

一种不平衡节点分类的数据增强算法 GraphSR

摘要

图神经网络 (GNN) 在节点分类任务中取得了巨大的成功。然而, 现有的 GNN 自然倾向于具有更多标记数据的多数类别, 而忽略那些具有相对较少标记数据的少数类别。传统的方法采用过采样的方法, 但往往会导致过拟合问题。最近, 一些工作提出从标记的节点中合成额外的少数类节点, 然而, 不能保证这些生成的节点是否能代表相应的少数类。事实上, 节点合成不当可能会导致算法泛化不足。为了解决这个问题, 本文寻求从图的大量未标记节点中自动增加少数类。具体来说, 本文提出了 GraphSR, 这是一种新的自我训练策略, 基于相似性的选择模块和强化学习 (RL) 选择模块, 以增加具有显著多样性的未标记节点的少数类。第一个模块找到与标记的少数节点最相似的未标记节点子集, 第二个模块通过 RL 技术进一步从子集中确定具有代表性和可靠的节点。此外, 基于 RL 的模块可以根据当前训练数据自适应确定采样尺度。这种策略是通用的, 可以很容易地与不同的 GNN 模型相结合。实验结果表明, 所提出的方法在各种类别不平衡数据集上优于最先进的基线。

关键词: 图不平衡学习; 图表示学习; 强化学习

1 引言

不平衡学习是指在分类问题中, 不同类别的样本数量差异较大, 导致模型训练和预测过程中出现偏倚的现象。在不平衡学习中, 通常会出现少数类样本数量较少, 而多数类样本数量较多的情况。在实际应用中, 很多情况下会出现正负样本比例严重不平衡的情况, 这会导致模型在训练和预测过程中对少数类样本的识别效果不佳。因此, 不平衡学习的目标是通过采用各种技术和算法来解决这种样本不平衡问题, 以提高模型对少数类样本的识别准确性。

2 相关工作

类不平衡表示学习是机器学习领域的一个经典话题, 已经得到了很好的研究 [1]。目标是在具有类不平衡分布的标记数据集上训练无偏分类器, 其中大多数类具有更多的样本, 而少数类具有更少的样本。相关工作包括重新加权和重新抽样方法。

2.1 基于重加权的方法

重新加权方法尝试通过提高少数类别的权重来修改损失函数, [2] [3] 或扩大少数群体的边缘 [4] [5] [6]。

2.2 基于重采样的方法

重抽样方法试图通过故意预处理训练样本来平衡数据分布，例如过度抽样少数类 [7]，抽样不足的多数类 [8]，以及两者的结合 [9]。随着神经网络的改进，重采样策略不仅通过采样技术 [10]，而且通过生成思想来增强少数类 [11] [12]。典型的 SMOTE 方法 [7] 通过对来自同一类的少数样本及其最近邻居使用插值技术来生成新样本。工作 [11] [12] 通过转移多数类的共同知识来合成少数类样本。然而，大多数现有的方法都是专门用于标识数据的，不能直接用于基于图的数据，在基于图的数据中应该考虑对象之间的关系。

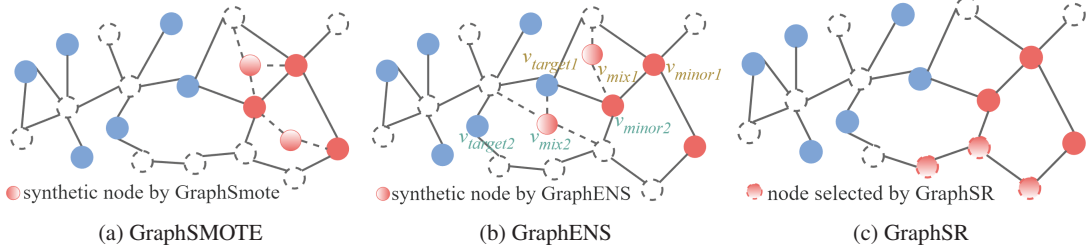


图 1. 这是一个带有有限标记节点和大量未标记节点的图，蓝色节点表示多数类，红色节点表示少数类，空白节点为未标记节点。(a) 合成节点由两个少数节点生成，边缘预测器生成边缘。(b) 混合节点及其单跳邻居由次要节点和目标节点生成。(c) 显示 GraphSR 对少数类的补充未标记节点。

3 本文方法

3.1 本文方法概述

在本节中，我们将介绍基于自我训练技术的 GraphSR 的细节。事实上，自我训练是一种经典的方法，在半监督学习中被广泛应用。原则上，该算法在可用的标记集上迭代训练模型，并使用训练好的模型对未标记的数据生成伪标记；然后，从未标记集中选择可信样本与训练集结合，进一步重新训练模型，直到收敛。

为了适应图中类不平衡的问题，基于自训练的思想，本文提出了两种组件来自适应地从未标记的数据中选择信息丰富且可靠的节点来补充少数类，如图 2 所示。首先，GraphSR 训练一个 GNN 模型基于标记集 V_L ，生成 U 中未标记节点的伪标签，然后设计基于相似性的选择模块，识别与少数类节点最相似的未标记节点，过滤出少数类的候选节点集 V_C 。其次，GraphSR 利用强化学习模块自适应选择信息可靠的节点，得到合适的补充集 V_L ，从而最有效地丰富少数类的多样性，最终增强训练集。利用增强的训练数据 $\{V_L, V_L\}$ ，我们可以训练一个类平衡节点分类器。在下文中，我们将展示每个组件的详细信息。

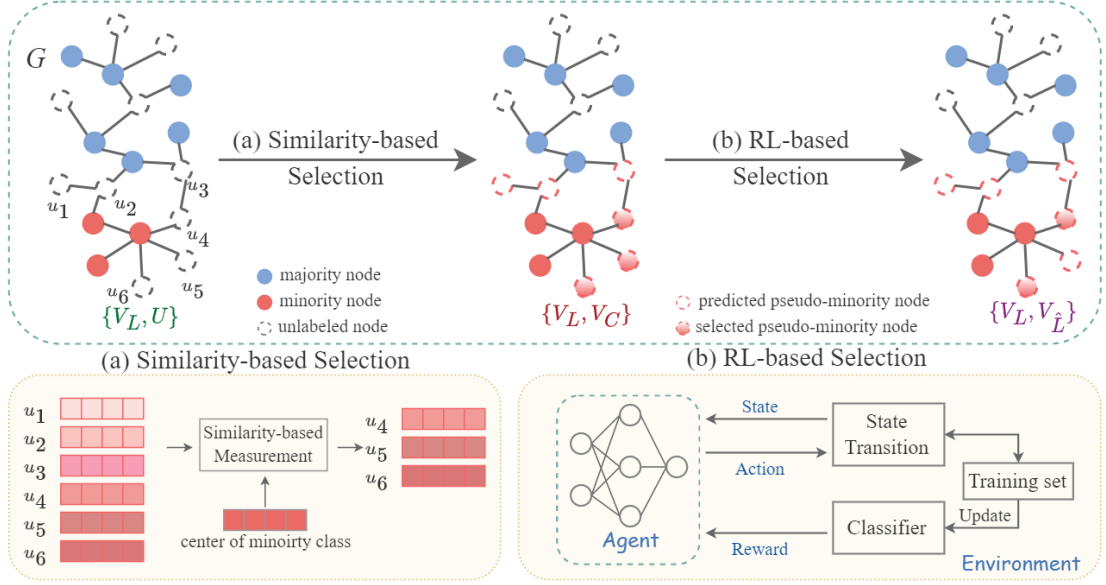


图 2. 方法示意图

3.2 相似性选择模块

在半监督设置中增强少数类的一种简单方法是从原始图中找到类似的未标记节点。一般来说，GNN 导出的节点表示可以反映节点的类型和类内关系，即同一类的节点在嵌入空间中会更接近，而来自不同类的节点在潜在空间中应该更远。因此，我们不是直接使用节点的原始属性来比较节点，而是在标记集上训练 GNN 模型来学习节点表示，它可以同时捕获节点的特征属性和拓扑信息。具体来说，我们在不平衡训练集 $\{V_L, Y_L\}$ 上训练 GNN 分类器 g ，其消息传递和融合过程表述为：

$$h_v^k = \sigma(W^k \cdot \text{CAT}(h_v^{k-1}, \text{AGG}(\{h_{v'}^{k-1}, \forall v' \in \mathcal{N}(v)\}))) \quad (1)$$

其中 $\text{AGG}(\cdot)$ 表示聚合邻域 $\mathcal{N}(v)$ 信息的聚合函数， $\text{CAT}(\cdot)$ 连接节点表示和邻域信息， W 是可学习的权重参数， σ 是非线性激活函数， h_v^k 表示具有 k 跳邻居的节点 v 的学习表示，并且 $h_v^0 = X[v, :]$ 。我们利用 z_v 来表示分类器 g 获得的节点 v 的嵌入。

此外，训练有素的 g 能够为未标记节点 u 生成伪标签 \hat{y}_u ，然后我们可以过滤掉一些可能不属于少数类的未标记节点以提高效率。为此，设 $M_i = \{u \in U | \hat{y}_u = i\}$ 表示被分类器 g 预测为少数类 i 的所有未标记节点的集合。从 M_i 中，我们仅通过基于相似性的模块选择那些距离嵌入空间中少数类中心足够近的节点。潜在空间中每个少数类的中心是根据标记节点计算的：

$$\text{cen}(i) = \frac{1}{|C_i|} \sum_{v \in C_i} z_v \quad (2)$$

其中 C_i 是类 i 的标记节点集。然后，该模块仅从 M_i 中选择距离中心 $\text{cen}(i)$ 最近的前 K 个节点作为候选集 V_C ，使得

$$V_C = \{u \in M_i | D(z_u, \text{cen}(i)) < \varphi \text{ and } |V_C| \leq K\} \quad (3)$$

其中 $D(\cdot, \cdot)$ 衡量嵌入空间中的相似度，我们采用欧氏距离作为衡量标准 $D(u, v) = \|z_u - z_v\|$ ， φ 表示 V_C 中的节点与 $\text{cen}(i)$ ， $\varphi = \max(z_u, \text{cen}(i) | u \in V_C)$ 。

通过这种选择，我们可以找到最有可能被预测为少数类的节点，但是， g 并不可靠，因为它是用不平衡数据进行训练的。为了解决这个问题，GraphSR 利用强化学习的另一个选择模块来提取能够准确提高分类器性能的可靠节点，并自适应地确定每个少数类的过采样规模。

3.3 基于强化学习的选择模块

该选择模块的关键任务是指定一个采样过程，该过程可以自适应地选择未标记的节点来补充少数类。由于缺乏未标记节点的监督信息，我们采用强化学习进行节点选择。我们设计了一个迭代采样程序，其表述为马尔可夫决策过程 (MDP)， $M = (S, A, R, T)$ 。生成平衡训练集的过程可以用轨迹 $(s_0, a_0, r_0, \dots, s_T, a_T, r_T)$ 来描述，其中初始状态 s_0 仅包含不平衡标记集，最后状态 s_T 包含最终补充平衡节点集。GraphSR 尝试使用强化学习算法来学习最佳策略，以允许代理决定在部分观察的环境中保留或丢弃未标记的节点。具体来说，代理（即选择器）顺序遍历 V_C 中的候选少数节点。对于每个节点 u_t ，代理根据当前状态通过由 π_θ 表示的策略网络采取行动，然后环境根据该行动分配奖励。代理根据奖励更新策略网络。在智能体和环境之间进行足够的交互之后，智能体可以学习最优策略来最优地选择未标记的节点来补充少数类。通过基于强化学习的选择，算法更容易推广到不同的数据集，而无需额外确定过采样规模。下面，我们详细讨论基于 RL 的选择模块的主要组成部分。

状态 我们定义环境的状态 s_t 来编码时间步 t 的中间训练集 V_t 和未标记节点 u_t 的信息。为了将 s_t 输入到策略网络中，我们需要固定 s_t 的维度，与 V_t 中的节点数量无关。受 [13] 的启发，我们使用 V_t 中节点嵌入的总和来表示 V_t 的信息，即 $z_{V_t} = \sum_{v \in V_t} z_v$ 。此外，利用 u_t 的嵌入来表示其信息。对于时间步 t ，状态 s_t 定义为 $s_t = (z_{V_t} z_{u_t}, \cdot)$ 。开始时， $V_0 = V_L$ 为不平衡标记集， u_0 为 V_C 中的第一个节点。

动作 a_t 的动作是在时间步 t 决定是否将 V_C 中当前未标记的节点 u_t 包含到当前训练集 V_t 中。特别地， $a_t \in \{0, 1\}$ ，其中 $a_t = 1$ 表示选择节点 u_t 来补充不平衡训练集，而 $a_t = 0$ 表示 u_t 不适用。此外，动作由策略函数 π_θ 生成，该函数将状态作为输入并由参数化。在这项工作中，策略网络表示动作的概率分布，并被指定为具有非线性激活函数的多层感知器 (MLP)，即：

$$a_t = P(a_t | s_t) = \pi_\theta(s_t) = MLP_\theta(s_t) \quad (4)$$

转换 在 a_t 采取行动后，环境状态应更改为 $s_t + 1$ 。在我们的工作中，状态由 V_t 和 u_t 组成，取 a_t 后，

$$V_{t+1} = \begin{cases} \{V_t \cup u_t\}, & a_t = 1 \\ V_t, & a_t = 0 \end{cases} \quad (5)$$

然后 $s_{t+1} = (z_{V_{t+1}} z_{u_{t+1}}, \cdot)$ 。当智能体完全遍历完候选集 V_C 一次后，转换就会终止。

奖励 环境给予的奖励 rt 是评估状态 st 时的动作。如果没有未标记节点的监督信息，很难准确找到少数节点并根据真实标签对其进行显式奖励。在这里，我们基于 $V_t \cup u_t$ 训练分类器，并在小型平衡验证集上评估其准确性。然而，准确性始终是非负的，直接将其用作奖励可能会阻碍智能体的收敛。奖励工程的思想是，如果添加 u_t 可以提高分类器的性能，则分配正奖

励，否则分配负奖励。因此，奖励函数设计为：

$$r_t = \begin{cases} +1, & acc_t \geq b_t \text{ and } a_t = 1 \\ -1, & acc_t < b_t \text{ and } a_t = 1 \\ +1, & acc_t < b_t \text{ and } a_t = 0 \\ -1, & acc_t \geq b_t \text{ and } a_t = 0 \end{cases} \quad (6)$$

其中, b_t 表示基线奖励, 是过去 10 个准确率的平均值, 即 $b_t = \text{mean}\{acc_{t-11}, \dots, acc_{t-1}\}$, acc_0 表示由标记节点集训练的初始分类器的准确率 V_L 。策略梯度训练代理的目标是训练一个可以最大化预期奖励的最优策略网络, 并且基于策略梯度的方法被广泛用于优化策略网络。在这项工作中, 我们使用近端策略优化 (PPO) [14] 来更新策略网络的参数 θ 。PPO 的目标函数定义为：

$$L^{CLIP}(\theta) = \mathbb{E}_t[\min(p_t(\theta)\hat{A}_t, \text{clip}(p_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t)] \quad (7)$$

其中 $p_t(\theta)$ 是概率比, $p_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$, 它被剪裁到范围 $[1 - \epsilon, 1 + \epsilon]$ 中, 从而形成了保守的策略迭代目标 [15] 和代理的探索更加稳定。 \hat{A}_t 是估计优势函数, 涉及折扣累积奖励和价值函数 V_{π} , 广泛应用于策略梯度算法中。

3.4 基于 GNN 的分类器

如前所述, 通过基于相似性的选择和基于强化学习的选择, GraphSR 从未标记的数据中抽取最具信息量和最可靠的节点, 以补充少数类进行训练。我们用标记集 $\{V_L, V_{\bar{L}}\}$ 得到最终的训练集 $\{(v_i, y_i)\}$, 并用伪标签 $\{(u_i, y_i)\}$ 补充训练集。基于新的训练数据, 我们可以根据消息传递过程为等式 1 和交叉熵损失函数来训练无偏 GNN 分类器 f 。GraphSR 的端到端训练过程概述在附录中。

4 复现细节

4.1 与已有开源代码对比

论文代码情况：只给出了一个测试用例, 包含一个数据集和测试所需要的代码。没有提供强化学习部分代码, 并且所提供的代码参数但是写的固定数字的参数, 需要对代码进行很大的改动才能够在别的数据集上进行训练。

论文中将候选集大小设置为 20, 并通过分类器进行初步筛选少数类节点候选集。但是由于训练集的不平衡, 分类模型可能会难以区分出不同类别的少数类节点, 那么该类别的候选集大小为 0。候选集为 0 就不能为后续的强化学习模块挑选准确的少数类节点加入训练集, 这导致模型难以运行成功。基于这种事实, 本文考虑采用带权重的交叉熵损失函数来解决这个问题。交叉熵损失函数如(8)所示, 为少数类增加权重, 使得分类模型能更好地区分出不同类别的少数类节点。

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C w_j y_{i,j} \log p_{i,j} \quad (8)$$

4.2 实验设置

4.2.1 数据集

实验采用 Cora、CiteSeer、PubMed [16]、Amazon-Photo (Photo)、Amazon-Computers (computer) 和 Coauthor-CS (CS) [17] 6 个基准数据集进行所有实验。这些数据集的统计数据如表1所示。

表 1. 数据集统计信息

数据集	节点数	边数	特征数	类别数
Cora	2708	10556	1433	7
CiteSeer	33271	9104	3703	6
PubMed	19717	88648	500	3
Photo	7650	119081	745	10
Computer	13752	245861	767	10
CS	183333	81894	6805	15

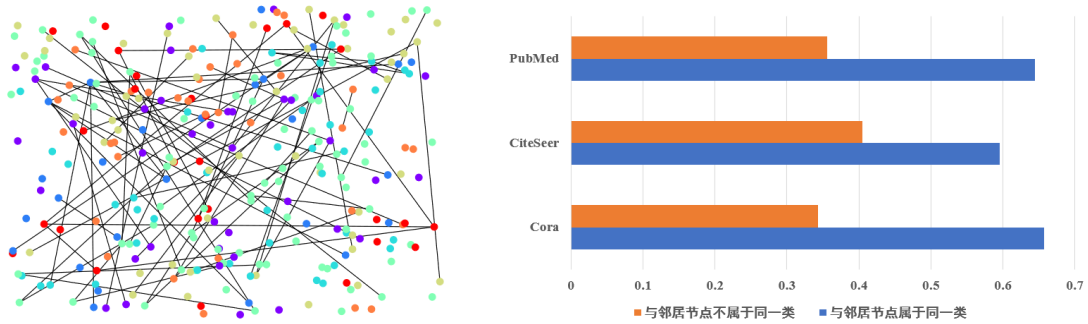


图 3. 各数据集节点是否与邻居节点属于同一类统计

4.2.2 实验参数设置

实验采用一个 3 层的 MLP 作为策略网络，隐藏层有 128 个维度，使用 ReLu 函数激活。使用 PPO 算法对具有默认超参数的智能体进行训练，学习率设为 0.005。使用 Adam 优化器对所有参数进行优化。对于 GNN 分类器，采用 GCN 和 GraphSAGE 两种基本架构。两者都由 2 个 GCN 或 SAGE 层和 128 个隐藏单元和 ReLU 激活组成。学习率设为 0.01，Dropout 率设为 0.5。GNN 分类器使用 SGD 优化器进行训练，该优化器具有早期停止机制，训练次数为 2000。具体来说，基于相似性的选择模块中候选集 K 的大小设为 20。对于所有比较的方法，报告 5 次运行的平均值和标准差。

多数类少数类的划分，根据数据集类别数多数的一半类别作为多数类，剩下的类别作为少数类。多数类的训练节点数量设置为 20，少数类的训练节点数量设置为 $20 \times \rho$ 其中 ρ 为不平衡率，设置为 0.3。每个类别的用于验证和测试的节点数量分别设置为 30，100。

4.3 创新点

针对训练智能体的过程时间复杂度较的局限，本文做出了一定的改进。如图 3 所示，通过对数据集的部分节点进行可视化，发现相同类的节点之间存在连接。进一步分析后发现，数据集与其邻居节点属于同一类的概率超过了 60%，基于此：

本文在进行相似化选择后，采用基于邻居的方法进行采样，有效降低了训练时间复杂度。

5 实验结果分析

表 2. 复现实验结果与论文实验结果对比

Method		Cora			CiteSeer			PubMed		
		ACC	F1	AUC-ROC	ACC	F1	AUC-ROC	ACC	F1	AUC-ROC
GCN	GraphSR(论文)	73.90	73.59	90.21	57.28	55.20	83.67	71.79	71.70	85.14
	GraphSR(复现)	72.57	73.78	93.71	59.46	56.12	82.73	70.84	70.71	86.38
GraphSAGE	GraphSR(论文)	78.78	78.36	94.92	54.30	51.15	84.21	74.13	74.36	89.33
	GraphSR(复现)	78.57	78.96	95.56	55.86	53.23	87.17	73.59	73.12	89.14

表 3. 改进后实验结果与论文实验结果对比

Method		Cora			CiteSeer			PubMed		
		ACC	F1	AUC-ROC	ACC	F1	AUC-ROC	ACC	F1	AUC-ROC
GCN	GraphSR(论文)	73.90	73.59	90.21	57.28	55.20	83.67	71.79	71.70	85.14
	GraphSR(改进)	72.17	72.27	93.82	57.33	57.26	82.00	69.33	70.01	84.62
GraphSAGE	GraphSR(论文)	78.78	78.36	94.92	54.30	51.15	84.21	74.13	74.36	89.33
	GraphSR(改进)	77.43	77.37	94.52	54.71	54.48	86.65	72.67	72.85	88.69

表 3 为实验结果，复现的实验结果与论文中性能基本一致，在 Citeseer 数据集上甚至超过了论文，但在 Pubmed 数据集上存在微弱的差距。

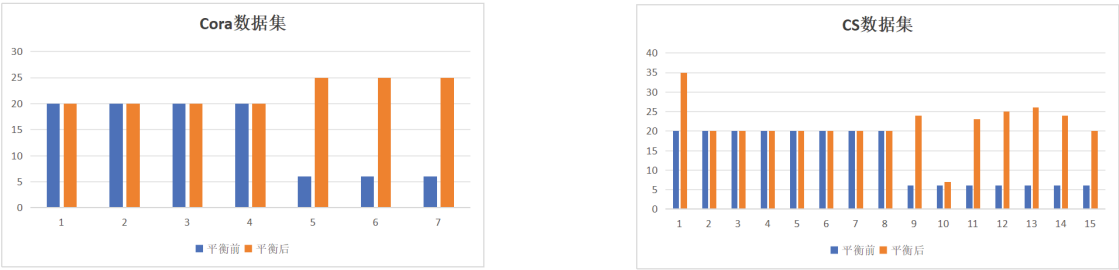


图 4. 两个数据集上选出的少数类节点统计

类别不多的数据集能够取得较好的实验效果，类别数多的数据集难以选取到合适少数类节点。这是因为随着类别增多，强化学习的动作空间变大了，难以去选取合适的节点。同时，节点嵌入的低维空间的特征，也难以保证是好的。

6 总结与展望

研究了半监督环境下具有类不平衡问题的节点分类问题。为了充分利用大量未标记节点的丰富信息，我们提出了一种新的数据增强策略 GraphSR，该策略可以通过基于相似性的选择模块和基于 RL 的选择模块自动补充大量未标记节点的少数类。此外，基于 rl 的模块可以自适应地确定不同少数类的过采样尺度。我们验证了所提出的模型可以有效地丰富少数类的多样性，并在一定程度上避免过拟合。实验结果证明了 GraphSR 在不同 GNN 架构的数据集上的有效性和鲁棒性。假设数据集中有足够数量的少数类节点可供补充少数类。训练智能体的过程异常耗时。考虑将与少数类节点相连接的无标签节点加入是否会更合适？

参考文献

- [1] H He and EA Garcia. Learning from imbalanced data *ieee transactions on knowledge and data engineering*, 2009.
- [2] 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.
- [3] 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/CVF conference on computer vision and pattern recognition*, pages 9268–9277, 2019.
- [4] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. *Advances in neural information processing systems*, 32, 2019.
- [5] Bingyu Liu, Weihong Deng, Yaoyao Zhong, Mei Wang, Jiani Hu, Xunqiang Tao, and Yaohai Huang. Fair loss: Margin-aware reinforcement learning for deep face recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10052–10061, 2019.
- [6] Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit, and Sanjiv Kumar. Long-tail learning via logit adjustment. *arXiv preprint arXiv:2007.07314*, 2020.
- [7] 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.
- [8] Miroslav Kubat, Stan Matwin, et al. Addressing the curse of imbalanced training sets: one-sided selection. In *Icml*, volume 97, page 179. Citeseer, 1997.
- [9] Gustavo EAPA Batista, Ronaldo C Prati, and Maria Carolina Monard. A study of the behavior of several methods for balancing machine learning training data. *ACM SIGKDD explorations newsletter*, 6(1):20–29, 2004.

- [10] Zhining Liu, Pengfei Wei, Jing Jiang, Wei Cao, Jiang Bian, and Yi Chang. Mesa: boost ensemble imbalanced learning with meta-sampler. *Advances in neural information processing systems*, 33:14463–14474, 2020.
- [11] Jaehyung Kim, Jongheon Jeong, and Jinwoo Shin. M2m: Imbalanced classification via major-to-minor translation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13896–13905, 2020.
- [12] Jianfeng Wang, Thomas Lukasiewicz, Xiaolin Hu, Jianfei Cai, and Zhenghua Xu. Rsg: A simple but effective module for learning imbalanced datasets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3784–3793, 2021.
- [13] Chuan Luo, Pu Zhao, Chen Chen, Bo Qiao, Chao Du, Hongyu Zhang, Wei Wu, Shaowei Cai, Bing He, Saravanakumar Rajmohan, et al. Pulns: Positive-unlabeled learning with effective negative sample selector. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 8784–8792, 2021.
- [14] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [15] Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the Nineteenth International Conference on Machine Learning*, pages 267–274, 2002.
- [16] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- [17] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. *arXiv preprint arXiv:1811.05868*, 2018.