.918±.003	.924±.004	.963±.003	.927±.004		
Domain-21	.936±.010	.931±.005	.938±.005	.930±.005	.926±.003	.935±.006	.977±.004	.940±.004		
Domain-22	.959±.005	.956±.004	.951±.007	.965±.002	.910±.003	.960±.004	.994±.002	.965±.002		
Domain-23	.936±.010	.931±.004	.942±.005	.931±.005	.963±.004	.947±.003	.981±.003	.943±.004		
Domain-24	.954±.005	.952±.003	.945±.003	.943±.003	.954±.003	.953±.002	.989±.003	.955±.002		
平均性能	.941	.939	.941	.939	.934	.943	.981	.946		
胜出/打平/打败		0/3/2	2/1/2	1/1/3	1/2/2	1/4/0	5/0/0	3/2/0		

个数据分布。在本文中，我们使用来自 *Comp*、*Rec*、*Sci* 和 *Talk* 四个顶层类别下的文档来构造数据集。

Landmine 数据集是一个地雷检测数据集，包括 29 个域和 9 个特征，从域 1 到域 5 的数据集是从一个绿植分布密集的区域收集的，从域 20 到域 24 是从沙漠环境中收集的。在本文的实验中，我们使用域 1 到域 5 的数据作为源域数据，域 20 到域 24 的数据作为 5 个目标域数据。

对于 *20newsgroup* 数据集，参考论文 [140]，我们随机选取 10% 的样本作为目标域的标注样本，采用 300 维最重要的特征用于训练，对于 *Landmine* 数据集，我们随机选取目标域中 5% 的样本作为标注样本，其余为无标注样本。

对比方法。我们将本文提出的 SafeW 方法与基线算法以及三种先进的领域自适应学习方法进行对比：

- LR (Logistic Regression) 方法：基线监督学习算法，直接在目标域的标注数据中以监督学习的方式训练逻辑回归模型。
- Original 方法：基线领域自适应学习方法，直接利用所有的源数据和目标域数据训练监督学习模型。
- MIDA (Maximum Independence Domain Adaptation) 方法 [141]：MIDA 是基于特征迁移的领域自适应学习算法。该方法首先学习一个在源域以及目标域之间不变的子空间，然后在子空间上训练监督学习模型，使其能够适应数据分布的变化。
- TCA (Transfer Component Analysis) 方法 [110]：TCA 也是一种基于特征迁移的领域自适应学习算法，在很多领域自适应任务中都取得了成功的应用。
- TrAdaBoost 方法 [31]：TrAdaBoost 是一种基于数据迁移的领域自适应学习算法，采用 Boosting 方法 [45] 来选择源域中对目标域的模式训练最有帮助的样本，TrAdaBoost 算法已经在多种应用证明了有效性。
- OpW 方法：与半监督学习实验中相同，利用最优的权重对基学习器进行集成，可以看作集成学习模型性能的上界。

参数设置。对于 MIDA 和 TCA 方法，核函数设置为线性核，子空间维数设置为 30；对于 MIDA、TCA 和 Original 方法，采用逻辑回归模型进行训练；对于 TrAdaBoost，采用支持向量机作为基学习器，迭代次数设置为 20；对于 OpW 方法和本文的 SafeW 方法，采用 MIDA、TCA 和 Original 方法作为基学习器，

SafeW 方法的参数 δ 在 $[0.5u, 0.7u]$ 范围内通过 5 折交叉验证进行选择。所有实验均重复 30 次，并报告分类正确率的均值与标准差。

实验结果。在领域自适应学习场景中的实验结果如表 4-5 所示，从实验结果中我们可以看到，Original、MIDA 和 TCA 方法在很多情况下都会出现性能退化的情况，而本文提出的 SafeW 方法则没有这样的问题，此外，在平均性能上，SafeW 方法取得了最好的结果。这些结果表明，我们的方法最优情况下能与最先进方法保持性能相当，同时可以在最坏情况下保证安全性。

4.5.3 多示例学习实验

多示例学习和半监督学习都属于弱监督学习的一种，区别在于半监督学习处理标注不完全的情况，即，只有一部分数据标注，其余数据无标注，而多示例学习处理标注不具体的情况，即，只知道若干个样本中存在正类样本，但不能具体知道哪一个样本的标注为正类。同样，多示例学习中也存在不安全的情况，即利用了更多的弱监督数据之后，性能不如简单的监督学习模型。

数据介绍。对于多示例学习任务，我们采用 7 种多示例数据集，包括 5 种基准数据集 *Musk1*、*Musk2*、*Elephant*、*Fox*、*Tiger*^①，和 2 种常用分类数据集 *Birds* [14]、*SIVAL* [113]。

对比方法。我们将本文提出的 SafeW 方法与如下先进的多示例学习方法相比较：

- SI-SVM 算法：基线监督学习算法，将每一个包的标记指派给包中的所有样本，然后转化为普通的监督学习进行模型训练。
- miSVM 算法 [1]：miSVM 是一种直推式支持向量机算法，首先将包的标注指派给包中的所有样本，然后训练支持向量机模型对数据集中的每个样本进行分类，然后利用新的标注指派重新训练模型，重复此过程，直到标注指派不再发生变动。
- C-kNN [131]：C-kNN 将 kNN 算法扩展到多示例学习任务中，该算法使用最小豪斯多夫距离（Hausdorff Distance）测量两个包之间的距离，该算法被广泛用于多示例分类任务，并取得了良好的性能表现 [162]。

^①<http://www.uco.es/grupos/kdis/momil/>

- CCE [163]: CCE 算法基于聚类和分类器集成, 首先, 该算法对特征空间进行聚类, 将每个包表示为二进制向量, 其中每个位对应于聚类的一个簇, 然后基于该二进制代码训练用于集成的基学习模型。
- MIBoosting [150]: 这种方法是梯度提升算法在多示例学习任务中的扩展, 将原有梯度提升算法的损失函数转换为基于包分类错误的损失函数。MIBoosting 对每个样本单独分类, 并将包内所有样本标注组合起来以获得包的标注。
- mi-Graph [161]: 该方法通过一个图来表示每个包, 其中样本对应图中的节点, 通过图中的联通分支调整样本权重, 属于较大分支的样本具有较低的权重, 因此当对样本进行平均时, 包中存在的每个类别都被平等地表示。然后, 该方法利用基于图的核方法描述包之间的相似性并训练支持向量机模型。
- Voting 算法: 与其它学习场景一致, 对多个基学习器预测结果进行平均。

参数设置。对于 *Birds* 和 *SIVAL* 数据集, 我们分别采用 *Brown Creeper* 和 *Apple* 作为目标类; 对于 C-kNN 算法, 将 refs 参数设置为 1, citers 参数设置为 5; 对于 SI-SVM 和 mi-SVM 算法, 我们采用 LibSVM^①作为具体实现, 并采用 RBF 核函数; 对于 CCE、MIBoosting 和 miGraph 方法, 我们采用默认的推荐参数; 对于 Voting 方法和本文提出的 SafeW 方法, 我们采用 SI-SVM、mi-SVM、C-kNN 和 mi-Graph 作为基学习器, SafeW 的参数 δ 在 $[0.3u, 0.8u]$ 范围内通过 5 折交叉验证选取。所有实验均重复 10 次, 并汇报分类准确率的均值与标准差。

实验结果。表 4-6 汇报了 SafeW 方法和对比方法在 7 个数据集上的实验结果, 从结果我们可以看出:

- 1) CCE、C-kNN 和 MIBoosting 方法在很多情况下都会出现性能退化的不安全现象, 而本文提出的 SafeW 方法则没有这样的问题;
- 2) miGraph 算法实现了最好的平均性能, 但 SafeW 方法和基线方法相比, 性能退化的次数最少;
- 3) 简单集成学习方法 Voting 综合考虑多个基分类器的性能, 但是仍然无法保障安全性, 出现了 1 次性能不如基线学习方法的情况。

上述结果充分验证了 SafeW 方法在多示例学习场景下的有效性。

^①<https://www.csie.ntu.edu.tw/~cjlin/libsvm/>

表 4-6: SafeW 方法和对比方法在多示例学习数据集上的分类正确率 (均值与标准差)。如果方法性能显著优于/弱于基线方法, 对应的数字会被标粗/框选。胜出/打平/打败表示相比基线监督学习方法性能提高、不变和退化的次数, 其中表现最优的方法被标粗。

数据集	SI-SVM	CCE	miSVM	C-kNN	MIBoosting	miGraph	Voting	SafeW
Musk1	.840 ± .119	.831 ± .027	.869 ± .120	.849 ± .143	.837 ± .120	.889 ± .073	.881 ± .079	.869 ± .101
Musk2	.853 ± .101	.723 ± .019	.838 ± .085	.875 ± .131	.790 ± .088	.903 ± .086	.879 ± .049	.884 ± .082
Fox	.546 ± .092	.599 ± .027	.582 ± .102	.576 ± .016	.638 ± .102	.616 ± .079	.590 ± .034	.590 ± .051
Elephant	.801 ± .088	.793 ± .021	.825 ± .073	.785 ± .016	.827 ± .073	.869 ± .078	.825 ± .049	.819 ± .053
Tiger	.778 ± .092	.758 ± .012	.789 ± .089	.757 ± .017	.784 ± .085	.801 ± .083	.779 ± .017	.790 ± .031
SIVAL	.761 ± .071	.715 ± .053	.771 ± .110	.735 ± .151	.715 ± .064	.756 ± .035	.737 ± .029	.755 ± .047
Birds	.720 ± .121	.690 ± .095	.720 ± .090	.707 ± .090	.643 ± .141	.663 ± .084	.713 ± .081	.713 ± .090
平均性能	.757	.730	.771	.755	.748	.785	.772	.774
胜出/打平/打败		1/2/4	4/3/0	2/2/3	2/2/3	5/1/1	4/2/1	5/2/0

4.5.4 标注噪声学习实验

标注噪声学习处理监督信息不正确的弱监督数据，也是弱监督学习领域的重要研究方向。在标注噪声学习中，也可能存在不安全的现象，即，利用了更多含有标注噪声的训练数据后模型性能反而不如只利用少量的正确标注数据。

表 4-7: 数据集统计信息介绍。

数据集	样本总数	特征维度	标注样本	噪声样本	测试样本
Australian	690	14	166	386	138
Breast-Cancer	683	10	165	382	136
Diabetes	768	8	184	431	153
Digit1	1,500	241	360	840	300
Heart	270	13	65	151	54
Ionosphere	351	34	85	196	70
Splice	3,175	60	762	1,778	635
USPS	9,298	256	2,231	5,208	1,859

数据介绍。对于标注噪声学习，我们采用了大量常用的分类数据集^①进行实验，包括 *Australian*、*Breast-Cancer*、*Diabetes*、*Digit1*、*Heart*、*Ionosphere*、*Splice*、*USPS*，这些数据集的统计特征如表 4-7 所示。对于每个数据集，我们随机划分 80% 的样本作为训练数据，其余样本作为测试数据，在训练数据集中，我们进一步随机选择 70% 的样本并以概率 p 随机反转其对应的标注作为标注噪声的样本， p 的取值范围是 $[0.1, 0.4]$ 。

对比方法。我们将本文提出的 *SafeW* 方法与如下方法进行比较：

- **Sup-SVM 算法：**基线算法，仅利用少量的正确标注样本以监督学习的方式的训练一个支持向量机模型。
- **Bagging 算法：**一种代表性的集成学习算法，在大量实验中证明 *Bagging* 算法对标注噪声具有很强的稳健性 [44]。
- **rLR (Robust Logistic Regression) 算法 [12]：**rLR 算法对标准的逻辑回归模型进行扩展，使其能够处理标注噪声的情况。
- **SVM, LR, k -NN：**经典的监督学习算法，在不考虑标注噪声处理的情况下，利用所有的数据训练模型。

^①<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

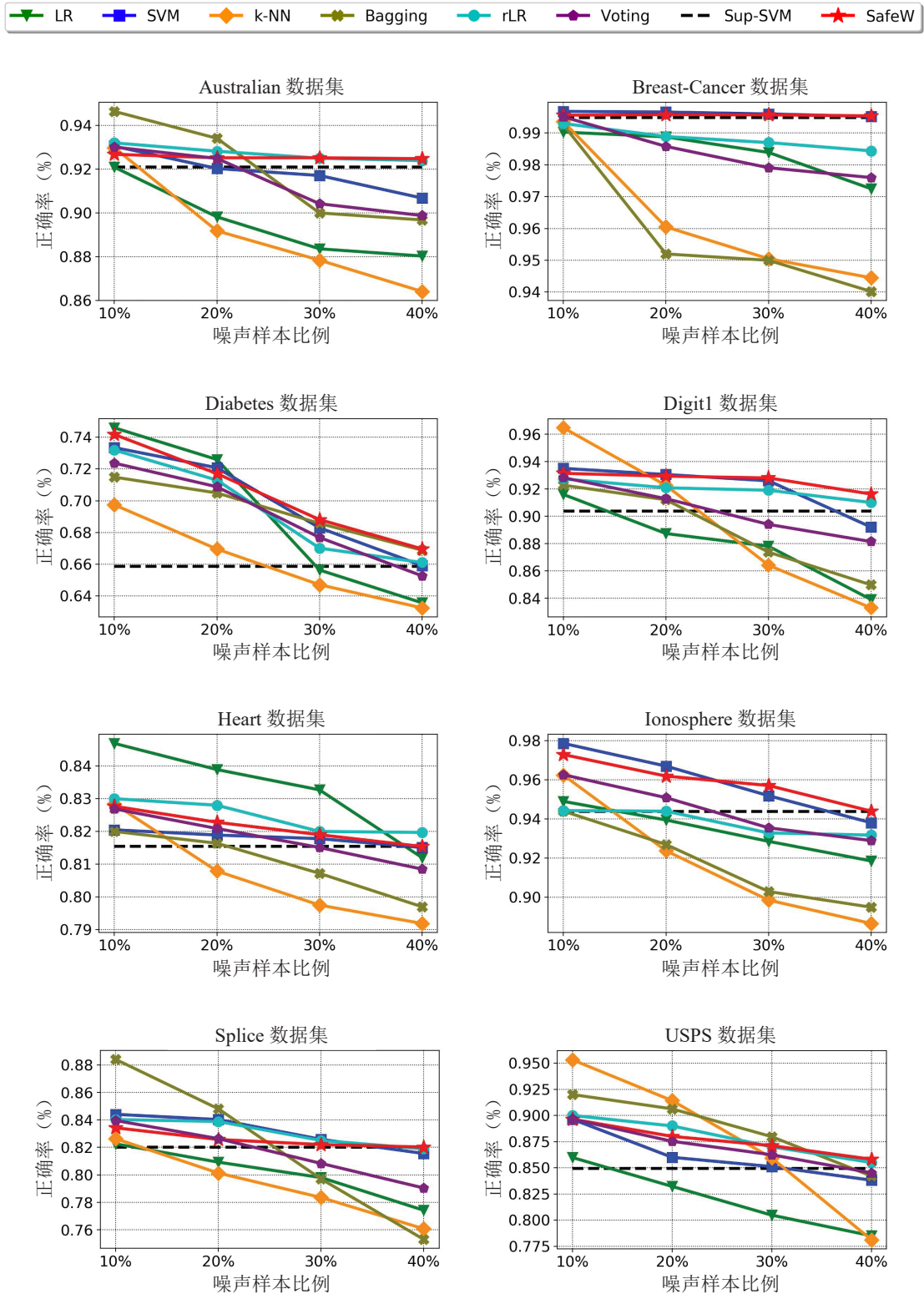


图 4-3: 随数据噪声比例变化, SafeW 和对比方法的性能变化。

表 4-8: SafeW 方法和对比方法在噪声学习数据集中的分类正确率。如果方法性能显著优于/弱于基线方法, 对应的数字会被标粗/框选。胜出/打平/打败表示相比基线监督学习方法性能提高、不变和退化的次数, 其中表现最优的方法被标粗。

数据集	Sup-SVM	LR	SVM	k-NN	Bagging	rLR	SafeW
Australian	.9216	.8958	.9088	.8910	.9193	.9271	.9264
Breast	.9949	.9839	.9961	.9623	.9587	.9883	.9956
Diabetes	.6585	.6910	.6991	.6616	.6936	.6940	.7041
Digit1	.9047	.8800	.9208	.8961	.8895	.9193	.9262
Heart	.8154	.8326	.8181	.8064	.8102	.8244	.8212
Ionosphere	.9435	.9383	.9589	.9179	.9173	.9384	.9590
USPS	.8491	.8204	.8614	.8770	.8870	.8788	.8764
Splice	.8202	.8010	.8314	.7930	.8205	.8307	.8255
平均性能	.8635	.8548	.8743	.8507	.8613	.8745	.8793
胜出/打平/打败		2/1/5	5/2/1	1/2/5	2/3/3	4/3/1	5/3/0

参数设置。对于 LR 算法, 我们采用 Matlab 中的 *glmfit* 函数进行实现; 对于 k -NN 方法, k 的取值设置为 3; 对于 Sup-SVM 和 SVM 方法, 我们采用 Libsvm 包 [21] 中的实现, 核函数设置为 RBF 核; 对于 Bagging 方法, 我们采用决策树作为基学习器; 对于 rLR 算法, 参数设置为原文推荐值; 对于本文提出的 SafeW 方法, 我们采用 LR、SVM 和 k -NN 作为基学习器, 参数 δ 在 $[0.5u, 0.7u]$ 取值范围内通过 5 折交叉验证进行设置。每个实验重复 30 次, 并汇报分类准确率的均值与标准差。

实验结果。图 4-3 汇报了随噪声数据比例的变化, SafeW 方法和对比方法的正确率变化, 表 4-8 汇报了在多种噪声数据比例下, 所有算法的平均正确率。从实验结果中我们可以看出:

1) 随着噪声样本比例的增加, 既有机器学习算法的性能普遍下降, 这表明噪声样本对机器学习性能的损害;

2) 与基线的监督学习算法相比, 在很多情况下对比方法的性能比基线的 Sup-SVM 算法更差, 尤其是当噪声比例增加时, 而本文提出的 SafeW 算法则没有遇到这种问题;

3) 在所有噪声比例下的平均性能, SafeW 算法是最优的。这些结果证明了 SafeW 方法在处理标注噪声学习问题时, 也能实现安全性。

表 4-9: 半监督学习中, 标注样本个数为 10 时, 平均绝对损失的均值与标准差。如果方法性能显著优于/弱于基线的 INN 方法, 对应的数字会被标粗/框选。胜出/打平/打败表示相比基线监督学习方法性能提高、不变和退化的次数, 其中表现最优的方法被标粗。

数据集	平均绝对损失					
	INN	Self-kNN	Self-LS	COREG	Voting	SafeW
abalone	.100 ± .025	.089 ± .020	.086 ± .018	.083 ± .015	.081 ± .018	.086 ± .019
bodyfat	.108 ± .013	.107 ± .018	.164 ± .026	.114 ± .015	.119 ± .023	.105 ± .018
cadata	.216 ± .037	.195 ± .022	.189 ± .016	.182 ± .023	.189 ± .019	.192 ± .023
cpusmall	.073 ± .014	.078 ± .007	.168 ± .010	.081 ± .008	.092 ± .008	.076 ± .007
eumite2001	.162 ± .023	.152 ± .027	.108 ± .016	.138 ± .018	.132 ± .021	.133 ± .017
housing	.137 ± .018	.135 ± .023	.140 ± .023	.135 ± .016	.131 ± .022	.132 ± .022
mg	.188 ± .029	.166 ± .025	.176 ± .017	.168 ± .026	.163 ± .023	.164 ± .025
mpg	.110 ± .014	.107 ± .018	.138 ± .029	.112 ± .020	.109 ± .022	.105 ± .018
pyrim	.105 ± .014	.107 ± .011	.174 ± .021	.111 ± .012	.095 ± .016	.099 ± .014
space_ga	.050 ± .005	.043 ± .005	.131 ± .004	.041 ± .005	.060 ± .004	.042 ± .004
平均性能	.125	.118	.147	.116	.117	.114
胜出/打平/打败		5/4/1	4/1/5	5/2/3	5/2/3	6/4/0

表 4-10: 噪声学习中, 所有噪声比例下交叉熵损失的均值。如果方法性能显著优于/弱于基线的 1NN 方法, 对应的数字会被标粗/框选。胜出/打平/打败表示相比基线监督学习方法性能提高、不变和退化的次数, 其中表现最优的方法被标粗。

交叉熵损失							
数据集	Sup-SVM	LR	SVM	k-NN	Bagging	rLR	SafeW
Australian	.0816	.1100	.0956	.1130	.0823	.0756	.0764
Breast-Cancer	.0051	.0162	.0049	.0372	.0402	.0115	.0048
Diabetes	.4177	.3679	.3596	.3658	.3652	.3614	.3590
Digit1	.1001	.1278	.0825	.0901	.1203	.0841	.0794
Heart	.2048	.1826	.2077	.2152	.2098	.1929	.1980
Ionosphere	.0582	.0606	.0419	.0856	.0863	.0636	.0419
USPS	.1638	.1990	.1492	.1335	.1189	.1292	.1320
Splice	.1982	.2230	.1845	.2320	.1978	.1835	.1917
平均性能	.1537	.1609	.1407	.1591	.1526	.1379	.1354
胜出/打平/打败		2/1/5	5/2/1	2/1/5	2/3/3	5/2/1	5/3/0

4.5.5 评价指标稳健性

在本文中, 对于回归问题和分类问题, 我们分别展示了损失函数为均方损失和铰链损失时 SafeW 的具体优化过程, 研究 SafeW 在其它性能指标下是否有效也是一个很重要的问题。因此, 我们进一步进行实验, 分析在回归任务中评价指标为平均绝对损失 (Mean Absolute Loss) 时 SafeW 的性能, 结果如表 4-9 所示, 以及在分类任务中评价指标为交叉熵损失时 SafeW 的性能, 结果如表 4-10 所示。实验结果表明, 尽管我们没有直接优化这些评价指标, 但 SafeW 方法仍然实现了安全性能, 这一结果表明 SafeW 对评价指标的变化具有一定的稳健性。

4.5.6 运行时间分析

对于回归任务和分类任务的常用损失函数, 本文分析了 SafeW 方法可以将原始最大最小优化问题转化为二次优化或线性优化问题进行解决, 我们在实验中采用了先进的优化求解工具 MOSEK 进行高效求解, 其运行时间开销主要取决于基学习器的数量。我们在表 4-11 中展示了 SafeW 方法在分类和回归任务上的运行时间, 结果表明本文提出的优化方法是快速有效的, 例如, 对于包含

表 4-11: SafeW 在分类任务和回归任务上的运行时间开销。

分类任务			
数据集	样本数目	特征维度	时间开销 (秒)
heart	270	3	0.002
ionosphere	351	34	0.004
breast-cancer	683	10	0.007
australian	690	4	0.007
diabetes	768	8	0.007
splICE	1,000	60	0.009
usps	1,500	241	0.011
回归任务			
bodyfat	252	14	0.007
mpg	392	7	0.009
housing	506	13	0.009
mg	1,385	6	0.009
abalone	4,177	8	0.010
cpusmall	8,192	12	0.010
cadata	2,0640	8	0.014

20,640 个样本的 Cadata 数据集，在 2.7Ghz CPU，4G 内存的普通 PC 上优化时间开销不到 0.02 秒。

4.6 小结

在本章中，我们研究了开放环境下领域知识不充分的安全半监督学习，目标是当领域知识不足以提供可靠的模型选择时，保证在通常情况下半监督学习模型可以实现性能提升，并且在最坏情况下，半监督学习模型也不会出现性能退化的问题。领域知识不充分的安全半监督学习是半监督学习领域至关重要的问题，但是现有的研究工作还非常少。我们在本章中提出了一种基于集成学习的安全半监督学习框架 SafeW，具体而言，SafeW 对多个半监督学习器的预测结果进行加权集成，并且我们考虑最坏情况下的集成结果，优化最坏情况下模型相比于基线监督学习模型的性能增益，在数学上得到了一个最大最小的优化目标。我们从理论上对该优化目标进行了分析，证明了当真实标注可以由多个基

学习器的预测结果加权组合得到时可以保证半监督学习的安全性，该假设条件比以往研究工作的条件要更容易实现。此外，该框架可以灵活的嵌入关于基学习器性能的先验知识。对于分类任务和回归任务常用的损失函数，如铰链损失、均方损失、交叉熵损失等，我们给出了具体的优化方法，可以将原始非凸的最大最小优化目标转化为容易解决的凸优化问题，如线性优化问题或二次规划问题，从而高效地获得全局最优解。值得一提的是，本章提出的 SafeW 框架具有很强的通用性，不仅能够处理安全的半监督学习问题，在其它弱监督学习场景中，如领域自适应学习、标注噪声学习、多示例学习，也可以取得安全的性能表现。大量的实验结果证明了 SafeW 的通用性和有效性。

在未来工作中，我们一方面会研究如何通过数据增强 (Data Augmentation) 提升模型的确定性，从而进行更可靠的模型选择，另一方面我们会研究更多开放场景下半监督学习模型的安全性，比如数据分布长尾的半监督学习、存在对抗样本的半监督学习等，从而促进半监督学习在更多现实领域的应用。

本章的主要工作已经成文发表，包括：

- Lan-Zhe Guo, Yu-Feng Li. A General Formulation for Safely Exploiting Weakly Supervised Data. In: **Proceedings of the 32nd AAAI conference on Artificial Intelligence (AAAI'18)**, New Orleans, LA, pp.3126-3133, 2018. (中国计算机学会 A 类会议，第一作者)
- Yu-Feng Li, Lan-Zhe Guo, Zhi-Hua Zhou. Towards Safe Weakly Supervised Learning. **IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)**, 43(1):334-346, 2021. (中国计算机学会 A 类期刊，导师外一作)

第五章 适于类别比例失衡的 稳健半监督学习

5.1 引言

封闭环境下的半监督学习假设数据分布中的类别比例是平衡的，即每个类别含有相似的样本数量。然而，在现实的开放环境中，该假设往往很难成立，例如，在物种分类任务中，“狗”、“猫”等物种天然就比“大熊猫”、“金丝猴”等物种更为常见；在金融欺诈检测任务中，“非欺诈”类别的样本要远远多于“欺诈”类别的样本数量。传统的半监督学习方法在类别比例不平衡的数据中会出现性能不稳健的现象，即只在样本数量较多的类别上表现良好，而在样本数量较少的类别上性能退化。该现象严重阻碍了半监督学习在真实开放环境中的应用，因此，亟需研究对类别比例不平衡稳健的半监督学习方法，当类别比例失衡时能够在所有类别上实现稳健的性能提升。

本章结合工业界应用网约车平台中面临的现实问题展开类别比例不平衡的半监督学习方法的研究。首先，在网约车平台中，每一笔网约车订单结束后，平台会给用户弹出评价问题，以收集乘客对司机的评价（好评或者差评），网约车智能评价系统尝试构建机器学习模型来预测用户对司机的评价结果。然而现实场景中，有大量用户在订单结束后并不会花时间进行评价，因此大量数据是无标注的，同时好评的订单数量要远远多于差评订单数量。网约车智能评价任务可以形式化为一个类别不平衡的二分类半监督学习问题，传统的未考虑类别不平衡的半监督学习方法不能很好的解决该问题。

为了解决上述问题，本文提出了一种基于 AUC (Area Under the ROC Curve) 优化的半监督学习方法 CWSL。相比于常用的正确率 (Accuracy) 指标，AUC 在类别不平衡问题中更加适用 [168]，例如，给定 1000 个样本，包括 998 个正类样本，但是负类样本只有 2 个，此时，只要机器学习算法返回一个将所有样本都预测为正类的模型，就能达到 99.8% 的正确率，但是显然这样的模型是没有

任何价值的，而 AUC 指标通过考虑模型预测结果的排序质量可以有效避免该问题。具体而言，我们为所有的无标注样本全部打上好评的伪标注，然后采用监督学习的方式进行模型训练，并且在训练过程中直接优化模型的 AUC 性能。考虑到虽然好评标注的样本远远多余差评标注，但是无标注样本中仍然会存在部分差评样本，我们进一步提出对样本进行赋权，降低噪声伪标注样本对模型性能的负面影响，并且在验证集性能的指导下优化样本权重，得到使 AUC 性能尽可能高的模型。我们在第二章提出的基于双层优化的样本赋权机制的基础上进一步引入高效的 AUC 优化方法，相比传统的半监督学习方法取得了显著的性能提升，并且在滴滴平台智能判责系统中成功应用。

其次，在网约车任务中，往往会出现司乘纠纷问题，当出现纠纷时网约车平台需要给出相应的责任判定，网约车智能判责系统尝试构建机器学习模型来预测判责结果并给出判责的原因。由于现实中可能遇到的判责原因非常多样，比如司机绕路、乘客不守时等，需要请专家人工进行标注，代价非常高，因此只能获得少量的标注数据，大量数据是无标注的。同时，在现实场景中一次纠纷可能同时与多个判责原因相关联，某些判责原因出现的频率也要远高于其它原因，所以网约车判责任务可以形式化为一个类别不平衡的半监督多标记学习问题。研究该问题的工作目前相对较少，所以我们首先在图 5-1 中给出了一个该问题的形式化例子以便于理解。对于该问题不仅需要考虑不平衡问题中常用的 AUC 指标，也要考虑多标记学习中常用的 F1-Score、Hamming Loss 等指标，以往的半监督学习方法并不能很好的解决该问题。

为了解决上述类别不平衡的半监督多标记学习问题，我们提出了一种新的学习框架 LIML (Learning from Imbalanced and Incomplete Supervision)。具体而言，LIML 包括三个核心模块：标注分离，关系挖掘和标注补全。在标注分离模块，我们采用两个模型，一个处理多数类标注，另一个模型处理少数类标注，通过标注分离的方式，可以有效防止模型在少数类标注上性能下降的问题。在获取两个模型的预测结果之后，我们在关系挖掘模块挖掘各个标注之间的相关性，标注相关性是多标记学习中非常核心的问题。最后我们设计无监督损失项对无标注样本进行标注补全，有效利用无标注数据。在网约车判责任务、图像分类任务、文本识别任务、基因检测任务等现实应用中的大量实验结果验证了 LIML 方法在类别平衡半监督多标记学习问题中，可以在多种评价指标，如 Hamming

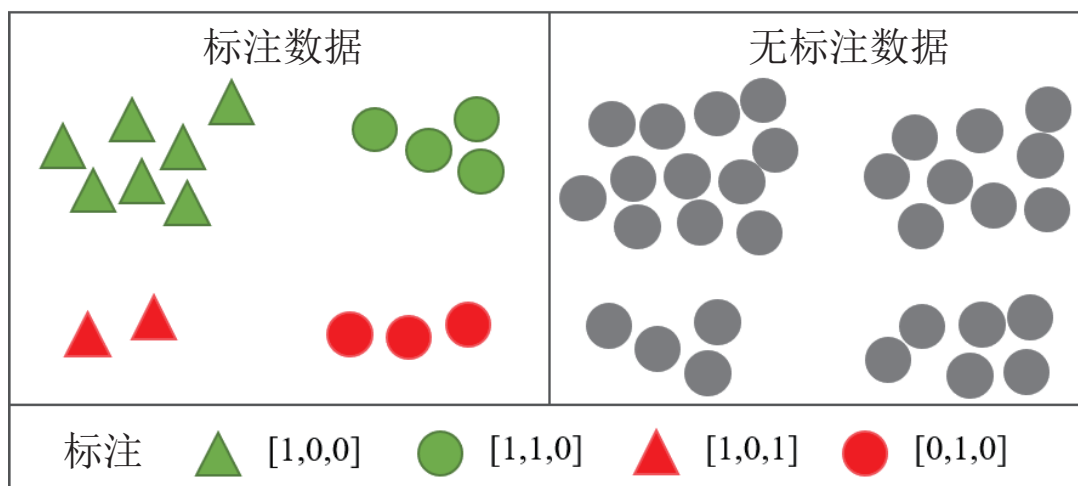


图 5-1: 类别不平衡的半监督多标记学习问题。左侧代表标注数据集，右侧代表无标注数据集。颜色（红色，绿色）和形状（三角形，圆形）的组合表示不同的标记组合，每个样本最多包含三个标记。在该例子中，可以看到三个标记出现的数量分别为 13，7，和 2，因此这是一个类别不平衡的半监督多标记学习问题。

Loss、Ranking Loss、One Error、Coverage、Average Precision、Macro AUC、Micro AUC、Macro F1、Micro F1 上实现稳健的学习性能。

5.2 相关工作

在本节，我们介绍与本章相关的研究工作，包括类别不平衡学习、类别不平衡多标记学习、半监督多标记学习。

5.2.1 类别不平衡学习

类别不平衡是指分类任务中不同类别的样本数量差别很大的情况 [86]，是现实任务中常见的现象。目前已经有大量针对类别不平衡的机器学习算法被提出 [15, 71]，其中最常用的一类算法是根据每个类别的样本数量对训练目标进行重新平衡 (re-balancing)，代表的方法包括：1) 重赋权法 (re-weighting)，对训练样本进行加权，通过赋予少数类样本较高的权重对训练目标进行平衡 [17, 30, 67, 72, 73, 94, 116, 66]；2) 重采样法 (re-sampling)，通过对训练数据进行采样，如对少数类样本进行过采样或者对多数类样本进行欠采样，从而构造类别平衡的训练数据 [26, 61, 16]。然而，这些方法均假设测试数据是类别平衡的，依然采用正确率作为评价指标，而在开放世界的真实任务，如网约车智能评价任务中，

测试环境也是极度不平衡的，需要考虑 AUC 指标，以上方法不再适用。

5.2.2 类别不平衡多标记学习

多标记学习问题中一个样本可能同时包含多个标记，由于标记的共现性，即多数类标记和少数类标记可能在一个训练样本上同时出现，不能直接采用单标记类别不平衡学习中的重采样或重赋权方法。有一些研究工作尝试将重采样方法扩展到多标记学习场景中，例如，MLRUS 算法 [24] 通过随机省略带有多数类标记的训练样本，缓解每个标记上正负类不平衡的问题；MLRUS 的孪生算法 MLROS [24] 通过复制与少数类标记相关的样本，增加少数类标记上的样本数量；MLeNN 算法 [23] 基于编辑最近邻规则 (Edited Nearest Neighbor) 进行欠采样，将容易被错误分类的多数类样本进行删除，从而缓解类别不平衡问题；MLSMOTE 算法 [25] 将经典的过采样方法 SMOTE (Synthetic Minority Oversampling Technique) 算法扩展至多标记场景中，通过随机选择包含少数类标记的样本及其邻近样本生成新的训练样本以缓解不平衡问题。也有一些工作尝试将重赋权法扩展至多标记学习场景，在考虑标记共现性的条件下引入不同的权重处理不同类别的训练样本。例如，COCOA 算法 [152] 将多标记数据按照每个标记转换为多个单标记的训练数据集，并针对每个数据集加权训练分类器；SOSHF 算法 [33] 利用代价敏感的聚类方法将多标记数据转换为类别不平衡的单标记数据，并通过 Helliger 决策树解决该学习任务。最近，[135] 提出一种改进的能够对数据分布进行平衡的二分类交叉熵损失，该损失在标准交叉熵损失的基础上考虑了多标记学习问题中的标记共现性，通过优化该损失得到的模型可以在多种类别不平衡的多标记数据集上实现最先进的性能。然而，以上方法需要大量的标注数据来估计类别不平衡的比例，无法直接应用到标注数据不足的半监督学习场景中。

5.2.3 半监督多标记学习

半监督学习在标注信息不足时通过利用大量易于获取的无标注数据提升模型性能。[97] 研究半监督场景下的多标记学习问题，基于平滑假设，即相似的样本应该具有相似的预测标注，将半监督多标记学习问题形式化为带约束的非

负矩阵分解问题。[134] 在挖掘标注样本中的标记相关性的同时最大化模型在无标注样本上预测结果的间隔，从而得到基于大间隔假设的半监督多标记学习器。[127] 提出 SMILE 算法，利用标注样本和无标注样本构建图结构，并基于图半监督学习算法训练半监督多标记学习模型。[149, 138] 将经典的协同训练 (Co-Training) 算法 [10] 扩展到多标记数据，通过最大化多样性优化两个不同的特征视图，然后在每个视图上训练分类器，并为另一个分类器提供伪标注样本，在多标记学习任务中伪标记的形式是样本与每个标记相关性的排序。[132] 提出 DRML 算法，该方法采用域适应 (Domain Adaptation) 技术探索特征分布和标记相关性，并为无标注样本生成伪标注。然而，以上半监督多标记学习算法忽略了现实场景中多标记数据天然存在的类别不平衡问题，导致模型出现在少数类上性能退化的现象。

综上所述，现有半监督学习方法不能很好的解决开放环境下类别不平衡的学习问题，亟需建立对类别不平衡稳健的单标记、多标记半监督学习方法。

5.3 本文工作

在本节，我们介绍本文的工作，包括问题设定、面向类别不平衡的半监督单标记学习方法、面向类别不平衡的半监督多标记学习方法，以及现实工业界任务，网约车智能评价系统和网约车智能判责系统中的应用。

5.3.1 问题设定

令 $\mathcal{X} = \mathbb{R}^d$ 表示 d 维特征空间， $\mathcal{Y} = \{0, 1\}^k$ 表示 k 维标记空间。给定包含 n 个样本的标注数据集 $\mathcal{D}_l = \{\mathbf{X}_l, \mathbf{Y}_l\} = \{(\mathbf{x}_i, \mathbf{y}_i) | 1 \leq i \leq n\}$ ，其中 $\mathbf{X}_l \in \mathbb{R}^{n \times d}$ 表示数据的特征矩阵， $\mathbf{Y}_l \in \{0, 1\}^{n \times k}$ 表示对应的标记矩阵， $\mathbf{x}_i \in \mathcal{X}$ ， $\mathbf{y}_i \in \mathcal{Y}$ ， $\mathbf{y}_i^j = 1$ 表示第 i 个样本与第 j 个标记相关，反之则不相关。对于单标记学习任务，每个样本只与一个标记相关联，而在多标记学习任务中，一个样本可以同时与多个标记相关联。此外，我们还可以获得包含 m 个样本的无标注数据集 $\mathcal{D}_u = \{\mathbf{X}_u\} = \{\mathbf{x}_i | 1 \leq i \leq m\}$ ，其中 $\mathbf{X}_u \in \mathbb{R}^{m \times d}$ 。通常来说， $n \ll m$ 。学习任务的目标是构建预测模型 $f(\mathbf{x}; \theta) : \mathcal{X} \rightarrow \mathcal{Y}$ ，其中 θ 代表模型参数，对于未见的测试样本 $\mathbf{x} \in \mathcal{X}$ ， f 能够预测出其准确的标记。

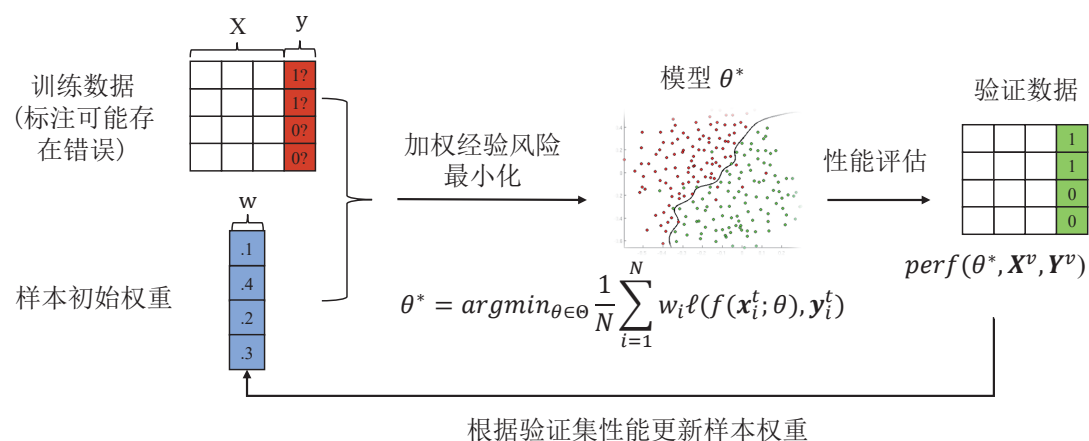


图 5-2: 类别不平衡的半监督单标记学习算法基本框架。

5.3.2 类别不平衡的半监督单标记学习

在网约车智能评价系统中，需要处理的任务是类别不平衡的半监督单标记学习问题，其主要难点是如何在有效利用无标注数据的同时，在类别不平衡场景下实现稳健的学习性能。我们提出 CWSL 算法来解决该问题，具体而言，我们考虑如下做法：首先，由于网约车智能评价数据中，好评样本远远多于差评样本，我们将所有的无标注数据全部视作好评样本构造得到包含 $N = n + m$ 个样本的的训练集 $\{(\mathbf{x}_i^t, y_i^t), 1 \leq i \leq N\}$ 。但是由于其中差评样本的存在导致标注存在噪声，我们进一步在传统经验风险最小化框架的基础上引入样本赋权的机制，降低噪声样本对学习性能的影响。模型的优化目标如下所示：

$$\theta^* = \operatorname{arg min}_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N w_i \ell_i(f(\mathbf{x}_i^t; \theta), y_i^t) \quad (5-1)$$

其中 w_i 表示第 i 个样本对应的权重， $\ell(\cdot)$ 表示采用的损失函数， Θ 表示模型的参数空间。

优化上述目标的难点是样本权重我们是无法得知的。传统的标注噪声学习方法通常根据模型的损失函数估计样本权重，其基本思想是损失较大的样本对应的权重应该较低 [44]。但是该思想在类别不平衡数据分布中不再适用，因为在类别不平衡任务中少数类的样本往往也会呈现较大的损失。在本文中，我们基于第二章所提出的基于双层优化的样本赋权机制，进一步引入包含 M 个样本的验证数据集 $\{\mathbf{X}^v, \mathbf{Y}^v\} = \{(\mathbf{x}_i^v, y_i^v), 1 \leq i \leq M\}$ ，提出根据模型在验证集上的性能

指标进行样本权重学习的算法框架，如下所示：

$$\mathbf{w}^* = \arg \max_{\mathbf{w} \in \mathcal{W}} \text{perf}(\theta^*, \mathbf{X}^v, \mathbf{Y}^v) \quad (5-2)$$

其中 $\text{perf}(\theta^*, \mathbf{X}^v, \mathbf{Y}^v)$ 表示当参数为 θ^* 时模型在验证数据上的性能， $\text{perf}(\cdot)$ 代表某种评价指标，如 AUC、F1-Score 等。

综合公式 5-1 和公式 5-2，可以得到整体的双层优化目标如下所示：

$$\begin{aligned} \max_{\mathbf{w} \in \mathcal{W}} \quad & \text{perf}(\theta^*, \mathbf{X}^v, \mathbf{Y}^v) \\ \text{s.t.} \quad & \theta^* = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N w_i \ell(f(\mathbf{x}_i^t; \theta), y_i^t) \end{aligned} \quad (5-3)$$

其中权重 \mathbf{w} 的取值范围 $\mathcal{W} = \{\mathbf{w} | 0 \leq \mathbf{w} \leq 1\}$ 。图 5-2 展示了该方法的整体框架。

为了应对网约车评价数据中类别不平衡的问题，我们采用 AUC 作为评价指标，因为 AUC 考虑预测结果排序的质量，与正负类的类别比例无关 [96]。AUC 指标被广泛应用在存在类别比例不平衡问题的现实任务中，比如推荐系统、欺诈检测、广告点击预测任务等 [100]。

具体来说，针对二分类问题，令 \mathbf{x}^+ 、 \mathbf{x}^- 表示从正类分布 \mathcal{P}^+ 和负类分布 \mathcal{P}^- 采样得到的正类样本和负类样本，AUC 指标的定义如下所示：

$$\text{AUC} = \mathbb{E}_{\substack{\mathbf{x}^+ \sim \mathcal{P}^+ \\ \mathbf{x}^- \sim \mathcal{P}^-}} [\mathbb{I}\{f(\mathbf{x}^+; \theta) - f(\mathbf{x}^-; \theta) > 0\}] \quad (5-4)$$

其中 $\mathbb{I}(\cdot)$ 为指示函数，如果条件成立返回 1，反之则返回 0。上式表示预测结果中正类样本排在负类样本前的概率。

因为 AUC 指标的非凸性和非连续性，直接优化 AUC 是 NP 难的问题 [48]。为了降低 AUC 优化的计算复杂度，我们采用凸替代损失函数作为优化目标。具体而言，将公式 5-4 中的 $\mathbb{I}\{f(\mathbf{x}^+; \theta) - f(\mathbf{x}^-; \theta) > 0\}$ 替换为替代损失 $\phi(f(\mathbf{x}^+; \theta) - f(\mathbf{x}^-; \theta))$ ，然后优化目标转换为：

$$\mathbb{E}_{\substack{\mathbf{x}^+ \sim \mathcal{P}^+ \\ \mathbf{x}^- \sim \mathcal{P}^-}} [\phi(f(\mathbf{x}^+; \theta) - f(\mathbf{x}^-; \theta))] \quad (5-5)$$

在 CWSL 方法中，我们采用配对均方误差 $\phi(t) = (1 - t)^2$ 作为替代损失函数，

因为该函数是可微凸函数并且已经有研究证明该函数与原始 AUC 函数具有一致性 [49]，即优化该损失函数等价于优化原始 AUC 函数。具体而言，将验证数据中正类样本和负类样本的集合记作 $\mathcal{S}^+ = \{\mathbf{x}_1^{v+}, \dots, \mathbf{x}_{m_+}^{v+}\}$ 和 $\mathcal{S}^- = \{\mathbf{x}_1^{v-}, \dots, \mathbf{x}_{m_-}^{v-}\}$ ，其中 m_+ 和 m_- 表示正类样本和负类样本的个数，令 $f(\mathbf{x}; \theta) = \theta^\top \mathbf{x}$ ，则对应的替代损失函数可以写作：

$$\begin{aligned}
\mathcal{L}(\theta) &= \frac{1}{m^+m^-} \sum_{i=1}^{m^+} \sum_{j=1}^{m^-} [(1 - (f(\mathbf{x}_i^{v+}; \theta) - f(\mathbf{x}_j^{v-}; \theta)))^2] \quad (5-6) \\
&= \frac{1}{m^+m^-} \sum_{i=1}^{m^+} \sum_{j=1}^{m^-} [(1 - (\theta^\top \mathbf{x}_i^{v+} - \theta^\top \mathbf{x}_j^{v-}))^2] \\
&= 1 - 2\theta^\top \left[\frac{1}{m^+m^-} \sum_{i=1}^{m^+} \sum_{j=1}^{m^-} (\mathbf{x}_i^{v+} - \mathbf{x}_j^{v-}) \right] \\
&\quad + \theta^\top \left[\frac{1}{m^+m^-} \sum_{i=1}^{m^+} \sum_{j=1}^{m^-} (\mathbf{x}_i^{v+} - \mathbf{x}_j^{v-})(\mathbf{x}_i^{v+} - \mathbf{x}_j^{v-})^\top \right] \theta \\
&= 1 - 2\theta^\top \mu_m + \theta^\top \Sigma_m \theta
\end{aligned}$$

其中

$$\mu_m = \frac{1}{m^+m^-} \sum_{i=1}^{m^+} \sum_{j=1}^{m^-} (\mathbf{x}_i^{v+} - \mathbf{x}_j^{v-}) \quad (5-7)$$

$$\Sigma_m = \frac{1}{m^+m^-} \sum_{i=1}^{m^+} \sum_{j=1}^{m^-} (\mathbf{x}_i^{v+} - \mathbf{x}_j^{v-})(\mathbf{x}_i^{v+} - \mathbf{x}_j^{v-})^\top \quad (5-8)$$

将公式 5-3 和公式 5-6 结合，以 AUC 作为评价指标时 CWSL 的优化目标可以写作如下形式：

$$\begin{aligned}
\min_{\theta \in \Omega} \quad & \mathcal{L}(\theta^*) \quad (5-9) \\
\text{s.t.} \quad & \theta^* = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N w_i \ell(f(\mathbf{x}_i^t; \theta), y_i^t)
\end{aligned}$$

其中 $\mathcal{L}(\theta^*) = -2\theta^{*\top} \mu_m + \theta^{*\top} \Sigma_m \theta^*$ 。在该任务中，我们采用负对数似然函数作为损失函数，即，

$$\ell(f(\mathbf{x}_i^t; \theta), y_i^t) = -y_i^t \theta^\top \mathbf{x}_i^t + \log(1 + e^{\theta^\top \mathbf{x}_i^t}) \quad (5-10)$$

算法 5.1 半监督单标记学习算法 CWSL 的基本流程。

输入： 训练数据 $\{\mathbf{x}_i^t, y_i^t\}_{i=1}^N$ ，验证数据 $\{\mathbf{x}_i^v, y_i^v\}_{i=1}^M$ ，权重向量的初始值 \mathbf{w}_0 ，模型参数的初始值 θ_0 ，优化迭代轮数 T 。

输出： 学得的权重向量 \mathbf{w}_T 和模型参数 θ_T

- 1: **for** $t = 1$ to T **do**
- 2: 更新模型权重: $\theta_t = \theta_{t-1} - \lambda_\theta \mathbf{w}_{t-1} \left. \frac{\partial \ell(\theta)}{\partial \theta} \right|_{\theta=\theta_{t-1}}$
- 3: 利用公式 5-13 和公式 5-14 计算梯度 $\frac{\partial g}{\partial \theta}$ 和 $\frac{\partial g}{\partial \mathbf{w}}$
- 4: 利用公式 5-12 计算雅可比矩阵 J
- 5: 利用公式 5-17 计算梯度 $\frac{\partial \mathcal{L}(\theta_t)}{\partial \mathbf{w}}$
- 6: 更新权重: $\mathbf{w}_t = \mathbf{w}_{t-1} - \lambda_w \left. \frac{\partial \mathcal{L}(\theta_t)}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}_{t-1}}$
- 7: 将权重 \mathbf{w}_t 映射到其取值空间 \mathcal{W}
- 8: **end for**
- 9: 返回 \mathbf{w}_T, θ_T 。

上式的外层优化问题可以用其 KKT 条件 (Karush-Kuhn-Tucker Condition) 等价替换 [13]，如下所示：

$$g(\mathbf{w}, \theta) \equiv \frac{1}{N} \sum_{i=1}^N w_i \frac{\partial \ell(f(\mathbf{x}_i^t; \theta), y_i^t)}{\partial \theta} = 0 \quad (5-11)$$

根据隐函数定理 [106]，我们可以得到如下的雅可比矩阵 (Jacobian Matrix)：

$$J = \frac{\partial \theta}{\partial \mathbf{w}} = - \begin{bmatrix} \frac{\partial g_1}{\partial \theta_1} & \dots & \frac{\partial g_1}{\partial \theta_d} \\ \vdots & & \vdots \\ \frac{\partial g_d}{\partial \theta_1} & \dots & \frac{\partial g_d}{\partial \theta_d} \end{bmatrix}^{-1} \begin{bmatrix} \frac{\partial g_1}{\partial w_1} & \dots & \frac{\partial g_1}{\partial w_N} \\ \vdots & & \vdots \\ \frac{\partial g_d}{\partial w_1} & \dots & \frac{\partial g_d}{\partial w_N} \end{bmatrix} \quad (5-12)$$

根据公式 5-11，梯度 $\frac{\partial g}{\partial \theta}$ 和 $\frac{\partial g}{\partial \mathbf{w}}$ 可以按照下列公式进行计算：

$$\frac{\partial g}{\partial \theta} = \mathbf{w}^\top \left(\frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^\top} \right) \quad (5-13)$$

$$\frac{\partial g}{\partial \mathbf{w}} = \frac{\partial \ell(\theta)}{\partial \theta} \quad (5-14)$$

上述两式中损失函数的一阶梯度和二阶 Hessian 矩阵如下所示：

$$\frac{\partial \ell(\theta)}{\partial \theta} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^t (-y_i^t + p(\mathbf{x}_i^t; \theta)) \quad (5-15)$$

$$\frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^\top} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^t \mathbf{x}_i^{t\top} p(\mathbf{x}_i^t; \theta) (1 - p(\mathbf{x}_i^t; \theta)) \quad (5-16)$$

其中 $p(\mathbf{x}_i^t; \theta)$ 当模型参数为 θ 时样本 \mathbf{x}_i^t 预测为正类的概率。

雅可比矩阵 J 反映了当权重 \mathbf{w} 变动时，模型参数 θ 如何变动的。现在，我们可以进一步采用链式法则计算整个双层优化的目标函数对样本权重 \mathbf{w} 的梯度，如下所示：

$$\frac{\partial \mathcal{L}(\theta^*)}{\partial \mathbf{w}} = \frac{\partial \mathcal{L}(\theta^*)}{\partial \theta} \frac{\partial \theta}{\partial \mathbf{w}} \quad (5-17)$$

综上，整个双层优化目标式可以按照如下的梯度下降方法进行更新：

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \lambda_{\mathbf{w}} \left. \frac{\partial \mathcal{L}(\theta^*)}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}_{t-1}} \quad (5-18)$$

其中 $\lambda_{\mathbf{w}}$ 表示梯度更新的步长。

此外，值得一提的是，上述更新过程中需要计算 Hessian 矩阵的逆矩阵，而在现实任务中逆矩阵的计算往往是开销比较大的，因此，我们进一步提出一种替代算法减少计算复杂度。具体而言，我们通过计算下列线性优化问题：

$$\left(\frac{\partial g}{\partial \theta} \right)_{\mathbf{x}} = - \frac{\partial g}{\partial \mathbf{w}} \quad (5-19)$$

来得到 $-(\frac{\partial g}{\partial \theta})^{-1} \frac{\partial g}{\partial \mathbf{w}}$ 。该线性优化问题只需要计算矩阵向量乘法，因此相比计算 Hessian 矩阵的逆矩阵更为高效。

CWSL 方法的基本流程如算法 5.1 所示。

5.3.3 类别不平衡的半监督多标记学习

在网约车智能判责任务中，需要处理的是类别不平衡的半监督多标记学习问题，现有的机器学习方法在这种问题上不能取得很好的效果。针对该问题，在本节我们提出了一个系统性的学习框架 LIM1，在多种常用的多标记学习评价指标中取得了稳健的性能表现。

类别不平衡的半监督多标记学习问题的三个主要挑战是：1) 少数类标记没有足够的样本，如何防止模型在少数类标记上性能退化？2) 多标记学习中，标记相关性对算法性能至关重要，如何充分挖掘标记相关性？3) 标注样本不足时，

如何有效利用无标注样本帮助模型提升性能？LIMI 针对以上难点，提供了一个系统的解决方案，包括三个主要模块：标注分离、相关性挖掘和标注补全，下面我们将具体介绍各模块的主要细节。

标注分离。对于类别不平衡的多标记数据，直接在所有数据上训练模型并不是最佳选择，因为模型会偏向多数类标记，在少数类标记上性能不佳。对整个数据集进行重采样或重加权可以在一定程度上缓解性能下降，但是这些方法需要大量的标注样本来估计类别不平衡的比例，因此在标注不足的场景中并不适用。为了缓解多数类标记和少数类标记之间的性能差异，我们设计了一种新颖的二分类器策略，使用不同的模型训练技术分别处理多数类标记和少数类标记。具体而言，对于多数类标记，我们可以通过最小化标准的 BCE (Binary Cross-Entropy) 损失直接训练一个神经网络 $C_h(\cdot)$ 作为多数类模型，因为直接使用原始数据训练分类器自然可以得到在多数类标记上表现良好的模型，其优化目标可以写为如下形式：

$$\min_{C_h} \ell_h(C_h(\mathbf{X}_l), \mathbf{Y}_l) = \frac{1}{n} \sum_{i=1}^n BCE(C_h(\mathbf{x}_i), \mathbf{y}_i) \quad (5-20)$$

BCE 损失的具体形式为：

$$BCE(C_h(\mathbf{x}_i), \mathbf{y}_i) = \frac{1}{k} \sum_{j=1}^k [y_i^j \log(C_h(\mathbf{x}_i)_j) + (1 - y_i^j) \log(1 - C_h(\mathbf{x}_i)_j)] \quad (5-21)$$

其中 $C_h(\mathbf{x}_i) \in [0, 1]^k$ 表示模型 $C_h(\cdot)$ 对样本 \mathbf{x} 的预测概率， $C_h(\mathbf{x}_i)_j$ 表示预测结果中样本 \mathbf{x}_i 是否具有第 j 个标记， y_i^j 为样本 \mathbf{x}_i 在标记 j 上的真实标注。

对于少数类模型，我们通过对损失函数进行重加权以提高模型在少数类标记上的性能，因为相比于同时所有标记上取得良好的性能，通过重加权的方法仅在少数类标记上实现良好的性能要容易的多。在不考虑标记共现性的情况下，对于每个具有标记 j 的样本 \mathbf{x}_i (即， $y_i^j = 1$)，标记级别的采样频率为：

$$P_j^C(x_i) = \frac{1}{k} \frac{1}{n_j} \quad (5-22)$$

其中 n_j 表示与第 j 个标记相关的训练样本的个数。

然而，在多标记学习问题中，忽略标记共现性简单的基于标记频率进行重加权并不能取得很好的性能，因为一个样本通常与多个标记相关联，不能独立

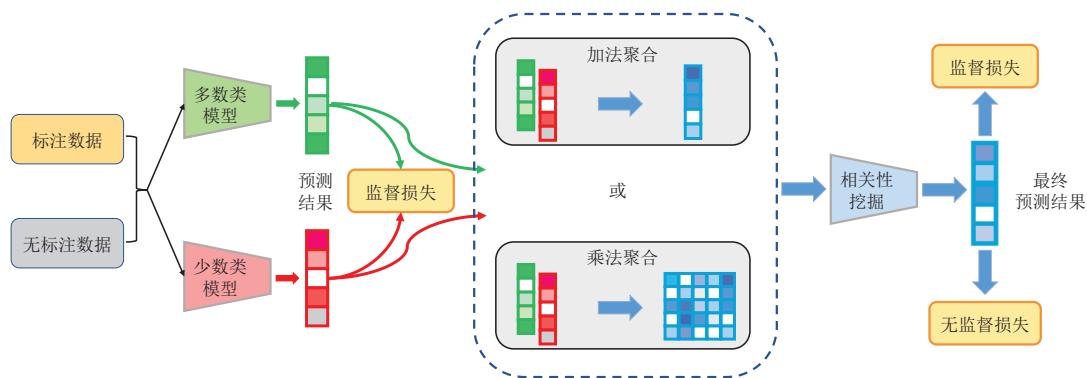


图 5-3: 类别不平衡的半监督多标记学习框架 LIMM 示意图。

的考虑每一个标记。因此，除了标记级别的采样频率外，我们还考虑了样本级别的采样频率，对于训练样本 \mathbf{x}_i 及其对应的标记向量 \mathbf{y}_i ，样本级别的采样频率可以估计为：

$$P^I(x_i) = \frac{1}{k} \sum_{y_i^j=1} \frac{1}{n_j} \quad (5-23)$$

基于上述两式，我们综合考虑 $P_j^C(\mathbf{x}_i)$ 和 $P^I(\mathbf{x}_i)$ 来获得最终的权重 r_i^j ：

$$\mathbf{r}_i^j = \frac{P_j^C(\mathbf{x}_i)}{P^I(\mathbf{x}_i)} \quad (5-24)$$

因此，对于少数类标记，并通过优化如下加权 BCE 损失训练得到少数类模型 $C_t(\cdot)$ ：

$$\ell_t(C_t(\mathbf{X}_l), Y_l) = \frac{1}{nk} \sum_{i=1}^n \sum_{j=1}^k BCE(C_t(\mathbf{x}_i), y_i^j) \cdot r_i^j \quad (5-25)$$

综上，我们可以通过联合优化如下目标得到多数类模型和少数类模型：

$$\min_{C_h, C_t} \ell_h(C_h(\mathbf{X}_l), Y_l) + \ell_t(C_t(\mathbf{X}_l), Y_l) \quad (5-26)$$

相关性挖掘。在得到多数类模型 $C_h(\cdot)$ 和少数类模型 $C_t(\cdot)$ 的预测结果之后，我们可以简单的将两个预测向量进行平均来获得最终的预测结果。然而，众所周知，在多标记学习中标记相关性对于学习性能至关重要 [153]。为了挖掘标记之间的相关性，我们进一步提出了一种新颖有效的相关性挖掘网络 $C_R(\cdot)$ ，该模型可以灵活采用如下两种聚合器来对模型 $C_h(\cdot)$ 和 $C_t(\cdot)$ 的预测结果进行聚合，然

后将聚合后的结果映射到新的标记空间，以实现自动探索标记相关性：

- 加法聚合器。第一种聚合策略为加法聚合器，通过对模型 $C_h(\cdot)$ 和 $C_t(\cdot)$ 的预测结果进行加权相加得到：

$$R_i = \mathbf{w}_h C_h(\mathbf{x}_i) + \mathbf{w}_t C_t(\mathbf{x}_i) \quad (5-27)$$

将上述结果 $R_i \in \mathbb{R}^{1 \times k}$ 通过相关性网络 $C_R(\cdot)$ 转换到新的标记空间，得到最终的预测结果：

$$C_R(\mathbf{w}_h C_h(\mathbf{x}) + \mathbf{w}_t C_t(\mathbf{x})) \quad (5-28)$$

- 乘法聚合器。进一步，我们考虑了更复杂的聚合策略，将模型 $C_h(\cdot)$ 和 $C_t(\cdot)$ 的预测结果相乘得到：

$$R_i = C_h(\mathbf{x}_i)^\top \times C_t(\mathbf{x}_i) \quad (5-29)$$

其中 $R_i \in \mathbb{R}^{k \times k}$ 表示标记相关性矩阵。

聚合后的结果 R_i 可以形式化为尺寸为 $\mathbb{R}^{1 \times k^2}$ 的向量，然后输入到相关性挖掘网络 $C_R(\cdot)$ 中， $C_R(\cdot)$ 进一步输出基于 R_i 的预测结果。

乘法聚合器可以看做标记对之间的点积相似性度量 [132]，因此 $C_R(\cdot)$ 可以根据该相似性挖掘训练数据中的标记相关性知识，并进一步优化 $C_h(\cdot)$ 和 $C_t(\cdot)$ 的预测结果，以提升模型性能。

模型 $C_R(\cdot)$ 可以通过优化如下目标训练得到：

$$\min_{C_R} \ell_R(C_R(\mathbf{X}_I), \mathbf{Y}_I) = \frac{1}{n} \sum_{i=1}^n BCE(C_R(R_i), \mathbf{y}_i) \quad (5-30)$$

在训练阶段，标记相关性挖掘网络 $C_R(\cdot)$ 和模型 $C_h(\cdot)$ ， $C_t(\cdot)$ 共同训练，其训练目标可以写作：

$$\min_{C_R, C_h, C_t} \ell_h + \ell_t + \ell_R \quad (5-31)$$

经验结果显示对于大规模标记数据，即，样本对应的标记数量非常大，加法聚合可以取得更好的结果，而当标记的总量较小时，乘法聚合器的性能较好。

标注补全。上述两个模块可以直接部署到标注数据集中，与此同时，如何有效利用无标注数据也是一个核心的问题。LIMI 框架可以灵活的嵌入多种半监督

损失辅助模型训练，基于第四章的研究内容我们可知，针对不同的数据分布需要选取合适的半监督学习模型假设，否则就会导致模型下降的问题 [92]。具体而言，在本章中我们考虑如下两种半监督损失：

- 一致性损失：一致性损失是深度半监督学习方法中最常用的无监督损失函数。一致性损失基于平滑假设，认为特征向量相似的样本应该具有相似的标注向量 [97]。具体而言，定义样本的相似性矩阵为 S ，基于 RBF 核定义的矩阵 S 如下所示：

$$S_{ij} = \exp \frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2} \quad (5-32)$$

LIMI 方法对样本 \mathbf{x} 最终输出的预测概率向量为 $f(\mathbf{x})$ ，即，

$$f(\mathbf{x}) = C_R(C_h(\mathbf{x}), C_t(\mathbf{x})) \quad (5-33)$$

一致性损失可以写作如下形式：

$$\begin{aligned} \ell_u(\mathbf{X}) &= \frac{1}{2} \sum_{i=1}^{n+m} \sum_{j=1}^{n+m} S_{ij} (f(\mathbf{x}_i) - f(\mathbf{x}_j))^2 \\ &= f(\mathbf{X})^\top \mathbf{L} f(\mathbf{X}) \end{aligned} \quad (5-34)$$

其中 $\mathbf{X} = [\mathbf{X}_l, \mathbf{X}_u]$ ，矩阵 $\mathbf{L} = \mathbf{D} - \mathbf{S}$ 表示图拉普拉斯矩阵 (Laplacian Matrix)，矩阵 \mathbf{D} 表示度矩阵 (Degree Matrix)， $D_{ii} = \sum_{j=1}^{n+m} S_{ij}$ ， $i = 1, \dots, n+m$ 。

- 熵最小化损失：根据“没有免费的午餐定理”，我们知道没有一种模型假设是适用所有的数据集和评价指标的。对于不满足平滑假设的数据集，采用一致性损失可能会导致模型性能退化，出现性能不安全的现象，在这种情况下，熵最小化损失是另一种常用的无监督损失函数。熵最小化损失基于低密度分割假设，认为模型的决策边界应该穿过数据密度比较低的区域 [20]，其具体形式如下所示：

$$\ell_u(\mathbf{X}) = - \sum_{i=1}^{n+m} \sum_{j=1}^k f(\mathbf{x}_i)_j \log(f(\mathbf{x}_i)_j) \quad (5-35)$$

综上所述，本章提出的类别不平衡半监督多标记学习框架 LIMI 包括三个核心组成部分：多数类模型 $C_h(\cdot)$ ，少数类模型 $C_t(\cdot)$ 和标记相关性挖掘网络 $C_R(\cdot)$ ，

LIMI 通过最小化如下损失函数同时优化这三个网络：

$$\min_{C_h, C_t, C_R} \ell_h + \ell_t + \ell_R + \lambda \ell_u \quad (5-36)$$

其中超参数 λ 表示标注样本上的监督损失与无标注样本上的无监督损失之间的平衡超参数。

LIMI 方法的整体框架图如图 5-3 所示。

5.4 实验验证

在本节，我们通过实验验证本文提出的类别不平衡的半监督单标记学习方法 CWSL 以及类别不平衡的半监督多标记学习方法 LIMIM 在机器学习常用数据集以及网约车共享平台滴滴出行真实任务数据上的性能表现。

5.4.1 UCI 数据集

我们首先在 UCI 数据集 *breast_cancer*^① 上进行实验来展示对于包含噪声伪标注的数据集，CWSL 学得的样本权重在训练过程中的变化。

Breast_cancer 是一个二分类的数据集，其中包括 569 个训练样本，每个训练样本具有 32 维的属性。我们将该数据集分成 3 部分：469 个训练样本，50 个验证样本以及 50 个测试样本。对于训练数据，我们随机选择 100 个样本，然后反转其类别标注，作为噪声样本。

在该实验中，我们采用一个两层神经网络作为分类模型，将二分类的交叉熵损失（Binary Cross-Entropy）作为模型训练的损失函数。在该数据中不存在类别不平衡的问题，因此我们没有采用 AUC 作为优化指标。模型训练的迭代轮数为 100，样本权重初始化为 0.5。

图 5-4 展示了训练过程中样本权重的变化。其中，图 (a) 的结果说明仅仅通过一轮优化，我们的方法就可以将正确标注的样本和噪声标注的样本区分开。在 10 轮优化之后，正确标注的样本权重已经集中在 1.0 附近。在 100 轮优化之后，正确标注的样本和噪声标注的样本权重已经几乎全部在 1.0 和 0.0 附近。这

^①<https://archive.ics.uci.edu>

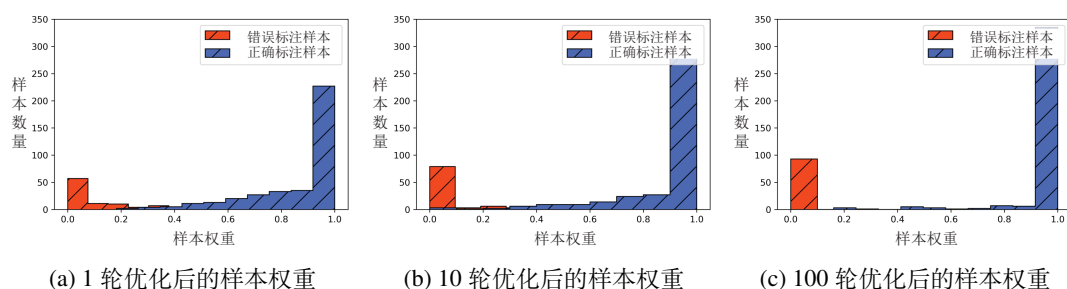


图 5-4: 经过 1 轮、10 轮、100 轮优化后的样本权重分布。样本权重的变化说明了我们的方法能够成功的将正确标注的样本和错误标注的样本区分开，通过赋予错误标注样本较低的权重，降低其对模型性能的损害。

证明了我们的方法可以有效的将正确标注的样本和噪声标注的样本区分开，通过赋予噪声样本较低的权重，降低其对模型性能的负面影响。

5.4.2 网约车智能评价任务

在网约车平台中，乘客输入出发地和目的地之后，平台将自动匹配附近的司机接载乘客。为了提升服务质量和司乘体验，在每次订单结束后，平台都会邀请乘客对司机进行评价，例如“车辆是否有异味”、“司机是否绕路”等。如果能够收集到负面的评价，则有助于平台进行司机管理，提升服务质量。因此，该任务的目标是建立机器学习模型，利用司乘历史信息预测司机可能出现差评的问题，进行差评聚焦。

在该任务中，我们基于滴滴平台 2019 年 3 月 1 日至 2019 年 4 月 1 日期间网约车订单的真实评价数据构建训练数据集。我们随机采样了 600,000 个样本作为训练数据，15,000 个专家验证的样本平均划分为验证集和测试集。数据分布中正负类样本的比例约为 1:30。评价数据的原始特征包括数百维司乘历史特征以及订单特征，司乘历史特征包括驾驶员年龄、性别、过去 10 天好评 / 差评次数等，订单特征包括订单所在区域、订单价格、目的地等。原始特征中既包括离散特征又包括连续特征，简单的在原始特征上训练机器学习模型效果不佳，而手动进行特征工程又需要耗费大量的时间精力。受 [62] 启发，我们采用 XGBOOST 算法 [27] 对原始数据进行特征转化，具体而言，我们将每棵单独的树视为一个分类特征，将样本最终落入的叶子节点的索引值作为该特征的取值，然后将转化后的特征进行归一化作为最终的特征。此类做法在广告推荐、点击率预测等同样需要复杂特征工程的现实任务中取得了成功的应用 [62]。

在实验中，我们将 CWSL 算法和如下算法进行对比，首先是常用的监督学习算法：

- XGBoost [27]: 一种高效的基于决策树集成的梯度提升树算法 [46]。XGBoost 算法的有效性已经在多种机器学习任务以及数据挖掘竞赛中得到了验证。对于 XGBoost 算法，我们将正类样本的权重设为 4 以缓解类别不平衡问题，其它参数均设置为默认值。
- LR (Logistic Regression) 算法: LR 算法可以看做分类任务的基线算法。对于 LR 算法，我们采用开源机器学习库 sklearn 中的实现，参数均设置为默认值。
- DNN (Deep Neural Network) 算法: 深度神经网络算法近年来在图像分类等任务中取得了优异的性能表现。我们采用了一个 3 层的神经网络作为分类模型，采用 ReLU 作为隐层激活函数，Dropout 的值设置为 0.5，每层大小设置为 64。我们采用交叉熵作为损失函数，并利用 SGD 算法进行模型的优化，其中 SGD 的优化步长设置为 0.01。

同时，为了缓解真实评价数据中类别不平衡的问题，我们采用预测概率校正法 [32] 对以上监督学习模型的预测结果进行概率校正。具体而言，令 β 表示训练数据中正类样本的采样比例， $0 < \beta < 1$ ， p_s 表示样本 \mathbf{x} 预测为负类的概率，则模型预测概率可以按照如下公式进行校正：

$$p = \frac{p_s \beta}{p_s \beta - p_s + 1} \quad (5-37)$$

此外，由于我们将无标注数据全部当做正类样本进行处理，因此训练数据中会包含标注错误的样本，为了验证我们方法的性能，我们进一步与能够处理标注噪声样本的噪声学习方法进行比较：

- Rank Pruning 算法 [108]: Rank Pruning 方法是一种先进的噪声标注筛选算法。对于二分类任务，该算法根据模型预测概率对样本进行排序，然后估计噪声的比例，根据噪声样本的比例将相应数量的预测置信度较低的样本进行删除。理论上，Rank Pruning 算法可以实现相当于利用不含噪声的数据集进行训练的模型性能。对于 Rank Pruning 算法，我们采用 XGBoost 作为基础的分类模型，参数设置为默认值。
- GLC (Golden Loss Correction) 算法 [64]: GLC 算法引入一个干净无偏的验证

集辅助训练，基于验证集估计出一个噪声纠正矩阵，并利用该矩阵进行模型预测概率的校正。GLC 算法在多种数据集，如 CIFAR-10, IMDB 中取得了最优的性能。对于 GLC 算法我们采用 DNN 模型作为基础的分类模型。

- LTR (Learning to Re-weight Examples) 算法 [116]: LTR 算法同样利用验证集进行模型训练，并提出了一种基于双层优化的样本赋权机制。与我们的方案不同，LTR 是基于验证集上的模型损失进行权重的学习，而我们的方法是直接优化模型的 AUC 性能。对于 LTR 算法，我们采用同样结构的 DNN 模型作为分类模型。

对于不能利用验证数据的算法，如 XGBoost、LR、DNN、Rank Pruning，我们将验证数据合并到训练集中辅助模型训练，以保证算法比较的公平性。所有的模型都是在转化后的特征上进行训练。对于 CWSL 算法，样本的初始权重设置为 0.5，模型优化轮数为 100，步长 λ_o 和 λ_w 分别设置为 0.1 和 0.4。

主要实验结果。实验的主要结果如表 5-1 所示，从表 5-1 中我们可以看出，本文提出的 CWSL 方法显著优于对比方法。具体而言，相比与 XGBoost 算法，CWSL 取得了超过 12% 的性能提升；相比 LTR 算法，CWSL 取得了近 6% 的性能提升。这证明了之前的算法在类别不平衡的机器学习问题上并不能取得理想的性能，而我们的方法通过直接优化 AUC，在类别不平衡的数据集上保证了性能的稳健性。我们同时汇报了各对比算法的运行时间，所有的实验均在同一台 10 核 2.20GHz 的 Intel Xeon (R) CPU 和 32GB 内存的计算机上运行。从实验结果我们可以看出，相比与其它两种基于双层优化的算法 LTR 和 GLC，CWSL 在运行效率上取得了显著的提升，而采用 XGBoost 作为基础分类模型的 Rank Pruning 算法运行效率最低；相比与三种监督学习算法，CWSL 运行效率和 DNN 类似，并且优于 XGBoost。这些结果证明了 CWSL 方法不仅在性能上取得了最佳的表现，同时也具有很高的运行效率，支持其处理大规模的数据集。

验证集数据分布。由于我们的方法依赖验证数据集，我们进一步研究验证数据集中类别不平衡的比例对模型性能的影响。我们在三种基于验证集的方法 (LTR, GLC 和我们提出的 CWSL 方法) 上进行实验，验证集中正负类的比例从 1:1 到超过 1:40 不等。图 5-5 中汇报了随验证集类别比例的变化模型 AUC 性能的变化情况。从实验结果可以看出，三种方法呈现出一致的变化趋势，当验证集中类别比例与真实数据分布的类别比例偏差较大时，三种方法性能表现较差，

表 5-1: 网约车智能评价任务 AUC 性能与时间开销。对于每种方法, 我们汇报了进行 10 次实验 AUC 指标的均值与标准差, 其中加粗数字代表最优的方法。

实验方法	AUC	时间开销 (s)
LR	78.64 ± 1.23	47.8
DNN	77.85 ± 1.17	90.4
XGBoost	79.55 ± 1.13	324.3
Rank Pruning	82.22 ± 0.53	601.4
GLC	82.63 ± 0.49	300.5
LTR	85.45 ± 0.51	211.6
CWSL	91.24 ± 0.47	107.1

当正负类比例超过 1: 10 之后, 模型的性能变化情况趋于稳定, 而当验证集中正负类的比例接近真实数据中正负类样本比例 (1: 30) 时, 三种算法均可实现较优的性能。同时, 不管类别比例如何变化, 我们的方法始终优于其他两种方法, 这进一步说明了我们方法的优越性。

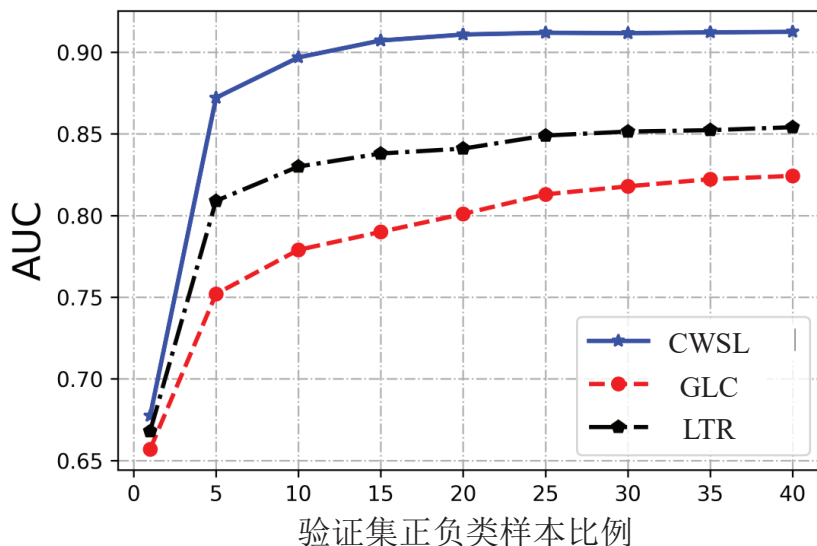


图 5-5: 随验证集正负类比例的变化, 模型 AUC 性能的变化情况。

评价指标稳健性。虽然我们的方法是基于 AUC 指标进行优化, 研究其在其它指标下的性能表现依然是有意义的。表 5-2 汇报了对比方法在 Precision、Recall 和 F1-Score 指标下的性能, 这些指标均是用于评估排序质量的常用指标。对于二分类问题, 令真实标注为正类的样本集合为 G , 预测标注为正类的样本集合

为 T ，则 Precision 和 Recall 的定义如下所示：

$$\text{Precision} = \frac{|G \cap T|}{|G|}, \quad \text{Recall} = \frac{|T \cap G|}{|T|} \quad (5-38)$$

Precision 和 Recall 是一对矛盾的度量，一般来说，Precision 高时，Recall 往往偏低，而 Recall 高时，Precision 往往偏低。F1-Score 综合考虑了这两种指标，是基于 Precision 和 Recall 的调和平均定义的，如下所示：

$$\text{F1-Score} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5-39)$$

从表 5-2 的结果来看，LTR 算法的 Recall 最高但是 Precision 最低，DNN、LR 和 GLC 方法在这三种评价指标上取得了相似的结果。本文提出的 CWSL 方法在 Precision 和 Recall 上取得了最好的折衷，并在 F1-Score 上取得了最好的结果。这些结果表明，尽管我们的方法针对 AUC 评价指标进行优化，但对评价指标的变化具有稳健性，在其它指标上仍然能取得良好的性能表现。

表 5-2: Precision, Recall 和 F1-Score 指标下的算法性能。

实验方法	Precision	Recall	F1-Score
LR	24.08%	34.69%	26.43%
DNN	23.26%	35.34%	25.71%
XGBoost	9.51%	69.80%	16.73%
Rank Pruning	12.26%	54.73%	18.27%
GLC	24.13%	34.24%	26.79%
LTR	8.63%	87.35%	15.70%
CWSL	16.50%	80.82 %	27.41%

收敛性验证。我们进一步进行实验分析我们采用的针对双层优化的交替优化方法与普通的随机梯度下降方法的收敛性。具体而言，针对 DNN 算法和本文提出的算法，训练集的交叉熵损失以及验证集上的 AUC 性能随优化迭代轮数的变化如图 5-6 所示。从实验结果中我们可以看到，我们的方法所需的收敛轮数和基于随机梯度下降的 DNN 算法所需的收敛轮数相当，这证明了我们提出的优化方法具有高效的收敛性。

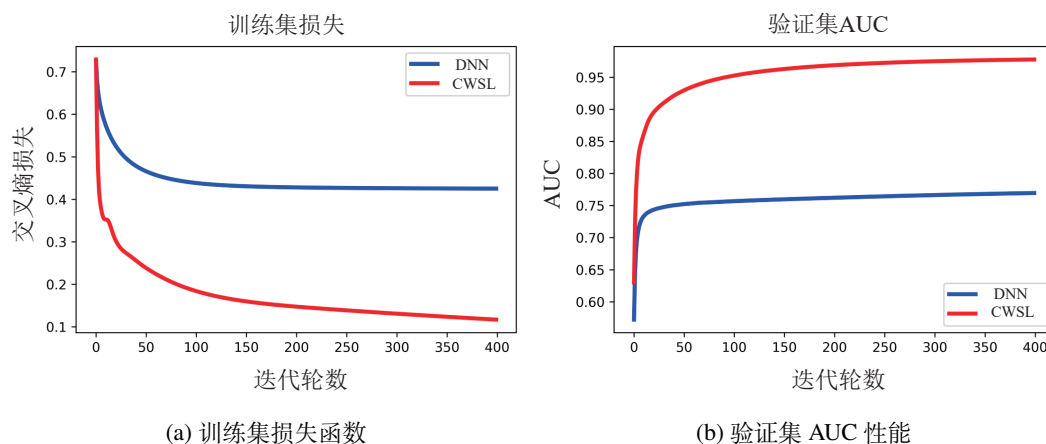


图 5-6: 训练集损失函数以及验证集 AUC 随训练轮数变化的曲线。

5.4.3 网约车智能判责任务

在本节，我们介绍本文提出的类别不平衡的半监督多标记学习算法 LIM1 在网约车智能判责任务中的实验结果。在网约车平台中，乘客和司机一旦产生投诉，平台需要对责任进行判定，从而进行处罚和管理，帮助平台提升服务质量和用户体验。该任务是一个多标记问题，需要模型同时在多种多标记评价指标上取得良好性能。

数据信息。在网约车智能判责任务中，我们以滴滴平台 2020 年 11 月 5 日至 2020 年 11 月 23 日期间产生的判责数据构造网约车智能判责数据集，共计 25,018 个样本。每个样本的特征由三部分组成：表格型特征、文本型特征和司乘对话信息，其中表格型特征包括司机乘客的历史信息、订单特征等，共计 85 维，对于文本型特征，我们采用 TextCNN 模型 [74] 进行预处理，获得 192 维的特征向量，对于对话信息，我们采用 HAN 模型 [145] 进行特征处理，获得 200 维度的特征向量，将上述三部分特征拼接在一起，每个样本的特征向量维度为 477。我们采用六种常见的判责原因作为样本所对应的标记，多标记数据的类别不平衡比例定义为每个标记上正负类样本比例的均值，在网约车智能判责数据中，类别不平衡的比例为 16.9。

评价指标。对于多标记学习任务，需要同时考虑样本级别的性能和标记级别的性能 [136]。在本文中我们采用了九种多标记学习常用评价指标：Hamming Loss、Ranking Loss、One Error、Coverage、Average Precision、Macro AUC、Micro AUC、Macro F1、Micro F1 对模型进行评估，其定义分别为：

- Hamming Loss:

$$\frac{1}{NK} \sum_{i=1}^N \sum_{j=1}^K \mathbb{I}(\hat{y}_{ij} \neq y_{ij}) \quad (5-40)$$

- Ranking Loss:

$$\frac{1}{N} \sum_{i=1}^N \frac{|\mathcal{S}_{rank}^i|}{|Y_i^+||Y_i^-|} \quad (5-41)$$

- One Error:

$$\frac{1}{N} \sum_{i=1}^N \mathbb{I}(\operatorname{argmax} f(x_i) \notin Y_i^+) \quad (5-42)$$

- Coverage:

$$\frac{1}{N} \sum_{i=1}^N \mathbb{I}(\max_{j \in Y_i^+} \operatorname{rank}_f(\mathbf{x}_i, j) - 1) \quad (5-43)$$

- Average Precision:

$$\frac{1}{N} \sum_{i=1}^N \frac{1}{|Y_i^+|} \sum_{j \in Y_i^+} \frac{|\mathcal{S}_{precision}^{ij}|}{\operatorname{rank}_f(\mathbf{x}_i, j)} \quad (5-44)$$

- Macro AUC:

$$\frac{1}{K} \sum_{j=1}^K \frac{|\mathcal{S}_{macro}^j|}{|Y_j^+||Y_j^-|} \quad (5-45)$$

- Micro AUC:

$$\frac{|\mathcal{S}_{micro}|}{(\sum_{i=1}^N |Y_i^+|) \cdot (\sum_{i=1}^N |Y_i^-|)} \quad (5-46)$$

- Macro F1:

$$\frac{1}{K} \sum_{j=1}^K \frac{2 \sum_{i=1}^N y_{ij} \hat{y}_{ij}}{\sum_{i=1}^N y_{ij} + \sum_{i=1}^N \hat{y}_{ij}} \quad (5-47)$$

- Micro F1:

$$\frac{2 \sum_{j=1}^K \sum_{i=1}^N y_{ij} \hat{y}_{ij}}{\sum_{j=1}^K \sum_{i=1}^N y_{ij} + \sum_{j=1}^K \sum_{i=1}^N \hat{y}_{ij}} \quad (5-48)$$

其中

$$\begin{aligned} \mathcal{S}_{rank}^i &= \{(u, v) | f(\mathbf{x}_i)_u \leq f(\mathbf{x}_i)_v, (u, v) \in Y_i^+ \times Y_i^-\} \\ \mathcal{S}_{precision}^{ij} &= \{k \in Y_i^+ | \operatorname{rank}_f(\mathbf{x}_i, k) \leq \operatorname{rank}_f(\mathbf{x}_i, j)\} \\ \mathcal{S}_{macro}^j &= \{(a, b) \in Y_j^+ \times Y_j^- | f(x_a)_j \geq f(x_b)_j\} \\ \mathcal{S}_{micro} &= \{(a, b, i, j) | (a, b) \in Y_i^+ \times Y_j^-, f(\mathbf{x}_a)_i \geq f(\mathbf{x}_b)_j\} \end{aligned}$$

对比方法。为了验证本章提出的 LIMM 方法的性能，我们和如下方法进行对比，包括基线监督学习方法、多标记学习方法、类别不平衡多标记学习方法以及半监督学习方法。

- FCN (Fully Connected Network) 算法：直接在所有标注数据上以监督学习的方式通过优化标准的 BCE 损失训练一个全连接神经网络。该方法可以看做是该任务的基线算法，既没有考虑到类别不平衡性，也没有利用无标注数据。
- CAMEL 算法 [43]：CAMEL 算法是代表性的多标记学习算法，该算法首先在标记空间中通过稀疏重构学习标记之间的相关关系，然后将相关关系融入到模型训练中，提升多标记学习模型的性能。
- DBL 算法 [135]：DBL 算法是一种能够处理类别不平衡的多标记学习算法，该算法通过对每个标记进行重新加权降低类别不平衡问题的影响，同时在加权过程中考虑到标记共现性，使加权方法能够适应多标记数据。
- DBL+NT 算法 [135]：在 DBL 方法的基础上，进一步提出了一个 NT (Negative Tolerant) 机制，以减轻 DBL 算法对于不相关标记的过度抑制。
- PL (Pseudo-Labeling) 算法 [84]：PL 算法是一种代表性的半监督学习算法，通过为无标注数据赋予预测置信度较高的伪标注，不断扩充标注数据集的规模，从而提升模型性能。
- DRML [132]：DRML 算法是一种代表性的半监督多标记学习方法，该算法采用两个分类器同时进行训练，以实现特征分布和标记分布的对齐，同时采用一个关联网络挖掘标记之间的相关性。
- DRML+DBL：DBL 算法只考虑了类别不平衡的多标记学习，DRML 算法只考虑了半监督的多标记学习，我们通过将 DRML 和 DBL 进行结合，综合考虑了类别不平衡的半监督多标记学习问题。

参数设置。对于数据集，我们将其按照 7:1:2 的比例划分为训练数据、验证数据和测试数据。在训练数据中，我们考虑两种半监督的设置，即随机采样 5% 或 10% 的训练样本作为标注数据，其余样本作为无标注数据。对于我们提出的 LIMM 方法，我们利用 Adam 算法优化 500 轮，算法学习率设置为 0.01，早停 (Early Stopping) 参数为 30，监督损失与无监督损失的平衡超参数 λ 为 1。对于多数类模型 $C_h(\cdot)$ 和少数类模型 $C_l(\cdot)$ ，我们采用神经网络作为实现，其网络结构为 $[d, 256, 64, k]$ 。聚合器模块和半监督正则项通过验证集的性能进行选择。对于

表 5-3: 网约车智能判责任务实验结果 (均值 \pm 标准差)。 \uparrow (\downarrow) 表示数值越大 (小) 性能越好, 加粗的数字表示最佳性能。

对比方法	5% 标注样本									
	Hamming Loss \downarrow	Ranking Loss \downarrow	One Error \downarrow	Coverage \downarrow	Ave. Precision \uparrow	Macro AUC \uparrow	Micro AUC \uparrow	Macro F1 \uparrow	Micro F1 \uparrow	
FCN	0.160 \pm 0.008	0.100 \pm 0.001	0.237 \pm 0.001	0.121 \pm 0.001	0.839 \pm 0.001	0.588 \pm 0.012	0.885 \pm 0.004	0.239 \pm 0.014	0.584 \pm 0.016	
CAMEL	0.159 \pm 0.006	0.115 \pm 0.003	0.243 \pm 0.002	0.135 \pm 0.002	0.831 \pm 0.001	0.593\pm0.006	0.878 \pm 0.002	0.239 \pm 0.006	0.585 \pm 0.017	
DBL	0.160 \pm 0.008	0.100 \pm 0.001	0.236\pm0.001	0.121 \pm 0.001	0.840\pm0.001	0.586 \pm 0.017	0.885 \pm 0.005	0.237 \pm 0.017	0.583 \pm 0.015	
DBL+NT	0.157 \pm 0.006	0.100 \pm 0.001	0.237 \pm 0.002	0.121 \pm 0.001	0.839 \pm 0.001	0.592 \pm 0.009	0.884 \pm 0.003	0.243\pm0.007	0.591 \pm 0.013	
PL	0.160 \pm 0.008	0.100 \pm 0.001	0.237 \pm 0.002	0.121 \pm 0.001	0.839 \pm 0.001	0.586 \pm 0.012	0.886 \pm 0.003	0.237 \pm 0.010	0.584 \pm 0.016	
DRML	0.161 \pm 0.013	0.108 \pm 0.002	0.264 \pm 0.003	0.127 \pm 0.002	0.825 \pm 0.002	0.524 \pm 0.016	0.870 \pm 0.005	0.217 \pm 0.007	0.579 \pm 0.026	
DRML+DBL	0.164 \pm 0.009	0.116 \pm 0.008	0.276 \pm 0.019	0.135 \pm 0.006	0.816 \pm 0.011	0.533 \pm 0.011	0.858 \pm 0.010	0.221 \pm 0.005	0.574 \pm 0.022	
LIMI	0.152\pm0.010	0.099\pm0.001	0.236\pm0.002	0.120\pm0.001	0.840\pm0.001	0.589 \pm 0.025	0.889\pm0.004	0.232 \pm 0.014	0.597\pm0.019	

对比方法	10% 标注样本									
	Hamming Loss \downarrow	Ranking Loss \downarrow	One Error \downarrow	Coverage \downarrow	Ave. Precision \uparrow	Macro AUC \uparrow	Micro AUC \uparrow	Macro F1 \uparrow	Micro F1 \uparrow	
FCN	0.155 \pm 0.006	0.098 \pm 0.001	0.237 \pm 0.002	0.120 \pm 0.001	0.841 \pm 0.001	0.607 \pm 0.016	0.889 \pm 0.004	0.249 \pm 0.015	0.596 \pm 0.013	
CAMEL	0.156 \pm 0.003	0.112 \pm 0.001	0.240 \pm 0.002	0.132 \pm 0.001	0.834 \pm 0.001	0.606 \pm 0.005	0.882 \pm 0.001	0.248 \pm 0.007	0.593 \pm 0.009	
DBL	0.156 \pm 0.006	0.098 \pm 0.002	0.236 \pm 0.001	0.120 \pm 0.001	0.841 \pm 0.001	0.604 \pm 0.019	0.889 \pm 0.004	0.249 \pm 0.013	0.596 \pm 0.014	
DBL+NT	0.154 \pm 0.006	0.098 \pm 0.001	0.236 \pm 0.001	0.120 \pm 0.001	0.841 \pm 0.001	0.601 \pm 0.018	0.889 \pm 0.004	0.245 \pm 0.013	0.599 \pm 0.012	
PL	0.154 \pm 0.005	0.098 \pm 0.001	0.237 \pm 0.002	0.120 \pm 0.001	0.840 \pm 0.001	0.613 \pm 0.007	0.891 \pm 0.002	0.255 \pm 0.007	0.600\pm0.009	
DRML	0.160 \pm 0.006	0.103 \pm 0.002	0.255 \pm 0.011	0.123 \pm 0.001	0.832 \pm 0.005	0.530 \pm 0.011	0.876 \pm 0.004	0.223 \pm 0.005	0.584 \pm 0.012	
DRML+DBL	0.160 \pm 0.008	0.105 \pm 0.002	0.252 \pm 0.011	0.125 \pm 0.002	0.831 \pm 0.005	0.533 \pm 0.013	0.871 \pm 0.004	0.222 \pm 0.005	0.585 \pm 0.017	
LIMI	0.150\pm0.005	0.096\pm0.002	0.234\pm0.001	0.115\pm0.002	0.843\pm0.002	0.614\pm0.023	0.892\pm0.010	0.260\pm0.010	0.599 \pm 0.012	

乘法聚合器，网络 $C_R(\cdot)$ 的结构为 $[k^2, k]$ ，对于加法聚合器，网络 $C_R(\cdot)$ 的结构为 $[2 \times k, k]$ 。为了降低随机误差，每个算法运行 10 次并汇报在不同指标下的均值和标准差。

实验结果。在网约车智能判责任务上的实验结果如表 5-3 所示。从实验结果我们可以看出，目前最先进的多标记学习算法，如 CAMEL，相比于基线监督学习方法 FCN 并没有取得明显的性能提升；类别不平衡的多标记学习方法，例如 DBL，在某些指标中能够取得良好的性能，但也在很多情况下比基线的 FCN 方法表现更差。主要原因是，这些方法依赖足够数量的标记数据，缺少稳健利用无标注数据的学习机制，导致其在标记数量不足的半监督学习场景中不能很好的工作。半监督多标记学习方法，如 DRML，在多个指标上表现不佳，主要原因是其无法处理多标记数据中类别不平衡的问题，在少数类标记上不能取得良好的性能。此外，简单地将半监督学习和类别不平衡学习算法进行结合也不能很好的解决该问题，反而会出现性能严重退化的情况。相比之下，本文所提出的 LIMM 方法几乎在每个指标上都实现了最优的性能。上述结果证明了我们的方法在类别不平衡半监督标记学习问题中的有效性，也证明了该方法能够很好地应用于工业界实际任务场景。

5.4.4 图像识别任务

为了进一步证明本文提出的 LIMM 方法在类别不平衡的半监督多标记学习问题上的有效性，我们还在多种基准的多标记学习数据集中进行实验。首先是机器学习中最常见的图像识别任务，我们采用多标记图像识别任务的基准数据集 CUB^①作为训练数据。CUB 数据集是一个鸟类识别数据集，包括 10,240 张图片，每张图片的标记维度为 312 维，数据集的类别不平衡比例为 57.72。对于每张图片，我们采用在 ImageNet 数据集 [34] 上预训练的 VGG 网络 [122] 作为特征提取器。

实验结果如表 5-4 所示。从结果中我们可以看到，目前最先进的对比方法，如 DBL、DRML 等，在多种评价指标上均遇到了相比基线的 FCN 算法性能下降的问题，这说明这些方法是性能不安全的，而本文提出的 LIMM 方法在多种评价指标上始终保持着性能的领先，这进一步证明了 LIMM 方法的有效性。

^①<http://www.vision.caltech.edu/visipedia/CUB-200-2011.html>

表 5-4: 图像识别任务实验结果 (均值 \pm 标准差)。 \uparrow (\downarrow) 表示数值越大 (小) 性能越好, 加粗的数字表示最佳性能。

5% 标注样本										
对比方法	Hamming	Ranking	One	Coverage	Ave.	Macro	Micro	Macro	Micro	
	\downarrow Loss	\downarrow Loss	\downarrow Error	\downarrow	\uparrow Precision	\uparrow AUC	\uparrow AUC	\uparrow F1	\uparrow F1	
FCN	0.123 \pm 0.004	0.142 \pm 0.002	0.191 \pm 0.016	0.597 \pm 0.007	0.503 \pm 0.005	0.679 \pm 0.005	0.852 \pm 0.002	0.225 \pm 0.007	0.413 \pm 0.011	
CAMEL	0.114\pm0.001	0.158 \pm 0.001	0.186 \pm 0.011	0.786 \pm 0.007	0.496 \pm 0.003	0.682 \pm 0.006	0.842 \pm 0.001	0.245\pm0.003	0.430 \pm 0.006	
DBL	0.121 \pm 0.004	0.140 \pm 0.003	0.178 \pm 0.021	0.587 \pm 0.007	0.504 \pm 0.006	0.684 \pm 0.007	0.856 \pm 0.003	0.230 \pm 0.011	0.414 \pm 0.013	
DBL+NT	0.123 \pm 0.003	0.142 \pm 0.003	0.150\pm0.008	0.589 \pm 0.007	0.503 \pm 0.006	0.680 \pm 0.008	0.854 \pm 0.003	0.225 \pm 0.009	0.410 \pm 0.010	
PL	0.118 \pm 0.002	0.138 \pm 0.002	0.201 \pm 0.014	0.590 \pm 0.007	0.508 \pm 0.005	0.683 \pm 0.004	0.852 \pm 0.003	0.239 \pm 0.003	0.424 \pm 0.008	
DRML	0.121 \pm 0.004	0.166 \pm 0.002	0.265 \pm 0.034	0.706 \pm 0.013	0.464 \pm 0.007	0.621 \pm 0.003	0.831 \pm 0.003	0.193 \pm 0.005	0.395 \pm 0.008	
DRML+DBL	0.121 \pm 0.002	0.156 \pm 0.002	0.161 \pm 0.017	0.625 \pm 0.008	0.484 \pm 0.003	0.645 \pm 0.006	0.841 \pm 0.002	0.200 \pm 0.004	0.399 \pm 0.005	
LIMI	0.114\pm0.002	0.134\pm0.002	0.165 \pm 0.018	0.585\pm0.007	0.512\pm0.007	0.688\pm0.004	0.858\pm0.004	0.240 \pm 0.004	0.436\pm0.009	
10% 标注样本										
对比方法	Hamming	Ranking	One	Coverage	Ave.	Macro	Micro	Macro	Micro	
	\downarrow Loss	\downarrow Loss	\downarrow Error	\downarrow	\uparrow Precision	\uparrow AUC	\uparrow AUC	\uparrow F1	\uparrow F1	
FCN	0.115 \pm 0.002	0.131 \pm 0.003	0.183 \pm 0.013	0.562 \pm 0.007	0.525 \pm 0.007	0.708 \pm 0.007	0.865 \pm 0.002	0.254 \pm 0.009	0.440 \pm 0.007	
CAMEL	0.112 \pm 0.001	0.151 \pm 0.002	0.174 \pm 0.004	0.767 \pm 0.005	0.510 \pm 0.003	0.691 \pm 0.005	0.849 \pm 0.002	0.256 \pm 0.005	0.444 \pm 0.004	
DBL	0.115 \pm 0.002	0.130 \pm 0.002	0.160 \pm 0.013	0.552 \pm 0.006	0.525 \pm 0.005	0.707 \pm 0.006	0.867 \pm 0.003	0.255 \pm 0.005	0.438 \pm 0.004	
DBL+NT	0.117 \pm 0.002	0.132 \pm 0.002	0.144\pm0.010	0.558 \pm 0.006	0.523 \pm 0.006	0.703 \pm 0.005	0.866 \pm 0.002	0.249 \pm 0.006	0.431 \pm 0.004	
PL	0.114 \pm 0.002	0.129\pm0.002	0.178 \pm 0.014	0.556 \pm 0.006	0.526 \pm 0.005	0.710\pm0.006	0.865 \pm 0.002	0.258\pm0.004	0.441 \pm 0.003	
DRML	0.116 \pm 0.001	0.157 \pm 0.004	0.277 \pm 0.032	0.676 \pm 0.014	0.475 \pm 0.006	0.641 \pm 0.008	0.840 \pm 0.003	0.214 \pm 0.009	0.418 \pm 0.003	
DRML+DBL	0.117 \pm 0.002	0.148 \pm 0.003	0.185 \pm 0.033	0.609 \pm 0.007	0.493 \pm 0.009	0.670 \pm 0.010	0.849 \pm 0.003	0.221 \pm 0.009	0.420 \pm 0.004	
LIMI	0.102\pm0.002	0.129\pm0.002	0.151 \pm 0.016	0.551\pm0.008	0.527\pm0.006	0.710\pm0.004	0.869\pm0.002	0.255 \pm 0.006	0.445\pm0.003	

表 5-5: 文本分类任务实验结果 (均值 \pm 标准差)。 \uparrow (\downarrow) 表示数值越大 (小) 性能越好, 加粗的数字表示最佳性能。

对比方法	5% 标注样本									
	Hamming Loss \downarrow	Ranking Loss \downarrow	One Error \downarrow	Coverage \downarrow	Ave. Precision \uparrow	Macro AUC \uparrow	Micro AUC \uparrow	Macro F1 \uparrow	Micro F1 \uparrow	
FCN	0.025 \pm 0.001	0.163 \pm 0.006	0.572 \pm 0.012	0.263 \pm 0.008	0.379 \pm 0.009	0.686 \pm 0.016	0.727 \pm 0.015	0.113 \pm 0.013	0.203 \pm 0.017	
CAMEL	0.024\pm0.001	0.176 \pm 0.005	0.587 \pm 0.015	0.290 \pm 0.008	0.366 \pm 0.008	0.803\pm0.006	0.825\pm0.005	0.210\pm0.008	0.284\pm0.009	
DBL	0.025 \pm 0.001	0.163 \pm 0.006	0.568 \pm 0.011	0.264 \pm 0.008	0.382 \pm 0.009	0.680 \pm 0.016	0.723 \pm 0.016	0.115 \pm 0.013	0.204 \pm 0.016	
DBL+NT	0.024\pm0.001	0.164 \pm 0.005	0.563\pm0.013	0.267 \pm 0.008	0.386 \pm 0.008	0.693 \pm 0.017	0.735 \pm 0.016	0.133 \pm 0.013	0.224 \pm 0.017	
PL	0.027 \pm 0.001	0.182 \pm 0.006	0.594 \pm 0.021	0.291 \pm 0.008	0.354 \pm 0.013	0.594 \pm 0.033	0.635 \pm 0.033	0.079 \pm 0.013	0.153 \pm 0.017	
DRML	0.024 \pm 0.002	0.343 \pm 0.009	0.802 \pm 0.042	0.490 \pm 0.012	0.161 \pm 0.029	0.571 \pm 0.035	0.618 \pm 0.018	0.060 \pm 0.019	0.144 \pm 0.033	
DRML+DBL	0.028 \pm 0.002	0.289 \pm 0.011	0.747 \pm 0.021	0.426 \pm 0.012	0.211 \pm 0.016	0.585 \pm 0.014	0.635 \pm 0.011	0.044 \pm 0.008	0.102 \pm 0.015	
LIMI	0.024\pm0.001	0.160\pm0.007	0.563\pm0.014	0.262\pm0.009	0.387\pm0.011	0.696 \pm 0.013	0.737 \pm 0.013	0.130 \pm 0.013	0.206 \pm 0.014	

对比方法	10% 标注样本									
	Hamming Loss \downarrow	Ranking Loss \downarrow	One Error \downarrow	Coverage \downarrow	Ave. Precision \uparrow	Macro AUC \uparrow	Micro AUC \uparrow	Macro F1 \uparrow	Micro F1 \uparrow	
FCN	0.021\pm0.001	0.127 \pm 0.003	0.508 \pm 0.008	0.210 \pm 0.004	0.438 \pm 0.006	0.755 \pm 0.009	0.795 \pm 0.008	0.156 \pm 0.009	0.264 \pm 0.014	
CAMEL	0.021\pm0.001	0.144 \pm 0.003	0.508 \pm 0.006	0.248 \pm 0.007	0.436 \pm 0.004	0.833\pm0.006	0.856\pm0.005	0.256\pm0.008	0.334\pm0.008	
DBL	0.021\pm0.001	0.127 \pm 0.003	0.508 \pm 0.009	0.210 \pm 0.004	0.440 \pm 0.006	0.751 \pm 0.009	0.791 \pm 0.008	0.155 \pm 0.010	0.263 \pm 0.015	
DBL+NT	0.021\pm0.001	0.128 \pm 0.003	0.507 \pm 0.008	0.214 \pm 0.005	0.444 \pm 0.006	0.757 \pm 0.010	0.797 \pm 0.009	0.173 \pm 0.010	0.280 \pm 0.014	
PL	0.022 \pm 0.001	0.133 \pm 0.004	0.519 \pm 0.008	0.218 \pm 0.008	0.428 \pm 0.006	0.705 \pm 0.017	0.747 \pm 0.016	0.125 \pm 0.014	0.228 \pm 0.018	
DRML	0.022 \pm 0.004	0.310 \pm 0.008	0.731 \pm 0.018	0.450 \pm 0.010	0.214 \pm 0.010	0.658 \pm 0.021	0.676 \pm 0.010	0.105 \pm 0.012	0.219 \pm 0.013	
DRML+DBL	0.024 \pm 0.002	0.237 \pm 0.011	0.671 \pm 0.017	0.352 \pm 0.013	0.272 \pm 0.013	0.665 \pm 0.023	0.713 \pm 0.015	0.088 \pm 0.016	0.181 \pm 0.027	
LIMI	0.021\pm0.001	0.124\pm0.002	0.506\pm0.006	0.206\pm0.204	0.446\pm0.006	0.785 \pm 0.010	0.821 \pm 0.010	0.178 \pm 0.008	0.277 \pm 0.011	

表 5-6: 基因检测任务实验结果 (均值 \pm 标准差)。 \uparrow (\downarrow) 表示数值越大 (小) 性能越好, 加粗的数字表示最佳性能。

对比方法	5% 标注样本									
	Hamming Loss \downarrow	Ranking Loss \downarrow	One Error \downarrow	Coverage \downarrow	Ave. Precision \uparrow	Macro AUC \uparrow	Micro AUC \uparrow	Macro F1 \uparrow	Micro F1 \uparrow	
FCN	0.282 \pm 0.014	0.203 \pm 0.007	0.266 \pm 0.014	0.492 \pm 0.012	0.721 \pm 0.009	0.601 \pm 0.020	0.793 \pm 0.010	0.370 \pm 0.017	0.536 \pm 0.022	
CAMEL	0.281 \pm 0.012	0.220 \pm 0.011	0.291 \pm 0.023	0.519 \pm 0.012	0.709 \pm 0.013	0.602 \pm 0.017	0.781 \pm 0.011	0.381 \pm 0.012	0.539 \pm 0.018	
DBL	0.282 \pm 0.013	0.205 \pm 0.009	0.291 \pm 0.053	0.486 \pm 0.012	0.713 \pm 0.018	0.600 \pm 0.023	0.789 \pm 0.015	0.370 \pm 0.018	0.535 \pm 0.021	
DBL+NT	0.281 \pm 0.013	0.206 \pm 0.007	0.281 \pm 0.032	0.491 \pm 0.010	0.717 \pm 0.012	0.601 \pm 0.020	0.789 \pm 0.012	0.373 \pm 0.016	0.537 \pm 0.021	
PL	0.281 \pm 0.012	0.201 \pm 0.006	0.267 \pm 0.018	0.487 \pm 0.010	0.722 \pm 0.007	0.601 \pm 0.018	0.794 \pm 0.009	0.372 \pm 0.012	0.537 \pm 0.017	
DRML	0.298 \pm 0.014	0.215 \pm 0.009	0.301 \pm 0.022	0.507 \pm 0.014	0.708 \pm 0.012	0.596 \pm 0.012	0.779 \pm 0.009	0.387 \pm 0.013	0.529 \pm 0.015	
DRML+DBL	0.298 \pm 0.017	0.220 \pm 0.010	0.308 \pm 0.024	0.512 \pm 0.013	0.706 \pm 0.012	0.603\pm0.012	0.775 \pm 0.010	0.388\pm0.013	0.525 \pm 0.019	
LIMI	0.280\pm0.014	0.197\pm0.006	0.258\pm0.006	0.477\pm0.012	0.726\pm0.008	0.600 \pm 0.018	0.797\pm0.007	0.371 \pm 0.017	0.540\pm0.021	

对比方法	10% 标注样本									
	Hamming Loss \downarrow	Ranking Loss \downarrow	One Error \downarrow	Coverage \downarrow	Ave. Precision \uparrow	Macro AUC \uparrow	Micro AUC \uparrow	Macro F1 \uparrow	Micro F1 \uparrow	
FCN	0.264 \pm 0.009	0.192 \pm 0.004	0.259 \pm 0.010	0.477 \pm 0.004	0.734 \pm 0.004	0.626 \pm 0.007	0.809 \pm 0.005	0.396 \pm 0.013	0.562 \pm 0.015	
CAMEL	0.264 \pm 0.008	0.207 \pm 0.005	0.280 \pm 0.015	0.505 \pm 0.007	0.722 \pm 0.008	0.627 \pm 0.015	0.796 \pm 0.005	0.405 \pm 0.012	0.564\pm0.013	
DBL	0.264 \pm 0.007	0.192 \pm 0.004	0.268 \pm 0.018	0.471 \pm 0.007	0.734 \pm 0.005	0.627 \pm 0.010	0.807 \pm 0.004	0.394 \pm 0.011	0.561 \pm 0.013	
DBL+NT	0.263 \pm 0.008	0.192 \pm 0.005	0.266 \pm 0.017	0.471 \pm 0.006	0.733 \pm 0.005	0.625 \pm 0.010	0.807 \pm 0.004	0.396 \pm 0.011	0.562 \pm 0.014	
PL	0.263 \pm 0.008	0.191 \pm 0.004	0.260 \pm 0.011	0.475 \pm 0.005	0.734 \pm 0.005	0.626 \pm 0.009	0.809 \pm 0.004	0.394 \pm 0.009	0.561 \pm 0.013	
DRML	0.301 \pm 0.014	0.210 \pm 0.010	0.297 \pm 0.021	0.499 \pm 0.012	0.716 \pm 0.013	0.622 \pm 0.014	0.791 \pm 0.010	0.410\pm0.015	0.537 \pm 0.016	
DRML+DBL	0.286 \pm 0.014	0.214 \pm 0.007	0.303 \pm 0.015	0.505 \pm 0.007	0.713 \pm 0.010	0.625 \pm 0.011	0.786 \pm 0.008	0.402 \pm 0.012	0.539 \pm 0.019	
LIMI	0.261\pm0.007	0.186\pm0.003	0.249\pm0.003	0.462\pm0.006	0.739\pm0.003	0.629\pm0.007	0.811\pm0.003	0.396 \pm 0.006	0.564\pm0.010	

5.4.5 文本分类任务

文本分类任务也是多标记学习中一种常见的任务，对于文本分类任务，我们在 BibTex 数据集^①中进行实验。BibTex 数据集是从公开的社交分享平台收集得到的，共包含 7,395 个样本，每个样本的特征维度为 1,836，标记维度为 159，数据集的类别不平衡比例为 32.25。

BibTex 数据集上的实验结果如表 5-5 所示。从实验结果我们可以发现，最先进的多标记学习方法 CAMEL 在多个性能指标上表现良好，尤其是在 Macro/Micro F1 和 Macro/Micro AUC 中，但是本文提出的 LIMM 方法仍然在超过一半的性能指标上取得了最佳性能，这也证明了我们提出的 LIMM 方法的有效性。

5.4.6 基因检测任务

我们进一步在基因检测任务的基准数据集 Yeast Saccharomyces^②上进行实验。Yeast 数据集 [41] 包括 2,417 个基因样本，每个样本包括 103 维微阵列表达特征，对应的标记维度为 14，数据集的类别不平衡比例为 2.78。

Yeast 数据集上的实验结果如表 5-6 所示，与其它机器学习任务类似，我们的方法相比于对比方法也实现了最优的性能。综上，在工业界真实应用场景的网约车智能判责任务以及多种常用的多标记学习基准任务上，LIMM 方法都可以在多种评价指标上同时取得良好的性能，这些实验结果证明了 LIMM 方法对类别不平衡问题的稳健性。

5.5 小结

在本章中，我们研究开放环境下数据集类别比例不平衡时，如何保证半监督学习算法在所有类别上实现稳健的性能。这是半监督学习应用到开放环境需要解决的一个关键问题，并且已有的研究工作较少。针对单标记的半监督类别不平衡学习问题，我们提出基于 AUC 优化的 CWSL 方法，首先通过对样本进行赋权有效地利用无标注样本并降低伪标注错误的样本对模型性能的影响，并且进一步通过在外层优化中直接优化对类别比例不敏感的 AUC 指标，保证了算法

^①<http://mulan.sourceforge.net/datasets-mlc.html>

^②<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multilabel.html>

在类别不平衡数据上性能的稳健性。针对 AUC 指标的非凸非连续性，CWSL 采用替代损失函数进行优化，在理论上能够实现与 AUC 指标一致性的同时得到了高效的优化方法。针对多标记的半监督类别不平衡学习问题，我们提出系统性的解决框架 LIM1，LIM1 由三个核心模块组成：标注分离、相关性挖掘以及标注补全，通过这三个模块实现了对类别不平衡的处理、标记关系的挖掘以及无标注样本的利用。大量多标记学习任务上的实验结果表明 LIM1 方法可以在 9 种常用的多标记学习评价指标中均取得最优的性能表现。值得一提的是，本文提出的两种算法已在工业界真实的任务场景，网约车智能评价任务与网约车智能判责任务中成功落地应用。

在未来工作中，我们将进一步探索类别比例不平衡情况下半监督学习性能的理论结果，例如，类别失衡对半监督学习泛化性能的影响，类别失衡对半监督模型训练收敛性的影响等，为开放环境下分布失配的半监督学习建立理论基础。

本章的主要工作已经成文发表，包括：

- Lan-Zhe Guo, Feng Kuang, Zhang-Xun Liu, Yu-Feng Li, Nan Ma, Xiao-Hu Qie. Weakly Supervised Learning Meets Ride-Sharing User Experience Enhancement. In: **Proceedings of the 34rd AAAI conference on Artificial Intelligence (AAAI'20)**, New York, NY, USA, pp.4052-4059, 2020. (中国计算机学会 A 类会议，第一作者)
- Lan-Zhe Guo, Zhi Zhou, Jie-Jing Shao, Qi Zhang, Feng Kuang, Gao-Le Li, Zhang-Xun Liu, Guo-Bin Wu, Nan Ma, Qun Li, Yu-Feng Li. Learning from Imbalanced and Incomplete Supervision with Its Application to Ride-Sharing Liability Judgment. In: **Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD'21)**, Virtual Event, Singapore, pp.487-495, 2021. (中国计算机学会 A 类会议，第一作者)

第六章 结束语

本文工作主要涉及国家自然科学基金创新群体项目“面向开放动态环境的机器学习”(61921006)和面上项目“适于图像分类和标注的安全机器学习技术研究”(61772262)等科研项目内容。

6.1 本文工作总结

半监督学习通过引入大量易于获取的无标注数据辅助模型训练,提升标注信息不足时机器学习的泛化性能,推动了机器学习方法在更多场景中的应用。既有半监督学习研究大多针对封闭环境,然而,现实任务往往是开放的,存在数据分布失配、数据动态流式、先验知识不足、类别比例失衡等挑战,既有方法泛化性能时好时坏,难以适用于开放环境。本文围绕开放环境下的半监督学习展开研究,提出了一套更稳健适应开放环境的半监督学习解决方案,具体而言,本文主要取得了以下创新成果:

第二章针对数据分布失配,提出了稳健使用无标注样本的方法 DS3L,通过双层优化进行样本赋权,避免分布外样本导致泛化性能下降。理论上证明了该方法的收敛率,并从经验风险和泛化风险角度分析了其性能,显著提升了既有方法对分布失配稳健性。

第三章针对数据动态流式,提出了稳健适应动态数据的方法 Record,通过影响力机制进行子集选择以适应受限的资源 and 变化的分布,避免数据分布变化导致泛化性能下降。该方法可与任意半监督方法结合,一致有效提升既有方法的稳健性。

第四章针对先验知识不足,提出了稳健选择半监督模型的方法 SafeW,通过优化潜在最坏情况下的性能提升,避免模型选择错误导致泛化性能下降。理论上给出了半监督学习实现稳健性的条件,显著提升了既有方法的稳健性,并易于扩展以提升更宽泛的弱监督学习稳健性。

第五章针对类别比例失衡,提出了稳健处理少数类标注的方法 CWSL 和 LIM1,通过 AUC 优化和标注分离策略,避免少数类样本过少导致泛化性能下降。

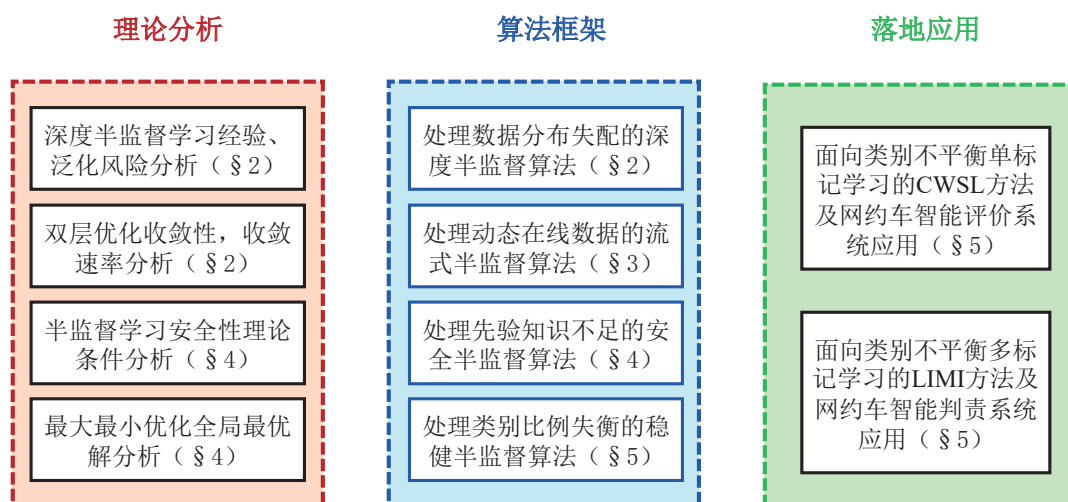


图 6-1: 从理论、算法、应用层面对本文工作的总结。

该方法在现实应用网约车智能评价和智能判责任务中成功落地转化, 显著提升了业界既有做法的稳健性。

图 6-1对本文的工作从理论、算法、应用的层面做了总结。理论上, 本文对深度半监督学习的经验风险及泛化风险进行了分析, 对提出的双层优化方法进行了收敛性及收敛速率的证明, 给出了半监督学习实现稳健性的理论条件, 对提出的最大最小优化目标进行了全局最优解的分析。算法上, 针对数据分布失配、数据动态流式、先验知识不足、类别比例失衡四种典型问题分别提出了稳健的半监督学习方法。应用上, 在滴滴出行网约车平台实际任务网约车智能评价系统和网约车智能判责系统中成功落地转化。

6.2 未来研究展望

本文考虑开放环境下的半监督学习, 并从数据输入-模型构建-标注输出三个角度针对性的提出解决方案, 后续考虑到更广泛的开放场景以及环境持续动态变化的情况, 有如下一些值得研究的方向:

1. 面向流数据的半监督学习。本文第三章中, 我们初步考虑了无标注数据流式到来且数据分布逐渐变化的情况, 在后续工作中, 我们将进一步考虑数据分布、类别空间、特征维度都有可能发生变化的情况, 设计稳健性强、复杂度低的流式半监督学习理论、方法和技术体系, 使半监督学习模型具备实时处理和适应变化的能力。

2. 利用环境交互的半监督学习。在开放动态环境中，往往存在大量复杂的领域知识和带噪的反馈信息。如何建立模型与环境的交互，针对交互中知识复杂与反馈带噪的问题，形成知识精化推理、带噪风险优化技术，降低对高质量监督信息的依赖，提升半监督学习的稳健性，是一个值得研究的问题。

3. 数据知识双驱动的半监督学习。除了无标注数据之外，逻辑知识也是很多任务中能够获得的信息。因此，在数据信息不足以训练可靠的半监督学习模型时，通过引入逻辑知识推理，在数据和知识的共同促进下提升半监督学习的性能是一个值得研究的问题。

4. 针对非结构化数据的深度半监督学习。现有深度半监督学习技术在图像、语音、文本等结构化数据中取得了成功的应用，但是在现实应用场景中，更为常见的数据类型是非结构化的表格型数据，现有深度半监督学习中常用的数据增广策略在这种非结构化数据中不再适用。因此，如何设计能够处理非结构化数据的深度半监督学习，是一个重要的研究方向。

此外，在合作者的帮助下，我们已经对开放环境下的半监督学习研究做了综述介绍 [58]，提出了多种半监督学习泛化性能下降的现实任务场景，并且开源了基于 Python 语言的半监督学习工具包及相应教程。未来我们会继续针对开放动态环境下的半监督学习进行研究和探索，完善理论体系、设计稳健方法、开源实用工具，进一步推动半监督学习的发展与应用。

参考文献

- [1] Stuart Andrews, Ioannis Tsochantaridis, and Thomas Hofmann. Support vector machines for multiple-instance learning. In *Advances in Neural Information Processing Systems*, pages 561–568, 2002.
- [2] A. Argyriou, A. Maurer, and M. Pontil. An algorithm for transfer learning in a heterogeneous environment. In *Proceedings of European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 71–85, 2008.
- [3] B. Bakker and T. Heskes. Task clustering and gating for bayesian multitask learning. *Journal of Machine Learning Research*, 4:83–99, 2003.
- [4] Akshay Balsubramani and Yoav Freund. Optimally combining classifiers using unlabeled data. In *Proceedings of the 28th Conference on Learning Theory*, pages 211–225, 2015.
- [5] Jonathan F Bard. *Practical bilevel optimization: algorithms and applications*. Springer Science and Business Media, 2013.
- [6] John M Bates and Clive WJ Granger. The combination of forecasts. *Operational Research*, 20(4):451–468, 1969.
- [7] David Berthelot, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Remixmatch: Semi-supervised learning with distribution matching and augmentation anchoring. In *Proceedings of the 8th International Conference on Learning Representations*, 2020.
- [8] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *Advances in Neural Information Processing Systems*, pages 5050–5060, 2019.
- [9] Albert Bifet, Geoffrey Holmes, Bernhard Pfahringer, and Ricard Gavaldà. Detecting sentiment change in twitter streaming data. In *Proceedings of the 2nd Workshop on Applications of Pattern Analysis*, pages 5–11, 2011.
- [10] Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the 11th Conference on Learning Theory*, pages 92–100, 1998.

- [11] Edwin V Bonilla, Kian M Chai, and Christopher Williams. Multi-task gaussian process prediction. In *Advances in Neural Information Processing Systems*, pages 153–160, 2008.
- [12] Jakramate Bootkrajang and Ata Kabán. Label-noise robust logistic regression and its applications. In *Proceedings of the European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 143–158, 2012.
- [13] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- [14] Forrest Briggs, Xiaoli Z Fern, and Raviv Raich. Rank-loss support instance machines for miml instance annotation. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 534–542, 2012.
- [15] Mateusz Buda, Atsuto Maki, and Maciej A Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 106:249–259, 2018.
- [16] Jonathon Byrd and Zachary Lipton. What is the effect of importance weighting in deep learning? In *Proceedings of the 36th International Conference on Machine Learning*, pages 872–881, 2019.
- [17] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Aréchiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In *Advances in Neural Information Processing Systems*, pages 1565–1576, 2019.
- [18] Marc-André Carbonneau, Veronika Cheplygina, Eric Granger, and Ghyslain Gagnon. Multiple instance learning: A survey of problem characteristics and applications. *Pattern Recognition*, 77:329–353, 2018.
- [19] Yair Censor and Stavros Andrea Zenios. *Parallel optimization: Theory, algorithms, and applications*. Oxford University Press, 1997.
- [20] Hakan Cevikalp, Burak Benligiray, and Ömer Nezih Gerek. Semi-supervised robust deep neural networks for multi-label image classification. *Pattern Recognition*, 100:107164, 2020.
- [21] Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2(3):27, 2011.

- [22] Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien. *Semi-supervised learning*. MIT Press, 2006.
- [23] Francisco Charte, Antonio J. Rivera, María José del Jesus, and Francisco Herrera. Mlenn: A first approach to heuristic multilabel undersampling. In *Proceedings of the 15th International Conference on Intelligent Data Engineering and Automated Learning*, pages 1–9, 2014.
- [24] Francisco Charte, Antonio J. Rivera, María José del Jesus, and Francisco Herrera. Addressing imbalance in multilabel classification: Measures and random resampling algorithms. *Neurocomputing*, 163:3–16, 2015.
- [25] Francisco Charte, Antonio J. Rivera, María José del Jesus, and Francisco Herrera. MLSMOTE: approaching imbalanced multilabel learning through synthetic instance generation. *Knowledge Based System*, 89:385–397, 2015.
- [26] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16:321–357, 2002.
- [27] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794, 2016.
- [28] Yanbei Chen, Xiatian Zhu, Wei Li, and Shaogang Gong. Semi-supervised learning under class distribution mismatch. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, pages 3569–3576, 2020.
- [29] Fabio Gagliardi Cozman, Ira Cohen, and M Cirelo. Unlabeled data can degrade classification performance of generative classifiers. In *Proceedings of the 15th International Conference of the Florida Artificial Intelligence Research Society*, pages 327–331, 2002.
- [30] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9268–9277, 2019.
- [31] Wenyuan Dai, Qiang Yang, Gui-Rong Xue, and Yong Yu. Boosting for transfer learning. In *Proceedings of the 24th International Conference on Machine Learning*, pages 193–200, 2007.

- [32] Andrea Dal Pozzolo, Olivier Caelen, Reid A Johnson, and Gianluca Bontempi. Calibrating probability with undersampling for unbalanced classification. In *IEEE Symposium Series on Computational Intelligence*, pages 159–166, 2015.
- [33] Zachary Alan Daniels and Dimitris N. Metaxas. Addressing imbalance in multi-label classification using structured hellinger forests. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence*, pages 1826–1832, 2017.
- [34] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- [35] Thomas G Dietterich. An experimental comparison of three methods for constructing ensembles of decision trees: Bagging, boosting, and randomization. *Machine Learning*, 40(2):139–157, 2000.
- [36] Thomas G. Dietterich. Steps toward robust artificial intelligence. *AI Magazine*, 38(3):3–24, 2017.
- [37] Gregory Ditzler and Robi Polikar. Semi-supervised learning in nonstationary environments. In *Proceedings of the International Joint Conference on Neural Networks*, pages 2741–2748, 2011.
- [38] John C. Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- [39] Karl B Dyer, Robert Capo, and Robi Polikar. Compose: A semi-supervised learning framework for initially labeled nonstationary streaming data. *IEEE Transactions on Neural Networks and Learning Systems*, 25(1):12–26, 2013.
- [40] Sandra Ebert, Mario Fritz, and Bernt Schiele. Semi-supervised learning on a budget: Scaling up to large datasets. In *Proceedings of the 11th Asian Conference on Computer Vision*, pages 232–245, 2012.
- [41] André Elisseeff and Jason Weston. A kernel method for multi-labelled classification. In *Advances in Neural Information Processing Systems*, pages 681–687, 2002.
- [42] Linwei Fan, Xuemei Li, Qiang Guo, and Caiming Zhang. Nonlocal image denoising using edge-based similarity metric and adaptive parameter selection. *SCIENCE CHINA Information Sciences*, 61(4):253–261, 2018.

- [43] Lei Feng, Bo An, and Shuo He. Collaboration based multi-label learning. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*, pages 3550–3557, 2019.
- [44] Benoît Frénay and Michel Verleysen. Classification in the presence of label noise: a survey. *IEEE Transactions on Neural Networks and Learning Systems*, pages 845–869, 2014.
- [45] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- [46] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of Statistics*, pages 1189–1232, 2001.
- [47] Anil Gaba and Robert L Winkler. Implications of errors in survey data: a bayesian model. *Management Science*, 38(7):913–925, 1992.
- [48] Wei Gao, Rong Jin, Shenghuo Zhu, and Zhi-Hua Zhou. One-pass auc optimization. In *Proceedings of the 30th International Conference on Machine Learning*, pages 906–914, 2013.
- [49] Wei Gao and Zhi-Hua Zhou. On the consistency of AUC pairwise optimization. In *Proceedings of the 24th International Joint Conference on Artificial Intelligence*, pages 939–945, 2015.
- [50] Liang Ge, Jing Gao, Hung Ngo, Kang Li, and Aidong Zhang. On handling negative transfer and imbalanced distributions in multiple source transfer learning. *Statistical Analysis and Data Mining*, 7(4):254–271, 2014.
- [51] Andrew B Goldberg, Ming Li, and Xiaojin Zhu. Online manifold regularization: A new learning setting and empirical study. In *Proceedings of the European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 393–407, 2008.
- [52] Andrew B Goldberg, Xiaojin Zhu, Alex Furger, and Jun-Ming Xu. Oasis: Online active semi-supervised learning. In *Proceedings of the 25th AAAI Conference on Artificial Intelligence*, pages 362–367, 2011.
- [53] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT Press, 2016.

- [54] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014.
- [55] Yves Grandvalet and Yoshua Bengio. Semi-supervised learning by entropy minimization. In *Advances in Neural Information Processing Systems*, pages 529–536, 2005.
- [56] Will Grathwohl, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. In *Proceedings of the 8th International Conference on Learning Representations*, pages 1–23, 2020.
- [57] Lan-Zhe Guo, Feng Kuang, Zhang-Xun Liu, Yu-Feng Li, Nan Ma, and Xiao-Hu Qie. Iwe-net: Instance weight network for locating negative comments and its application to improve traffic user experience. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, pages 4052–4059, 2020.
- [58] Lan-Zhe Guo, Zhi Zhou, and Yu-Feng Li. Robust deep semi-supervised learning: A brief introduction. *CoRR*, abs/2202.05975, 2022.
- [59] Ahsanul Haque, Latifur Khan, and Michael Baron. Sand: Semi-supervised adaptive novel class detection and classification over data stream. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence*, pages 1652–1658, 2016.
- [60] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2001.
- [61] Haibo He and Edwardo A Garcia. Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 21(9):1263–1284, 2009.
- [62] Xinran He, Junfeng Pan, Ou Jin, Tianbing Xu, Bo Liu, Tao Xu, Yanxin Shi, Antoine Atallah, Ralf Herbrich, Stuart Bowers, et al. Practical lessons from predicting clicks on ads at facebook. In *Proceedings of the 8th International Workshop on Data Mining for Online Advertising*, pages 1–9, 2014.
- [63] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *Proceedings of the 5th International Conference on Learning Representations*, pages 1–12, 2017.

- [64] Dan Hendrycks, Mantas Mazeika, Duncan Wilson, and Kevin Gimpel. Using trusted data to train deep networks on labels corrupted by severe noise. In *Advances in Neural Information Processing Systems*, pages 10456–10465, 2018.
- [65] Ray J Hickey. Noise modelling and evaluating learning from examples. *Artificial Intelligence*, 82(1-2):157–179, 1996.
- [66] Zhiting Hu, Bowen Tan, Ruslan Salakhutdinov, Tom M. Mitchell, and Eric P. Xing. Learning data manipulation for augmentation and weighting. In *Advances in Neural Information Processing Systems*, pages 15738–15749, 2019.
- [67] Chen Huang, Yining Li, Chen Change Loy, and Xiaoou Tang. Deep imbalanced learning for face recognition and attribute prediction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(11):2781–2794, 2019.
- [68] Lei Huang, Xianglong Liu, Binqiang Ma, and Bo Lang. Online semi-supervised annotation via proxy-based local consistency propagation. *Neurocomputing*, 149:1573–1586, 2015.
- [69] Jisoo Jeong, Seungeui Lee, Jeessoo Kim, and Nojun Kwak. Consistency-based semi-supervised learning for object detection. In *Advances in Neural Information Processing Systems*, pages 10759–10768, 2019.
- [70] Thorsten Joachims. Transductive inference for text classification using support vector machines. In *Proceedings of the 16th International Conference on Machine Learning*, pages 200–209, 1999.
- [71] Justin M Johnson and Taghi M Khoshgoftaar. Survey on deep learning with class imbalance. *Journal of Big Data*, 6(1):1–54, 2019.
- [72] Salman Khan, Munawar Hayat, Syed Waqas Zamir, Jianbing Shen, and Ling Shao. Striking the right balance with uncertainty. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 103–112, 2019.
- [73] Salman H Khan, Munawar Hayat, Mohammed Bennamoun, Ferdous A Sohel, and Roberto Togneri. Cost-sensitive learning of deep feature representations from imbalanced data. *IEEE Transactions on Neural Networks and Learning Systems*, 29(8):3573–3587, 2017.
- [74] Yoon Kim. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, pages 1746–1751, 2014.

- [75] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations*, 2015.
- [76] Josef Kittler and Cemre Zor. Delta divergence: A novel decision cognizant measure of classifier incongruence. *IEEE Transactions on Cybernetics*, 49(6):2331–2343, 2018.
- [77] Josef Kittler, Cemre Zor, Ioannis Kaloskampis, Yulia Hicks, and Wenwu Wang. Error sensitivity analysis of delta divergence—a novel measure for classifier incongruence detection. *Pattern Recognition*, 77:30–44, 2018.
- [78] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning*, pages 1885–1894, 2017.
- [79] Georg Kreml, Indre Zliobaite, Dariusz Brzeziński, Eyke Hüllermeier, Mark Last, Vincent Lemaire, Tino Noack, Ammar Shaker, Sonja Sievi, Myra Spiliopoulou, et al. Open challenges for data stream mining research. *SIGKDD Explorations*, 16(1):1–10, 2014.
- [80] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. In *Proceedings of the 5th International Conference on Learning Representations*, 2017.
- [81] Ken Lang. Newsweeder: Learning to filter netnews. In *Proceedings of the 12th International Conference on Machine Learning*, pages 331–339, 1995.
- [82] Trung Le, Phuong Duong, Mi Dinh, Tu Dinh Nguyen, Vu Nguyen, and Dinh Q Phung. Budgeted semi-supervised support vector machine. In *Proceedings of the 32nd Conference on Uncertainty in Artificial Intelligence*, pages 377–386, 2016.
- [83] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- [84] Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Proceedings of the 30th International Conference on International Conference on Machine Learning Workshop*, page 896, 2013.

- [85] Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. In *Proceedings of the 6th International Conference on Learning Representations*, pages 1–16, 2018.
- [86] Joffrey L. Leevy, Taghi M. Khoshgoftaar, Richard A. Bauder, and Naeem Seliya. A survey on addressing high-class imbalance in big data. *Journal of Big Data*, 5:42, 2018.
- [87] Lianghao Li, Xiaoming Jin, and Mingsheng Long. Topic correlation analysis for cross-domain text classification. In *Proceedings of the 26th AAAI Conference on Artificial Intelligence*, pages 998–1004, 2012.
- [88] Pei-Pei Li, Xindong Wu, and Xuegang Hu. Mining recurring concept drifts with limited labeled streaming data. In *Proceedings of the 2nd Asian Conference on Machine Learning*, pages 241–252, 2010.
- [89] Yu-Feng Li, James T Kwok, and Zhi-Hua Zhou. Towards safe semi-supervised learning for multivariate performance measures. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence*, pages 1816–1822, 2016.
- [90] Yu-Feng Li, Shao-Bo Wang, and Zhi-Hua Zhou. Graph quality judgement: A large margin expedition. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence*, pages 1725–1731, 2016.
- [91] Yu-Feng Li, Han-Wen Zha, and Zhi-Hua Zhou. Learning safe prediction for semi-supervised regression. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence*, pages 2217–2223, 2017.
- [92] Yu-Feng Li and Zhi-Hua Zhou. Towards making unlabeled data never hurt. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(1):175–188, 2015.
- [93] Shiyu Liang, Yixuan Li, and R Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *Proceedings of the 6th International Conference on Learning Representations*, pages 1–15, 2018.
- [94] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2980–2988, 2017.

- [95] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. In *Advances in Neural Information Processing Systems*, pages 21464–21475, 2020.
- [96] Xu-Ying Liu, Jianxin Wu, and Zhi-Hua Zhou. Exploratory undersampling for class-imbalance learning. *IEEE Transactions on Systems, Man, and Cybernetics*, 39(2):539–550, 2008.
- [97] Yi Liu, Rong Jin, and Liu Yang. Semi-supervised multi-label learning by constrained non-negative matrix factorization. In *Proceedings of the 21st AAAI Conference on Artificial Intelligence*, pages 421–426, 2006.
- [98] Marco Loog. Contrastive pessimistic likelihood estimation for semi-supervised classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(3):462–475, 2015.
- [99] Naresh Manwani and PS Sastry. Noise tolerance under risk minimization. *IEEE Transactions on Cybernetics*, 43(3):1146–1151, 2013.
- [100] H Brendan McMahan, Gary Holt, David Sculley, Michael Young, Dietmar Ebner, Julian Grady, Lan Nie, Todd Phillips, Eugene Davydov, Daniel Golovin, et al. Ad click prediction: a view from the trenches. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1222–1230, 2013.
- [101] Lilyana Mihalkova, Tuyen Huynh, and Raymond J Mooney. Mapping and revising markov logic networks for transfer learning. In *Proceedings of the 22nd AAAI Conference on Artificial Intelligence*, pages 608–614, 2007.
- [102] David J Miller and Hasan S Uyar. A mixture of experts classifier with learning based on both labelled and unlabelled data. In *Advances in Neural Information Processing Systems*, pages 571–577, 1997.
- [103] Tom Mitchell. *Machine learning*. McGraw-Hill, 1997.
- [104] Takeru Miyato, Andrew M. Dai, and Ian J. Goodfellow. Adversarial training methods for semi-supervised text classification. In *Proceedings of the 5th International Conference on Learning Representations*, 2017.
- [105] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: A regularization method for supervised and semi-supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(8):1979–1993, 2018.

- [106] Kazuhisa Nakasho, Yuichi Futa, and Yasunari Shidama. Implicit function theorem. part I. *Formalized Mathematics*, pages 269–281, 2017.
- [107] Kamal Nigam, Andrew Kachites McCallum, Sebastian Thrun, and Tom Mitchell. Text classification from labeled and unlabeled documents using em. *Machine Learning*, 39(2):103–134, 2000.
- [108] Curtis G. Northcutt, Tailin Wu, and Isaac L. Chuang. Learning with confident examples: Rank pruning for robust classification with noisy labels. In *Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence*, 2017.
- [109] Avital Oliver, Augustus Odena, Colin A Raffel, Ekin Dogus Cubuk, and Ian Goodfellow. Realistic evaluation of deep semi-supervised learning algorithms. In *Advances in Neural Information Processing Systems*, pages 3235–3246, 2018.
- [110] Sinno Jialin Pan, Ivor W Tsang, James T Kwok, and Qiang Yang. Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks*, 22(2):199–210, 2011.
- [111] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10):1345–1359, 2010.
- [112] Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural Networks*, 12(1):145–151, 1999.
- [113] Rouhollah Rahmani, Sally A Goldman, Hui Zhang, John Krettek, and Jason E Fritts. Localized content based image retrieval. In *Proceedings of the ACM SIGMM International Workshop on Multimedia Information Retrieval*, pages 227–236, 2005.
- [114] Rajat Raina, Alexis Battle, Honglak Lee, Benjamin Packer, and Andrew Ng. Self-taught learning: Transfer learning from unlabeled data. In *Proceedings of the 24th International Conference on Machine Learning*, pages 759–766, 2007.
- [115] Soumya Ray and Mark Craven. Supervised versus multiple instance learning: An empirical comparison. In *Proceedings of the 22nd International Conference on Machine Learning*, pages 697–704, 2005.
- [116] Mengye Ren, Wenyan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *Proceedings of the 35th International Conference on Machine Learning*, pages 4331–4340, 2018.

- [117] M. T. Rosenstein, Z. Marx, and L. P. Kaelbling. To transfer or not to transfer. In *NIPS Workshop on “Inductive Transfer: 10 Years Later”*, 2005.
- [118] Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In *Advances in Neural Information Processing Systems*, pages 1163–1171, 2016.
- [119] Behzad M Shahshahani and David A Landgrebe. The effect of unlabeled samples in reducing the small sample size problem and mitigating the Hughes phenomenon. *IEEE Transactions on Geoscience and remote sensing*, 32(5):1087–1095, 1994.
- [120] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge University Press, 2014.
- [121] Ira Shavitt and Eran Segal. Regularization learning networks: Deep learning for tabular datasets. In *Advances in Neural Information Processing Systems*, pages 1386–1396, 2018.
- [122] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proceedings of the 3rd International Conference on Learning Representations*, 2015.
- [123] Ankur Sinha, Pekka Malo, and Kalyanmoy Deb. A review on bilevel optimization: from classical to evolutionary approaches and applications. *IEEE Transactions on Evolutionary Computation*, 22(2):276–295, 2018.
- [124] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fix-match: Simplifying semi-supervised learning with consistency and confidence. In *Advances in Neural Information Processing Systems*, pages 596–608, 2020.
- [125] Nasim Souly, Concetto Spampinato, and Mubarak Shah. Semi supervised semantic segmentation using generative adversarial network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5688–5696, 2017.
- [126] V. M. A. Souza, D. F. Silva, J. Gama, and G. E. A. P. A. Batista. Data stream classification guided by clustering on nonstationary environments and extreme verification latency. In *Proceedings of the SIAM International Conference on Data Mining*, pages 873–881, 2015.

- [127] Qiaoyu Tan, Yanming Yu, Guoxian Yu, and Jun Wang. Semi-supervised multi-label classification using incomplete label information. *Neurocomputing*, 260:192–202, 2017.
- [128] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *Advances in Neural Information Processing Systems*, pages 1195–1204, 2017.
- [129] Apoorv Vyas, Nataraj Jammalamadaka, Xia Zhu, Dipankar Das, Bharat Kaul, and Theodore L Willke. Out-of-distribution detection using an ensemble of self supervised leave-out classifiers. In *Proceedings of the European Conference on Computer Vision*, pages 550–564, 2018.
- [130] Tal Wagner, Sudipto Guha, Shiva Kasiviswanathan, and Nina Mishra. Semi-supervised learning on data streams via temporal label propagation. In *Proceedings of the 35th International Conference on Machine Learning*, pages 5095–5104, 2018.
- [131] Jun Wang and Jean-Daniel Zucker. Solving the multiple-instance problem: A lazy learning approach. In *Proceedings of the 17th International Conference on Machine Learning*, pages 1119–1126, 2000.
- [132] Lichen Wang, Yunyu Liu, Can Qin, Gan Sun, and Yun Fu. Dual relation semi-supervised multi-label learning. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, pages 6227–6234, 2020.
- [133] C.J Willmott and K. Matsuura. Advantages of the mean absolute error (mae) over the root mean square error (rmse) in assessing average model performance. *Climate Research*, 30(1):79–82, 2005.
- [134] Le Wu and Min-Ling Zhang. Multi-label classification with unlabeled data: An inductive approach. In *Proceedings of the 5th Asian Conference on Machine Learning*, pages 197–212, 2013.
- [135] Tong Wu, Qingqiu Huang, Ziwei Liu, Yu Wang, and Dahua Lin. Distribution-balanced loss for multi-label classification in long-tailed datasets. In *Proceedings of the 16th European Conference on Computer Vision*, pages 162–178, 2020.
- [136] Xi-Zhu Wu and Zhi-Hua Zhou. A unified view of multi-label performance measures. In *Proceedings of the 34th International Conference on Machine Learning*, pages 3780–3788, 2017.

- [137] Qizhe Xie, Zihang Dai, Eduard H. Hovy, Thang Luong, and Quoc Le. Un-supervised data augmentation for consistency training. In *Advances in Neural Information Processing Systems*, pages 6256–6268, 2020.
- [138] Yuying Xing, Guoxian Yu, Carlotta Domeniconi, Jun Wang, and Zili Zhang. Multi-label co-training. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 2882–2888, 2018.
- [139] Xin Xu and Eibe Frank. Logistic regression and boosting for labeled bags of instances. In *Proceedings of the 8th Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 272–281, 2004.
- [140] Gui-Rong Xue, Wenyuan Dai, Qiang Yang, and Yong Yu. Topic-bridged pls for cross-domain text classification. In *Proceedings of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 627–634, 2008.
- [141] Ke Yan, Lu Kou, and David Zhang. Domain adaptation via maximum independence of domain features. *arXiv preprint arXiv:1603.04535*, 2016.
- [142] Haiqin Yang, Kaizhu Huang, Irwin King, and Michael R Lyu. Maximum margin semi-supervised learning with irrelevant data. *Neural Networks*, 70:90–102, 2015.
- [143] Haiqin Yang, Shenghuo Zhu, Irwin King, and Michael R Lyu. Can irrelevant data help semi-supervised learning, why and how? In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management*, pages 937–946, 2011.
- [144] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. *CoRR*, abs/2110.11334, 2021.
- [145] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alexander J. Smola, and Eduard H. Hovy. Hierarchical attention networks for document classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1480–1489, 2016.
- [146] David Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. In *Proceedings of the 33rd Annual Meeting on Association for Computational Linguistics*, pages 189–196, 1995.

- [147] Qing Yu and Kiyoharu Aizawa. Unsupervised out-of-distribution detection by maximum classifier discrepancy. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9517–9525, 2019.
- [148] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *Proceedings of the British Machine Vision Conference*, 2016.
- [149] Wang Zhan and Min-Ling Zhang. Inductive semi-supervised multi-label learning with co-training. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1305–1314, 2017.
- [150] Cha Zhang, John C Platt, and Paul A Viola. Multiple instance boosting for object detection. In *Advances in Neural Information Processing Systems*, pages 1417–1424, 2006.
- [151] Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. Mixup: Beyond empirical risk minimization. In *Proceedings of the 6th International Conference on Learning Representations*, 2018.
- [152] Min-Ling Zhang, Yu-Kun Li, and Xu-Ying Liu. Towards class-imbalance aware multi-label learning. In *Proceedings of the 24th International Joint Conference on Artificial Intelligence*, pages 4041–4047, 2015.
- [153] Min-Ling Zhang and Zhi-Hua Zhou. A review on multi-label learning algorithms. *IEEE Transactions on Knowledge and Data Engineering*, 26(8):1819–1837, 2013.
- [154] Qi Zhang and Sally A. Goldman. EM-DD: an improved multiple-instance learning technique. In *Advances in Neural Information Processing Systems*, pages 1073–1080, 2001.
- [155] Peng Zhao, Xinqiang Wang, Siyu Xie, Lei Guo, and Zhi-Hua Zhou. Distribution-free one-pass learning. *IEEE Transactions on Knowledge and Data Engineering*, 3(3):951–963, 2019.
- [156] Zhi-Hua Zhou. *Ensemble methods: Foundations and algorithms*. CRC Press, 2012.
- [157] Zhi-Hua Zhou. A brief introduction to weakly supervised learning. *National Science Review*, 5(1):44–53, 2017.

- [158] Zhi-Hua Zhou and Ji Feng. Deep forest: Towards an alternative to deep neural networks. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 3553–3559, 2017.
- [159] Zhi-Hua Zhou and Ming Li. Semi-supervised regression with co-training. In *Proceedings of the 19th International Joint Conference on Artificial Intelligence*, pages 908–913, 2005.
- [160] Zhi-Hua Zhou, Michael Ng, Qiao-Qiao She, and Yuan Jiang. Budget semi-supervised learning. In *Proceedings of the 13rd Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 588–595, 2009.
- [161] Zhi-Hua Zhou, Yu-Yin Sun, and Yu-Feng Li. Multi-instance learning by treating instances as non-iid samples. In *Proceedings of the 26th International Conference on Machine Learning*, pages 1249–1256, 2009.
- [162] Zhi-Hua Zhou, Xiao-Bing Xue, and Yuan Jiang. Locating regions of interest in cbr with multi-instance learning techniques. In *Proceedings of the 18th Australasian Joint Conference on Artificial Intelligence*, pages 92–101, 2005.
- [163] Zhi-Hua Zhou and Min-Ling Zhang. Solving multi-instance problems with classifier ensemble based on constructive clustering. *Knowledge and Information Systems*, 11(2):155–170, 2007.
- [164] Xiaojin Zhu, Zoubin Ghahramani, and John D. Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *Proceedings of 20th International Conference on Machine Learning*, pages 912–919, 2003.
- [165] Yong-Nan Zhu and Yu-Feng Li. Semi-supervised streaming learning with emerging new labels. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, pages 7015–7022, 2020.
- [166] 周志华. 机器学习与数据挖掘. 中国计算机学会通讯, 3(12):35–44, 2007.
- [167] 周志华. 基于分歧的半监督学习. 自动化学报, 39(11):1871–1878, 2013.
- [168] 周志华. 机器学习. 清华大学出版社, 2016.
- [169] 周志华, 王魏, 高尉, 张利军. 机器学习理论导引. 机械工业出版社, 2020.
- [170] 徐文华, 覃征, 常扬. 基于半监督学习的数据流集成分类算法. 模式识别与人工智能, 25(2):8, 2012.

致 谢

行文至此，思绪万千。研究生时光转瞬即逝，回首五年的博士求学光阴，目之所及，皆是感恩，今借此文，聊表谢意。

时穷节乃见，一一垂丹青。求学期间不幸遇到了全球新冠疫情，感谢国家、南京大学在疫情流行期间为我提供了一个和平、安全的环境，让我能够顺利的完成学业，感谢所有为此付出、牺牲的人们。

饮水思其源，学成念吾师。衷心感谢我的导师李宇峰老师，李老师在做学术和为人处事方面都是我学习的榜样。在学术上，李老师经常帮我反复修改论文，完善方法，指导我逐渐建立起做科研的方法论，李老师也注重培养我做科研的价值观，他经常让我们思考什么问题是值得研究的，什么工作是有价值的，提升了我做学术的品味。在为人处事上，李老师谦虚严谨，考虑周全，追求尽善尽美，为学生树立了良好的榜样。和李老师一起工作的这些年受益良多，有许多尚没有学到的东西，需要在未来进一步去体会。此外，也要感谢 LAMDA 提供的丰富的学习资源和浓厚的科研氛围，感谢 LAMDA 的每一位老师，你们都是我学术和人生道路上的榜样。

学贵在得师，亦贵在得友。感谢博士求学期间所有帮助过我的朋友。感谢魏通师兄，作为同门的博士大师兄，他在科研和生活上给了我很多的帮助。感谢所有合作过的优秀的同学们，包括周值、邵杰晶、韩韬、张震宇等，和各位讨论工作的时间让我们一起共同成长。

父母恩情重，三生报答轻。感谢我的父母，他们提供最好的一切供我求学，一直以来尊重我的选择，支持我对学术的追求，我会继续努力成为你们的骄傲。

同声若鼓瑟，合韵似鸣琴。感谢我的妻子张昕女士，我们已经相濡以沫走过了十年，一直相互依靠，共同前进。

最后感谢自己，没有过人的天赋，依然在求学之路上坚持了下来。

祝愿师长工作顺利，朋友前程似锦，亲人身体健康！

谨以此文纪念我的学生时代。

附录 A 攻读博士学位期间学术成果及 参加科研项目情况

一、攻读博士学位期间撰写的主要论文

中国计算机学会 A 类论文：

1. Lan-Zhe Guo, Yu-Feng Li. Class Imbalanced Semi-Supervised Learning with Adaptive Thresholding. In: **Proceedings of the 39th International Conference on Machine Learning (ICML'22)**, Virtual Event, In Press, 2022. (中国计算机学会 A 类会议, 第一作者)
2. Yu-Feng Li, Lan-Zhe Guo, Zhi-Hua Zhou. Towards Safe Weakly Supervised Learning. **IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)**, 43(1):334-346, 2021. (中国计算机学会 A 类期刊, 导师外一作)
3. Lan-Zhe Guo, Zhi Zhou, Jie-Jing Shao, Qi Zhang, Feng Kuang, Gao-Le Li, Zhang-Xun Liu, Guo-Bin Wu, Nan Ma, Qun Li, Yu-Feng Li. Learning from Imbalanced and Incomplete Supervision with Its Application to Ride-Sharing Liability Judgment. In: **Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD'21)**, Virtual Event, Singapore, pp.487-495, 2021. (中国计算机学会 A 类会议, 第一作者)
4. Zhi Zhou, Lan-Zhe Guo (co-first author), Zhan-Zhan Cheng, Yu-Feng Li, Shi-Liang Pu. STEP: Out-of-Distribution Detection in the Presence of Limited In-Distribution Labeled Data. In: **Advances in Neural Information Processing Systems (NeurIPS'21)**, Virtual Event, pp.29168-29180, 2021. (中国计算机学会 A 类会议, 共同一作)
5. Lan-Zhe Guo, Zhen-Yu Zhang, Yuan Jiang, Yu-Feng Li, Zhi-Hua Zhou. Safe Deep Semi-Supervised Learning for Unseen-Class Unlabeled Data. In: **Proceedings of the 37th International Conference on Machine Learning (ICML'20)**, Virtual

- Event, pp.3897-3906, 2020. (中国计算机学会 A 类会议, 第一作者)
6. Lan-Zhe Guo, Zhi Zhou, Yu-Feng Li. RECORD: Resource Constrained Semi-Supervised Learning under Distribution Shift. In: **Proceedings of the 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD'20)**, San Diego, CA, USA, pp.1636-1644, 2020. (中国计算机学会 A 类会议, 第一作者)
 7. Lan-Zhe Guo, Feng Kuang, Zhang-Xun Liu, Yu-Feng Li, Nan Ma, Xiao-Hu Qie. Weakly Supervised Learning Meets Ride-Sharing User Experience Enhancement. In: **Proceedings of the 34rd AAAI conference on Artificial Intelligence (AAAI'20)**, New York, NY, USA, pp.4052-4059, 2020. (中国计算机学会 A 类会议, 第一作者)
 8. Lan-Zhe Guo, Yu-Feng Li. A General Formulation for Safely Exploiting Weakly Supervised Data. In: **Proceedings of the 32nd AAAI conference on Artificial Intelligence (AAAI'18)**, New Orleans, LA, pp.3126-3133, 2018. (中国计算机学会 A 类会议, 第一作者)
 9. Chang-Jian Chen, Zhao-Wei Wang, Jing Wu, Xi-Ting Wang, Lan-Zhe Guo, Yu-Feng Li, Shi-Xia Liu. Interactive Graph Construction for Graph-Based Semi-Supervised Learning. **IEEE Transactions on Visualization and Computer Graphics (TVCG)**, 27(9): 3701-3716, 2021. (中国计算机学会 A 类期刊, 第五作者)

中国计算机学会 B、C 类论文:

1. Miao Xu, Lan-Zhe Guo. Learning from Group Supervision: The Impact of Supervision Deficiency on Multi-Label Learning. **Science China Information Science (SCIS)**, 64(3):1-13, 2021. (《中国科学: 信息科学》英文版, “中国科技期刊卓越行动计划” 重点期刊, 第二作者, 通讯作者)
2. Tong Wei, Lan-Zhe Guo, Yu-Feng Li, Wei Gao. Learning Safe Multi-Label Prediction for Weakly Labeled Data. **Machine Learning Journal (MLJ)**, 107(4):703-725, 2018. (机器学习领域重要期刊, 第二作者)
3. Lan-Zhe Guo, Tao Han, Yu-Feng Li. Robust Semi-Supervised Representation Learning for Graph-Structured Data. In: **Proceedings of the 23rd Pacific-Asia**

Conference on Knowledge Discovery and Data Mining (PAKDD'19), Macau, China, pages:131-143, 2019. (数据挖掘领域重要会议, 第一作者)

4. Lan-Zhe Guo, Shao-Bo Wang, Yu-Feng Li. Large Margin Graph Construction for Semi-Supervised Learning. In: **IEEE International Conference on Data Mining Workshops (ICDMW)**, Singapore, pages:1030-1033, 2018. (数据挖掘领域重要会议, 第一作者)

在准备论文:

1. SU-SSL: Maximize Performance on Unseen Classes and Maintain Safeness on Seen Classes. (投稿至 NeurIPS 2022, 中国计算机学会 A 类会议, 共同一作)
2. Active Model Adaptation Under Changed Environments. (投稿至 NeurIPS 2022, 中国计算机学会 A 类会议, 第二作者)
3. Robust Deep Semi-Supervised Learning: A Brief Introduction. (投稿至 IEEE Transactions on Pattern Analysis and Machine Intelligence, 中国计算机学会 A 类期刊, 第一作者)

二、攻读博士学位期间申请的专利

1. 李宇峰, 郭兰哲。一种安全可靠的图像分类半监督机器学习方法及装置, 发明专利, CN201910565453.1。

三、攻读博士学位期间参加的主要科研项目

1. 国家自然科学基金创新群体项目“面向开放动态环境的机器学习”(61921006)
2. 国家自然科学基金面上项目“适于图像分类和标注的安全机器学习技术研究”(61772262)
3. 科技部“云计算和大数据”重点研发计划项目“智能无人集群系统全局规划及协同行为管控”(2017YFB1001903)

附录 B 攻读博士学位期间获奖及学术 活动情况

一、攻读博士学位期间获奖情况

1. 第九届百度奖学金（奖金 20 万元，全球华人共 10 名），2021 年
2. 第一届华为-南京大学人工智能联合实验室突出贡献奖，2021 年
3. 微软学者提名（大陆地区共 16 人），2021 年
4. 腾讯奖学金，2021 年
5. 博士研究生国家奖学金，2020 年
6. 博士研究生国睿奖学金，2019 年
7. PAKDD 学生旅行奖（Student Travel Award），2019 年
8. 硕士研究生国家奖学金，2018 年
9. AAAI 学生旅行奖（Student Travel Award），2018 年

二、攻读博士学位期间学术活动情况

国内外重要学术会议组织委员会成员：

1. 第 13 届亚洲机器学习会议流程主席（The 13rd Asian Conference on Machine Learning, ACML 2021, Workflow Chair, 国际人工智能重要会议）
2. 第 18 届中国机器学习会议流程主席（The 18th China Conference on Machine Learning, CCML 2021, Workflow chair, 国内人工智能重要会议）

顶级国际学术会议高级程序委员会成员：

1. 第 30 届国际人工智能联合大会（The 30th International Joint Conference on Artificial Intelligence, IJCAI 2021, 人工智能领域顶级会议，CCF-A 类会议）
2. 第 13 届亚洲机器学习会议（The 13rd Asian Conference on Machine Learning, ACML 2021, 人工智能领域重要会议）

顶级国际学术会议程序委员会成员:

1. 国际机器学习大会 (International Conference on Machine Learning, ICML, CCF-A 类会议): 2020, 2021, 2022
2. 神经信息处理系统大会 (Neural Information Processing Systems, NeurIPS, CCF-A 类会议): 2020, 2021, 2022
3. AAAI 人工智能大会 (AAAI Conference on Artificial Intelligence, AAAI, CCF-A 类会议): 2020, 2021, 2022
4. 国际人工智能联合大会 (International Joint Conference on Artificial Intelligence, IJCAI, CCF-A 类会议): 2020, 2021, 2022
5. 国际表示学习大会 (International Conference on Learning Representations, ICLR, CCF-A 类会议): 2021, 2022
6. ACM SIGKDD 知识发现与数据挖掘大会 (ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD, CCF-A 类会议): 2021
7. 亚太区知识发现与数据挖掘大会 (Pacific-Asia Conference on Knowledge Discovery and Data Mining, PAKDD, 数据挖掘重要会议): 2021

顶级国际学术期刊审稿人

1. IEEE Transactions on Knowledge and Data Engineering (IEEE TKDE)
2. Machine Learning Journal (MLJ)
3. Artificial Intelligence Journal (AIJ)

顶级国内学术会议及期刊审稿人

1. 第 18 届中国机器学习会议 (CCML2021)
2. 软件学报
3. Frontiers of Computer Science