

Class and Attribute-Aware Logit Adjustment for Generalized Long-Tail Learning

Xiaoling Zhou^{1,2}, Ou Wu^{1,3*}, and Nan Yang¹

¹Center for Applied Mathematics, Tianjin University, China

²National Engineering Research Center for Software Engineering, Peking University, China

³HIAS, University of Chinese Academy of Sciences, Hangzhou, China

{xiaolingzhou, yny, wuou}@tju.edu.cn

Abstract

Compared to conventional long-tail learning, which focuses on addressing class-wise imbalances, generalized long-tail (GLT) learning considers that samples within each class still conform to long-tailed distributions due to varying attributes, known as attribute imbalance. In the presence of such imbalance, the assumption of equivalence between the class-conditional probability densities of the training and testing sets is no longer tenable. Existing GLT approaches typically employ regularization techniques to avoid directly modeling the class-conditional probability density (CCPD) ratio between training and test data, leading to suboptimal performance. This study aims to directly estimate this ratio, for which a novel class-attribute aware logit-adjusted (CALA) loss incorporating both the CCPD ratio and the class priors is presented. Two new GLT learning methods, named Heuristic-CALA and Meta-CALA, are then proposed, which estimate the CCPD ratio in the CALA loss by leveraging the neighborhood information of samples. Extensive experiments across diverse scenarios susceptible to class and attribute imbalances showcase the state-of-the-art performance of Meta-CALA. Furthermore, while Heuristic-CALA exhibits inferior performance compared to Meta-CALA, it incurs only negligible additional training time compared to the Cross-Entropy loss, yet surpasses existing methods by a significant margin.

Introduction

Long-tail (LT) learning is a common challenge in many real applications, where only a few categories are represented by a large number of instances while many others are represented by only a few (Cui et al. 2019; Zhou, Yang, and Wu 2023). A popular technique to address this challenge is logit adjustment (Menon et al. 2021; Wang et al. 2024; Zhao et al. 2022). However, existing methods (Tao et al. 2023; Menon et al. 2021; Li et al. 2021) typically assume that the primary difference between training and test data lies in the prior probabilities over categories, where $p_{\text{tr}}(y) \neq p_{\text{te}}(y)$, while the class-conditional probabilities for the training and test data remain the same, i.e., $p_{\text{tr}}(\mathbf{x}|y) = p_{\text{te}}(\mathbf{x}|y)$. Additionally, these methods presume uniform prior probabilities over classes when evaluating model performance. Therefore, the adjustment terms in existing methods (Menon et al. 2021;

Cao et al. 2019) are primarily based on $p_{\text{tr}}(y)$. Additionally, several imbalanced benchmarks, such as CIFAR-LT (Cui et al. 2019), are manually constructed based on these assumptions, making existing algorithms well-suited for these datasets but limiting their generalizability to others.

Recently, Tang et al. (2022) emphasized that the assumption of identical class-conditional probability densities between training and test data cannot be guaranteed for real-world datasets. They have thus identified a new type of imbalance known as attribute imbalance. For instance, concerning the color attribute, the training set may consist mostly of white doves, whereas the number of white and dark doves may be equal in the testing set. Attribute imbalance can lead to poor performance of samples with rare attributes and compromise the generalization ability of deep learning models. Consequently, they formulated a new research problem, named generalized long-tail (GLT) learning, which encompasses both class and attribute imbalances. An example of GLT learning is shown in Fig. 1. Given the poor performance of existing LT baselines for GLT learning, they proposed a regularization technique to learn invariant features. Despite demonstrating improved performance, the challenge of attribute imbalance remains unsolved because the fundamental issue, $p_{\text{tr}}(\mathbf{x}|y) \neq p_{\text{te}}(\mathbf{x}|y)$, has not been adequately addressed.

This study pioneers the estimation of the class-conditional probability density (CCPD) ratio between training and test data. We first introduce a modified Cross-Entropy (CE) loss, termed class-attribute aware logit-adjusted (CALA) loss, which incorporates both class priors and the ratio of class-conditional probability densities as adjustment terms to address class and attribute imbalances. Next, we develop a novel GLT method called Heuristic-CALA, which utilizes neighborhood information of samples to estimate the CCPD ratio within the CALA loss. Notably, Heuristic-CALA serves as a generalization of several conventional LT approaches. Finally, leveraging the strong performance of meta-learning, we propose another GLT method named Meta-CALA. This method employs an adjustment network optimized through meta-learning to estimate the CCPD ratio based on the neighborhood-related training characteristics of samples. We conduct extensive experiments across various learning scenarios prone to class and attribute imbalances: LT learning, GLT learning, and subpopulation shift learn-

*Corresponding author.

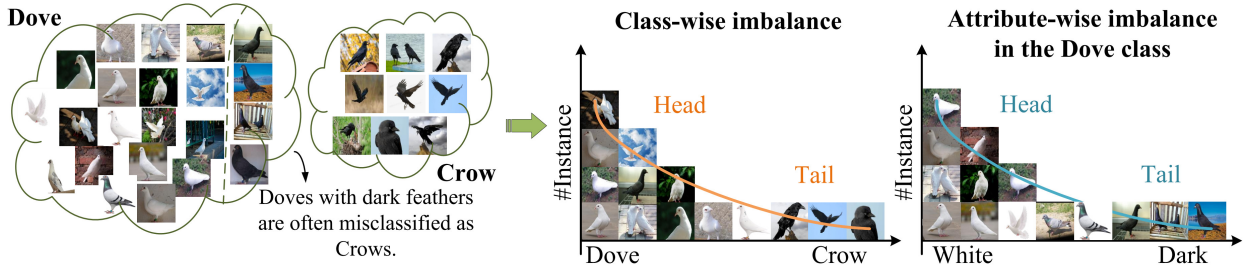


Figure 1: Illustration of imbalances at both the class and attribute levels. Dove and Crow represent the head and tail classes, respectively. Due to the prevalence of white feathers among doves, there exists an attribute imbalance within the Dove class.

ing. The results demonstrate that our methods consistently achieve state-of-the-art (SOTA) performance by effectively addressing both class and attribute imbalances.

Our main contributions can be summarized as follows:

- We conduct a pioneering exploration by directly utilizing the CCPD ratio for logit calibration. A novel logit adjustment loss (termed CALA) that accounts for both class priors and the ratio of class-conditional probability densities between training and test data is then presented.
- We propose two new GLT learning methods, Heuristic-CALA and Meta-CALA, which employ K -neighborhood-based and meta-learning-based estimation approaches, respectively, to estimate the CCPD ratio in the CALA loss.
- We conduct extensive experiments across three learning scenarios susceptible to class and attribute imbalances. The results conclusively demonstrate the effectiveness of our approaches in enhancing the generalization and robustness of deep learning models.

Related Work

Long-Tail Classification Despite success in various applications, deep neural networks still struggle with long-tailed datasets (De Alvis and Seneviratne 2024; Mao, Fan, and Li 2023). Different approaches have been proposed to address this issue, including algorithms based on resampling (Lin, Tsai, and Lin 2023; Tripathi, Chakraborty, and Kopparapu 2021; Yan et al. 2019), reweighting (Wan et al. 2023; Cui et al. 2019), knowledge distillation (Zhang et al. 2023), data augmentation (Li et al. 2021; Zhou et al. 2024; Zhou and Wu 2023), multiple experts (Wang et al. 2020; Xiang, Ding, and Han 2020), and contrastive learning (Cui et al. 2021). Among these methods, logit adjustment-based approaches (Wang et al. 2024; Li, Cheung, and Lu 2022; Menon et al. 2021) have gained popularity and demonstrated their effectiveness. For instance, LA (Menon et al. 2021) perturbs the logits of samples to encourage a large relative margin between logits of rare versus dominant labels. More recently, ALA (Zhao et al. 2022) introduces an adaptive adjustment term that consists of two complementary factors: a quantity factor and a difficulty factor. However, existing methods (Wang et al. 2024; Tao et al. 2023) primarily focus on addressing class-wise imbalances, while overlooking the imbalanced attribute distribution within each class.

Generalized Long-Tail Classification Tang et al. (2022) argued that imbalanced classifications suffer from both class- and attribute-wise imbalances and, therefore, proposed the GLT learning task. They subsequently presented an invariant feature learning (IFL) method to tackle GLT learning by maintaining the feature center of each class across different environments. Apart from this approach, there are currently few dedicated solutions available to address the emerging GLT problem. Nevertheless, some methods tailored for addressing subpopulation shift (Deng et al. 2024; Liang and Zou 2022; Koh et al. 2021) and spurious correlation (Chen et al. 2023; Srivastava, Hashimoto, and Liang 2020; Agarwal, Shetty, and Fritz 2020) are deemed effective for tackling GLT learning by implicitly mitigating the issue of attribute imbalance. However, nearly all existing studies overlook the direct modeling of attribute distributions within each class, leading to subpar performance when dealing with attribute imbalance.

Class-Attribute Aware Logit-Adjusted Loss

Following prior studies, the classification model is formulated as $p(y|\mathbf{x})$, which predicts the label y from the input \mathbf{x} . The training and test data are drawn from different joint distributions, namely $p_{tr}(\mathbf{x}, y)$ and $p_{te}(\mathbf{x}, y)$, respectively. Utilizing Bayes' Rule, we have $p_{tr}(y|\mathbf{x}) \propto p_{tr}(\mathbf{x}|y)p_{tr}(y)$ and $p_{te}(y|\mathbf{x}) \propto p_{te}(\mathbf{x}|y)p_{te}(y)$. Hence, we arrive at

$$p_{tr}(y|\mathbf{x}) \propto p_{te}(y|\mathbf{x}) \cdot \frac{p_{tr}(\mathbf{x}|y)}{p_{te}(\mathbf{x}|y)} \cdot \frac{p_{tr}(y)}{p_{te}(y)}. \quad (1)$$

To simplify Eq. (1), existing methods commonly rely on the following two assumptions:

Assumption 1 *The class-conditional probability densities of the training and testing sets are equal: $\forall \mathbf{x}, y, p_{tr}(\mathbf{x}|y)/p_{te}(\mathbf{x}|y) \equiv 1$.*

Assumption 2 *The class priors $p_{te}(y)$ are assumed to be identical when evaluating the model's performance.*

Given the two assumptions mentioned above, the objective of Eq. (1) can be expressed as

$$\arg \max p_{te}(y|\mathbf{x}) = \arg \max p_{tr}(y|\mathbf{x})/p_{tr}(y). \quad (2)$$

From Eq. (2), the adjustment terms should be determined by the class prior $p_{tr}(y)$, which has been adopted by existing LT learning proposals, such as LDAM (Cao et al. 2019) and LA (Menon et al. 2021).

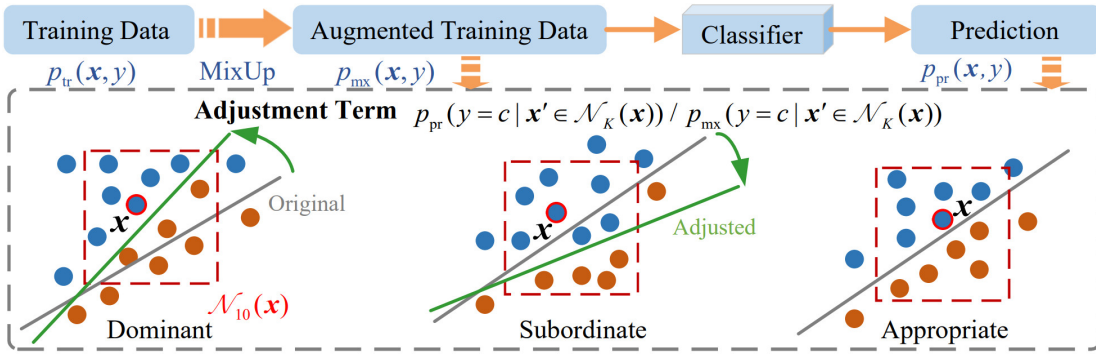


Figure 2: Illustration of Heuristic-CALA. Blue and orange dots represent samples from two classes. The gray and green lines denote the original and adjusted classifiers. The red boxes indicate the ten-nearest neighborhood $\mathcal{N}_{10}(\mathbf{x})$. If $p_{pr}(y = y_{\mathbf{x}} | \mathbf{x}' \in \mathcal{N}_K(\mathbf{x})) > (<, =) p_{mx}(y = y_{\mathbf{x}} | \mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$, then \mathbf{x} is in a dominant (subordinate, appropriate) position in the training set. Our adjustment term will make \mathbf{x} easier (or harder, or leave it unchanged) than before, facilitating better adjusting the classifier.

Assumption 2 evidently promotes fairness across different categories. However, we challenge the validity of Assumption 1 in practical learning scenarios and argue that the objective function in Eq. (2) is overly simplified, rendering it incapable of resolving a number of issues, such as the misclassification of samples with rare attributes. As illustrated in Fig. 1, although the number of samples in the Dove class exceeds that of the Crow class, doves with dark feathers are often misclassified as crows due to the fact that $p_{tr}(\mathbf{x}_{\text{feather}} = \text{white} | y = \text{Dove}) \gg p_{tr}(\mathbf{x}_{\text{feather}} = \text{dark} | y = \text{Dove})$ and $p_{tr}(\mathbf{x}_{\text{feather}} = \text{white} | y = \text{Crow}) \ll p_{tr}(\mathbf{x}_{\text{feather}} = \text{dark} | y = \text{Crow})$ ¹. Even if the issue of class imbalance is addressed, the presence of attribute imbalance remains and substantially impairs the generalization performance of models.

Actually, the training objective without guarantying Assumption 1 should be

$$\arg \max p_{te}(y|\mathbf{x}) = \arg \max p_{tr}(y|\mathbf{x}) \cdot \frac{p_{te}(\mathbf{x}|y)}{p_{tr}(\mathbf{x}|y)} \cdot \frac{1}{p_{tr}(y)}. \quad (3)$$

Eq. (3) suggests that the adjustment terms should be determined by both the CCPD ratio and the class priors. Accordingly, building on the inference method used in LA (Menon et al. 2021), we derive a novel logit-adjusted loss that incorporates these two terms to mitigate attribute imbalance and class bias. With $\tau_1, \tau_2 > 0$, the resulting loss, termed CALA, is as follows:

$$\begin{aligned} \ell_{CALA}(\mathbf{x}) &= -\log \frac{\exp[f_y(\mathbf{x}) + \tau_1 \log p_{tr}(y) + \tau_2 \log \frac{p_{tr}(\mathbf{x}|y)}{p_{te}(\mathbf{x}|y)}]}{\sum_{y' \in [C]} \exp[f_{y'}(\mathbf{x}) + \tau_1 \log p_{tr}(y') + \tau_2 \log \frac{p_{tr}(\mathbf{x}|y')}{p_{te}(\mathbf{x}|y')}]}, \\ &= \log \left[1 + \sum_{y' \neq y} \left(\frac{p_{tr}(y')}{p_{tr}(y)} \right)^{\tau_1} \cdot \left(\frac{p_{tr}(\mathbf{x}|y') p_{te}(\mathbf{x}|y)}{p_{te}(\mathbf{x}|y') p_{tr}(\mathbf{x}|y)} \right)^{\tau_2} \cdot \frac{e^{f_{y'}(\mathbf{x})}}{e^{f_y(\mathbf{x})}} \right], \end{aligned} \quad (4)$$

where $f(\cdot)$ and C represent the classifier and the number of classes. We then explain the CALA loss from a regularization perspective using Taylor expansion. Our analysis indicates that the CALA loss imposes significant penalties on

¹It is worth noting that $p_{te}(\mathbf{x}_{\text{feather}} = \text{white} | y = \text{Dove}) = p_{te}(\mathbf{x}_{\text{feather}} = \text{dark} | y = \text{Dove})$ is assumed to be established for fairness. Consequently, $p_{tr}(\mathbf{x} | y = \text{Dove}) \neq p_{te}(\mathbf{x} | y = \text{Dove})$.

samples from tail classes (e.g., Crow) and those with rare attributes (e.g., dark doves). This increased penalization amplifies the impact of these samples during model training, thereby enhancing their prediction performance. However, directly obtaining the CCPD ratio $p_{tr}(\mathbf{x}|y)/p_{te}(\mathbf{x}|y)$ is impossible due to the unknown $p_{te}(\mathbf{x}|y)$. To this end, we propose two estimation approaches, as stated in the subsequent sections.

Learning with CALA Loss

To estimate the CCPD ratio in the CALA loss, we propose two methods: one based on K -neighborhood and the other based on meta-learning. Consequently, two logit adjustment approaches, named Heuristic-CALA and Meta-CALA, are devised for GLT learning.

Heuristic-CALA Framework

To simplify the notation, the CALA loss is expressed as

$$\ell_{CALA}(\mathbf{x}) = -\log \frac{\exp[f_y(\mathbf{x}) + \tau_1 u(y) + \tau_2 v(\mathbf{x}, y)]}{\sum_{y' \in [C]} \exp[f_{y'}(\mathbf{x}) + \tau_1 u(y') + \tau_2 v(\mathbf{x}, y')]}, \quad (5)$$

where $u(y) = \log p_{tr}(y)$ and $v(\mathbf{x}, y) = \log[p_{tr}(\mathbf{x}|y)/p_{te}(\mathbf{x}|y)]$. The neighborhood information of each sample reflects its local distribution in the feature space, thus providing an estimate of the CCPD values. Accordingly, we establish the following relationship, with detailed inference provided in Section B.I of the Appendix:

$$p_{te}(y|\mathbf{x}) \approx p_{tr}(y|\mathbf{x}) \cdot \frac{p_{te}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))}{p_{tr}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))}, \quad (6)$$

where $\mathcal{N}_K(\mathbf{x})$ represents the K -nearest neighbors of \mathbf{x} . Eq. (6) manifests that the adjustment terms should be determined by $p_{te}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))/p_{tr}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$. However, the distribution of the test data is not readily available. Zhang et al. (2018) proposed a generic vicinal distribution called MixUp, which effectively estimates the unknown data distribution and significantly improves models' generalization performance. Consequently, we employ the training data augmented using MixUp, denoted as $p_{mx}(\mathbf{x}, y)$, to

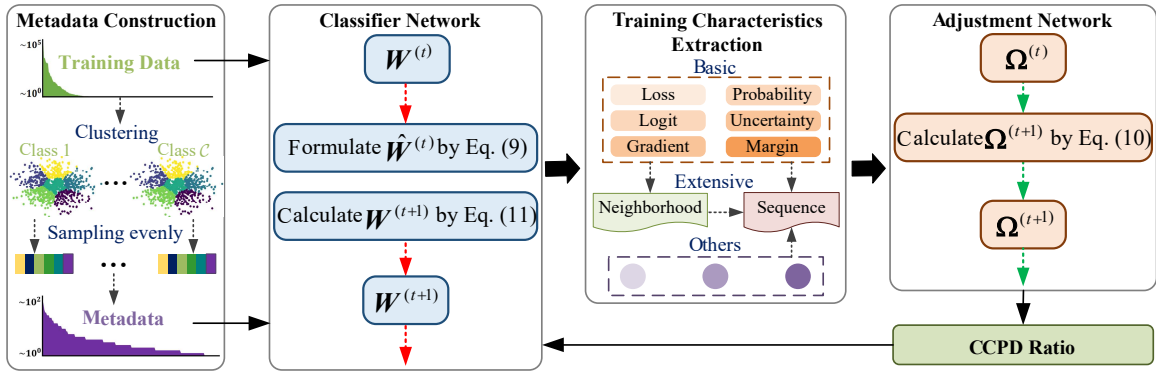


Figure 3: Illustration of Meta-CALA, which contains four main components, including the metadata construction module, the classifier network, the training characteristics extraction module, and the adjustment network.

approximate the test data. However, MixUp applied to imbalanced datasets fails to balance the label distribution. To address this, we utilize the ratio $p_{\max}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))/p_{\text{tr}}(y)$ to substitute $p_{\text{te}}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$. Furthermore, $p_{\text{tr}}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$ is replaced by the predicted probability $p_{\text{pr}}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$ to more precisely and dynamically adjust the classifier during training. Thus, Eq. (6) transforms into

$$p_{\text{te}}(y|\mathbf{x}) := p_{\text{tr}}(y|\mathbf{x}) \cdot \frac{p_{\max}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))}{p_{\text{pr}}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))} \cdot \frac{1}{p_{\text{tr}}(y)}. \quad (7)$$

We observe that Eq. (7) employs the ratio $p_{\text{pr}}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))/p_{\max}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$ to replace the CCPD ratio $p_{\text{tr}}(\mathbf{x}|y)/p_{\text{te}}(\mathbf{x}|y)$ in Eq. (3). Consequently, the adjustment term $v(\mathbf{x}, y)$ can be calculated as

$$v(\mathbf{x}, y) = \log \left[\frac{p_{\text{pr}}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))}{p_{\max}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))} \right], \quad (8)$$

where $p_{\max}(y = c|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$ represents the label distribution of samples within the neighborhood. Additionally, $p_{\text{pr}}(y = c|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$ denotes the prediction distribution within the neighborhood, which can be calculated as $\sum_{\mathbf{x}' \in \mathcal{N}_K(\mathbf{x})} q_{\mathbf{x}', c}/K$, where $\mathbf{q}_{\mathbf{x}'} = \text{Softmax}(f(\mathbf{x}'))$ represents the probability vector. Furthermore, since neighborhood information can be sensitive to border effects and outliers, we employ class averages of the corresponding values within the neighborhood to smooth the values of $p_{\text{pr}}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$ and $p_{\max}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$.

We then validate the rationality of our approach by elaborating on its meaning. The term $p_{\text{pr}}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))/p_{\max}(y|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$ reflects the dominance of \mathbf{x} in the training data. As shown in Fig. 2, if $p_{\text{pr}}(y = y_{\mathbf{x}}|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x})) > (<, =) p_{\max}(y = y_{\mathbf{x}}|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))$, then \mathbf{x} is in a dominant (subordinate, appropriate) position in the training set. Our adjustment term will make \mathbf{x} easier (or harder, or leave it unchanged) than before, resulting in a smaller (or larger, or unchanged) impact on the model training. Generally, samples from tail classes and those with rare attributes occupy a subordinate position, and their influence will be amplified by our approach.

Furthermore, we demonstrate that several typical LT base-lines can be viewed as special cases of Heuristic-CALA. Indeed, the adjustment term for class c in Heuristic-CALA is

determined by $p_{\text{tr}}(y = c) \cdot \frac{p_{\text{pr}}(y = c|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))}{p_{\max}(y = c|\mathbf{x}' \in \mathcal{N}_K(\mathbf{x}))}$. As when $K = 0$, $\mathcal{N}_0(\mathbf{x}) = \mathbf{x}$, we have the following findings:

- If $K = 0$ and $p_{\text{pr}}(y = y_{\mathbf{x}}|\mathbf{x}) \equiv 1$, as $p_{\max}(y = y_{\mathbf{x}}|\mathbf{x}) \equiv 1$, then only the adjustment term for class $y_{\mathbf{x}}$ is non-zero and equals to $p_{\text{tr}}(y = y_{\mathbf{x}})$, relying on the class proportion to adjust the logit of $y_{\mathbf{x}}$. Therefore, Heuristic-CALA is equivalent to LDAM (Cao et al. 2019) in this scenario.
- If $K = 0$ and $p_{\text{pr}}(y = y_{\mathbf{x}}|\mathbf{x}) \neq 1$, as $p_{\max}(y = y_{\mathbf{x}}|\mathbf{x}) \equiv 1$, then Heuristic-CALA's adjustment term for class $y_{\mathbf{x}}$ is $p_{\text{tr}}(y = y_{\mathbf{x}})p_{\text{pr}}(y = y_{\mathbf{x}}|\mathbf{x})$, considering both the class proportion and model prediction. This adjustment term is equivalent to that of ALA (Zhao et al. 2022).
- When $K = +\infty$, $p_{\max}(y = c|\mathbf{x}' \in \mathcal{N}_{+\infty}(\mathbf{x}))$ approximates $p_{\text{tr}}(y = c)$. Thus, the adjustment term for the c th class in Heuristic-CALA is $p_{\text{pr}}(y = c)$ which is similar to LA (Menon et al. 2021). The distinction is that Heuristic-CALA in this case employs the class proportions of the predicted labels, which vary with model performance. We have verified that this approach is superior and more rational compared to LA. The comparisons are detailed in Section D.VII of the Appendix.

Meta-CALA Framework

Leveraging the universal approximation capability of deep neural networks, we introduce an adjustment network to estimate the CCPD ratio of samples. The classifier and the adjustment network are alternately updated using a meta-learning-based optimization strategy. Consequently, another GLT method, termed Meta-CALA, is presented. Fig. 3 illustrates the pipeline of the Meta-CALA framework, which consists of four primary components: metadata construction, classifier network, training characteristics extraction, and adjustment network.

We first construct a metadata set with balanced classes and attributes to represent the meta-knowledge of the ground-truth distribution. This metadata set is then utilized to train the adjustment network. To ensure class and attribute balance as much as possible, samples from each class in the training data are clustered into six groups using KMeans with a pre-trained ResNet-50 model (He et al. 2016). We

Dataset	CIFAR10-LT		CIFAR100-LT	
	100:1	10:1	100:1	10:1
Class-Balanced CE (Cui et al. 2019)	72.68%	86.90%	38.77%	57.57%
Class-Balanced Focal (Cui et al. 2019)	74.57%	87.48%	39.60%	57.99%
LDAM-DRW (Cao et al. 2019)	78.12%	88.37%	42.89%	58.78%
De-confound-TDE (Tang et al. 2020)	80.60%	88.50%	44.10%	59.60%
LA (Menon et al. 2021)	77.67%	88.93%	43.89%	58.34%
MiSLAS* (Zhong et al. 2021)	82.10%	90.00%	47.00%	<u>63.20%</u>
LADE (Hong et al. 2021)	81.17%	89.15%	45.42%	61.69%
GLC (Li, Cheung, and Lu 2022)	<u>82.68%</u>	89.81%	<u>48.71%</u>	62.97%
ALA (Zhao et al. 2022)	77.65%	88.32%	43.67%	58.92%
LDAM-DRW-SAFA (Hong et al. 2022)	80.48%	88.94%	46.04%	59.11%
BKD (Zhang et al. 2023)	82.50%	89.50%	46.50%	62.00%
CSA (Shi et al. 2023)	82.53%	90.80%	46.61%	62.60%
Heuristic-CALA (Ours)	83.91%	91.78%	50.53%	64.34%
Meta-Weight-Net (Shu et al. 2019)	73.57%	87.55%	41.61%	58.91%
MetaSAug (Li et al. 2021)	<u>80.54%</u>	<u>89.44%</u>	<u>46.87%</u>	<u>61.73%</u>
Meta-CALA (Ours)	84.79%	92.47%	52.34%	65.51%

Table 1: Accuracy comparison on the CIFAR-LT benchmark. Bold and underlined numbers are the best and second-best results, respectively.

then evenly sample instances from each group and class, as illustrated in the first box of Fig. 3.

Considering that the CCPD ratio can be reflected by the neighborhood information of samples, we extract a series of neighborhood-related training characteristics from the classifier and feed them into the adjustment network to estimate the CCPF ratio $v(x, y)$, thereby obtaining the adjustment vector $\delta_x = [v(x, y_1), \dots, v(x, y_C)]$. The characteristics extraction module is depicted in the third box of Fig. 3. We first extract six basic characteristics that reflect the learning difficulty of samples, including sample loss, logit vector, loss gradient, probability vector, uncertainty which is quantified by the information entropy of the Softmax output, and sample margin. Regarding the logit and probability characteristics, we utilize their values specific to the ground-truth category. Subsequently, we consider the neighborhood extensions of the six basic characteristics. First, we compute the mean values of these characteristics for the samples in the neighborhood. Second, we establish the disparities between the sample’s characteristic values and the neighborhood’s average values. Additionally, we incorporate three other characteristics: 1) the ratio of samples sharing the same label in the neighborhood, 2) the ratio of heterogeneous samples with the highest proportion in the neighborhood, and 3) the cosine distance between the deep feature of the sample and the average feature of the samples in the neighborhood. Furthermore, all characteristics can be extended through the sequence by considering the differences in these training characteristics between the current and previous epochs. In summary, a total number of 42 characteristics are finally extracted. The calculations of all characteristics are detailed in Section C.I of the Appendix.

As the training characteristics are tabular data, we employ a two-layer Multilayer Perceptron as the adjustment network. Moreover, a meta-learning-based learning strategy is proposed to alternatively update the parameters in the classifier \mathbf{W} and the adjustment network Ω , as shown in the second and fourth boxes of Fig. 3. Denote the training data as $\mathcal{D}^{\text{tr}} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ and the metadata as

Dataset	iNat 2018	Places-LT
CE loss	65.76%	30.20%
Decoupling (Kang et al. 2020)	69.49%	37.62%
LA (Menon et al. 2021)	66.36%	34.23%
DisAlign (Zhang et al. 2021)	70.06%	39.30%
MisLAS (Zhong et al. 2021)	71.51%	40.15%
LADE (Hong et al. 2021)	70.00%	38.87%
GCL (Li, Cheung, and Lu 2022)	<u>72.01%</u>	<u>42.64%</u>
LDAM-DRS-SAFA (Hong et al. 2022)	69.78%	41.53%
BKD (Zhang et al. 2023)	71.20%	38.92%
Heuristic-CALA (Ours)	73.23%	43.42%
Meta-Weight-Net (Shu et al. 2019)	67.95%	37.14%
MetaSAug (Li et al. 2021)	<u>68.75%</u>	<u>39.83%</u>
Meta-CALA (Ours)	74.05%	43.97%

Table 2: Accuracy comparison on the iNat 2018 and Places-LT benchmarks.

$\mathcal{D}^{\text{meta}} = \{\mathbf{x}_i^{\text{meta}}, y_i^{\text{meta}}\}_{i=1}^M$. First, a batch of training samples $\{\mathbf{x}_i, y_i\}_{i=1}^n$ is selected, where n is the batch size and the updating of \mathbf{W} is formulated as

$$\hat{\mathbf{W}}^{(t)} \leftarrow \mathbf{W}^{(t)} - \eta_1 \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{W}} \ell_{\text{CALA}} \left(f(\mathbf{x}_i), y_i; \delta_i^{(t)} \right), \quad (9)$$

where η_1 is the step size and δ_i represents the adjustment vector for sample \mathbf{x}_i . Then, the parameters of the adjustment network Ω can be updated on a minibatch of metadata $\{\mathbf{x}_i^{\text{meta}}, y_i^{\text{meta}}\}_{i=1}^m$, with the following formula:

$$\Omega^{(t+1)} \leftarrow \Omega^{(t)} - \eta_2 \frac{1}{m} \sum_{i=1}^m \nabla_{\Omega} \ell_{\text{CE}} \left(f_{\hat{\mathbf{W}}}(\mathbf{x}_i^{\text{meta}}), y_i^{\text{meta}} \right), \quad (10)$$

where m and η_2 are the minibatch size of metadata and the step size, respectively. Finally, the parameters of the classifier are updated using the resulting adjustment terms:

$$\mathbf{W}^{(t+1)} \leftarrow \mathbf{W}^{(t)} - \eta_1 \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{W}} \ell_{\text{CALA}} \left(f(\mathbf{x}_i), y_i; \delta_i^{(t+1)} \right). \quad (11)$$

Utilizing the aforementioned steps, both the classifier and the adjustment network can be effectively optimized.

Experimental Investigation

We evaluate the performance of our methods in addressing class imbalance, attribute imbalance, and their combination across three typical learning scenarios, including LT learning, subpopulation shift learning, and GLT learning. All experiments are repeated three times using different seeds. Due to space constraints, details regarding the comparison methods and datasets are included in Section D of the Appendix.

Experiments for Class Imbalance

Three LT benchmarks are evaluated: CIFAR-LT (Cui et al. 2019), Places-LT (Liu et al. 2019), and iNaturalist (iNat) 2018. The imbalance ratios for the CIFAR-LT benchmark are set to 100:1 and 10:1. For all experiments, we utilize the SGD optimizer with a momentum of 0.9. For CIFAR-LT, we primarily follow Cao et al. (2019) and train all models with a ResNet-32 (He et al. 2016) backbone on a single GPU, employing a multistep learning rate schedule that

Dataset	Waterbirds		CMNIST	
	Avg.	Worst	Avg.	Worst
CORAL (Sun and Saenko 2016)	90.3%	79.8%	71.8%	69.5%
IRM (Arjovsky et al. 2019)	87.5%	75.6%	72.1%	70.3%
GroupDRO (Sagawa et al. 2020)	91.8%	90.6%	72.3%	68.6%
DomainMix (Xu et al. 2020)	76.4%	53.0%	51.4%	48.0%
IB-IRM (Ahuja et al. 2021)	88.5%	76.5%	72.2%	70.7%
V-REx (Krueger et al. 2021)	88.0%	73.6%	71.7%	70.2%
Fish (Shi et al. 2022)	85.6%	64.0%	46.9%	35.6%
LISA (Yao et al. 2022)	91.8%	89.2%	74.0%	73.3%
COSMOS (Chen et al. 2023)	91.7%	89.3%	73.5%	72.4%
PDE (Deng et al. 2024)	92.4%	90.5%	78.1%	75.9%
Heuristic-CALA (Ours)	94.3%	91.8%	79.5%	77.0%

Table 3: Comparison of the average and worst-group accuracy on two subpopulation shift datasets.

reduces the learning rate by a factor of 0.01 at the 160th and 180th epochs. For Places-LT and iNat 2018, we mainly follow Kang et al. (2020) and use the cosine learning rate schedule (Loshchilov and Hutter 2016) to train the ResNet-152 and ResNet-50 backbones, respectively. For the hyperparameters in CALA, the neighborhood size K is selected from $\{20, 40, 60, 80, 100\}$ for all experiments unless noted. τ_1 and τ_2 are set to 1.5 and 1, respectively. The metadata size is 3,000 for CIFAR-LT. For iNat 2018 and Places-LT, one image is selected per class and group to construct the metadata. In Meta-CALA, the adjustment network is optimized using Adam with an initial learning rate of 1×10^{-3} .

Results. Table 1 presents the comparison results on CIFAR-LT, while Table 2 shows the results on the iNat 2018 and Places-LT datasets, with some results sourced from the original papers. The results are divided into two groups based on the usage of meta-learning. Our proposed methods consistently achieve SOTA performance across various datasets and imbalance ratios. Specifically, Heuristic-CALA surpasses the best compared baselines by 1.11% and 1.48% for CIFAR10-LT and CIFAR100-LT, respectively. Moreover, it outperforms the best compared baselines by 1.22% and 0.78% for the iNat 2018 and Places-LT benchmarks, respectively. Notably, Meta-IADA achieves even superior performance compared to Heuristic-CALA, as it leverages the metadata distribution to adjust the model during training. The accuracy of each class for the three methods, including CE loss, LA, and Heuristic-CALA, is compared in Fig. 4. While LA improves the accuracy of the tail classes, it compromises the performance of some head classes. Conversely, our approach demonstrates optimal performance without adversely affecting the head classes. Additionally, we utilize the Wilcoxon signed-rank test to establish the significance of our performance improvement. The obtained p -value of 0.03 signifies a statistically significant enhancement.

The superior performance of CALA compared to other LT learning methods provides evidence of attribute imbalances in LT datasets, an aspect that has usually been overlooked by previous LT baselines. Furthermore, our approach surpasses De-confound-TDE, which utilizes causal intervention during training and counterfactual reasoning during inference, demonstrating the effectiveness of CALA in mitigating spu-

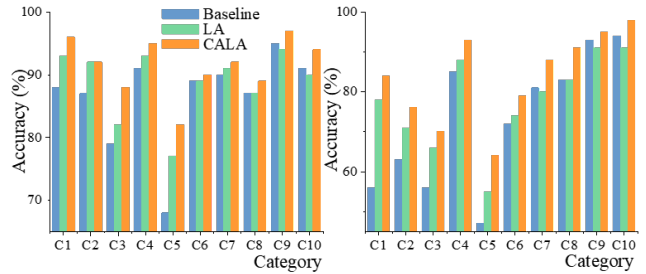


Figure 4: Comparison of class-wise accuracy among baseline (CE loss), LA, and Heuristic-CALA on CIFAR10 with imbalance ratios of 10:1 (left) and 100:1 (right). Moving left to right, the categories progress from tail to head.

rious correlations induced by class and attribute imbalances. We then analyze the distinct characteristics of Heuristic-CALA and Meta-CALA. Although Meta-CALA necessitates an additional metadata set and increases time complexity, it achieves SOTA performance by adjusting model training using a high-quality meta dataset. In contrast, Heuristic-CALA does not require a metadata set and incurs only a marginal increase in computational time compared to the CE loss. Although Heuristic-CALA demonstrates lower performance compared to Meta-CALA, it significantly outperforms existing methods. Detailed comparisons of training times are provided in Section D.VIII of the Appendix.

Experiments for Attribute Imbalance

Three subpopulation shift datasets are adopted: CMNIST (Arjovsky et al. 2019), Waterbirds (Sagawa et al. 2020), and CelebA (Liu et al. 2015), each exhibiting significant attribute imbalances within classes. Taking the Waterbirds dataset as an example, two groups (“land”, “waterbird”) and (“water”, “landbird”) are minority groups. The experimental settings follow Yao et al. (2022), utilizing the ResNet-50 model as the backbone network. Since none of the compared methods rely on meta-learning, we include only Heuristic-CALA in the comparison. We set both τ_1 and τ_2 to 1, and K to 10. Performance is evaluated using both average and worst-group accuracy metrics.

Results. Table 3 presents the comparison results on the Waterbirds and CMNIST datasets, while those for CelebA are provided in the Appendix. Our approach surpasses other invariant learning methods in both average and worst-group accuracy. Specifically, Heuristic-CALA surpasses the second-best baselines by 1.65% in average accuracy and 1.15% in worst-group accuracy. This demonstrates its effectiveness in improving model generalization and enhancing performance for samples with rare attributes.

Experiments for Class & Attribute Imbalance

We consider two GLT benchmarks, ImageNet-GLT and MSCOCO-GLT (Tang et al. 2022). Each benchmark consists of three protocols: CLT, ALT, and GLT, involving changes in the class, attribute, and both class and attribute distributions from training to testing. The experimental settings follow those in Tang et al. (2022), utilizing ResNeXt-50 (Xie

Benchmark	ImageNet-GLT						MSCOCO-GLT					
	CLT		ALT		GLT		CLT		ALT		GLT	
	Acc.	Prec.	Acc.	Prec.	Acc.	Prec.	Acc.	Prec.	Acc.	Prec.	Acc.	Prec.
CE loss	42.52%	47.92%	41.73%	41.74%	34.75%	40.65%	72.34%	76.61%	50.17%	50.94%	63.79%	70.52%
MixUp (Zhang et al. 2018)	38.81%	45.41%	42.11%	42.42%	31.55%	37.44%	74.22%	78.61%	48.90%	49.53%	64.45%	71.13%
LDAM (Cao et al. 2019)	46.74%	46.86%	42.66%	41.80%	38.54%	39.08%	75.57%	77.70%	55.52%	56.21%	67.26%	70.70%
cRT (Kang et al. 2020)	45.92%	45.34%	41.59%	41.43%	37.57%	37.51%	73.64%	75.84%	49.97%	50.37%	64.69%	68.33%
De-confound-TDE (Tang et al. 2020)	45.70%	44.48%	41.40%	42.36%	37.56%	37.00%	73.79%	74.90%	50.76%	51.68%	66.07%	68.20%
BLSoftmax (Ren et al. 2020)	45.79%	46.27%	41.32%	41.37%	37.09%	38.08%	72.64%	75.25%	49.72%	50.65%	64.07%	68.59%
BBN (Zhou et al. 2020)	46.46%	49.86%	43.26%	43.86%	37.91%	41.77%	73.69%	77.35%	51.83%	51.77%	64.48%	70.20%
RandAug (Cubuk et al. 2020)	46.40%	52.13%	46.29%	46.32%	38.24%	44.74%	76.81%	79.88%	53.69%	54.71%	67.71%	72.73%
LA (Menon et al. 2021)	46.53%	45.56%	41.73%	41.74%	37.80%	37.56%	75.50%	76.88%	50.17%	50.94%	66.17%	68.35%
IFL (Tang et al. 2022)	45.97%	52.06%	45.89%	46.42%	37.96%	44.47%	74.31%	78.90%	52.86%	53.49%	65.31%	72.24%
RISDA (Chen et al. 2022)	46.31%	51.24%	43.65%	43.23%	38.45%	42.77%	74.34%	78.27%	51.58%	52.28%	66.85%	71.36%
CSA (Shi et al. 2023)	46.49%	50.77%	43.03%	44.05%	37.22%	42.01%	74.25%	78.56%	52.34%	52.11%	64.78%	69.10%
BKD (Zhang et al. 2023)	46.51%	50.15%	42.17%	41.83%	37.93%	41.50%	75.82%	78.23%	51.88%	51.23%	65.48%	70.59%
Heuristic-CALA (Ours)	54.13%	58.38%	51.88%	52.75%	44.71%	50.82%	79.14%	82.04%	56.67%	57.78%	69.04%	75.51%
MetaSAug (Li et al. 2021)	50.53%	55.21%	49.12%	48.56%	41.27%	47.38%	77.89%	79.45%	54.87%	54.78%	67.83%	73.05%
Meta-CALA (Ours)	55.14%	59.47%	52.76%	53.66%	45.83%	51.36%	80.05%	82.98%	58.99%	59.56%	71.05%	76.21%

Table 4: Comparison of accuracy and precision of the CLT, GLT, and ALT protocols on ImageNet-GLT and MSCOCO-GLT.

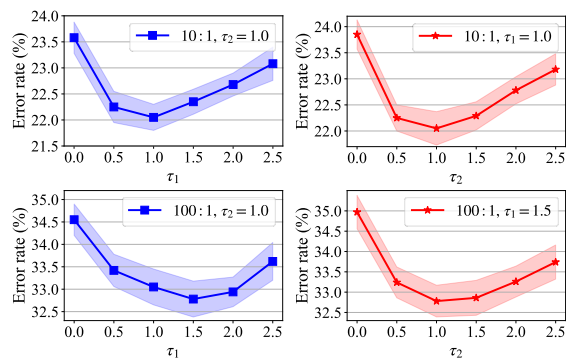


Figure 5: Influence of different τ_1 and τ_2 values on CIFAR-10 with imbalance ratios of 10:1 and 100:1.

et al. 2017) as the backbone network. For Meta-CALA, we optimize the adjustment network using Adam with an initial learning rate of 1×10^{-3} . To construct the metadata, we randomly select two samples per group and class from the training data. Additionally, we set τ_1 and τ_2 to 1.5 and 1 for the CLT protocol. Both τ_1 and τ_2 are set to 1 for the ALT and GLT protocols. We report both accuracy and precision to provide a comprehensive evaluation.

Results. The comparison results for ImageNet-GLT and MSCOCO-GLT are presented in Table 4, with some results sourced from the IFL (Tang et al. 2022) paper. As observed, there is a significant performance decline from the CLT protocol to the GLT protocol, highlighting the challenge posed by attribute imbalance. Heuristic-CALA, which incorporates two adjustment terms to address both class and attribute imbalances, achieves substantial improvements over other methods across all three protocols. Furthermore, Meta-CALA attains SOTA performance by leveraging metadata information to adjust the model during training.

Sensitivity and Ablation Studies

We perform sensitivity analyses on τ_1 and τ_2 , which govern the influence of the two adjustment terms. The results for

Setting	ImageNet-GLT		MSCOCO-GLT	
	Acc.	Prec.	Acc.	Prec.
Heuristic-CALA	44.71%	50.82%	69.04%	75.51%
w/o $u(y)$	42.51%	46.28%	67.49%	72.34%
w/o $v(x, y)$	37.80%	37.56%	66.17%	68.35%

Table 5: Accuracy and precision on the GLT protocol of the ImageNet-GLT and MSCOCO-GLT benchmarks.

Heuristic-CALA are shown in Fig. 5. $\tau_2 = 1$ yields the best performance across different imbalance ratios. For τ_1 , the optimal value is 1 under a 10:1 ratio, while a value of 1.5 is optimal under a 100:1 ratio. These findings suggest that as the class imbalance becomes more pronounced, a larger τ_1 is preferable. Additionally, we perform ablation studies on the CALA loss, considering two settings that remove $u(y)$ and $v(x, y)$ separately. The results, presented in Table 5, indicate that both terms are necessary and crucial for addressing class and attribute imbalances. Furthermore, the role of the CCPD ratio $v(x, y)$ is generally more significant than that of the class priors $u(y)$ under GLT learning.

Conclusion

This study underscores the importance of directly estimating the CCPD ratio between the training and test data in addressing GLT learning. We first introduce a novel logit-adjusted loss function, termed CALA, which incorporates both class priors and the CCPD ratio as adjustment terms. Subsequently, we propose two methods for estimating the CCPD ratio in the CALA loss: a K -neighborhood-based approach and a meta-learning-based approach. These methods give rise to two logit adjustment techniques, Heuristic-CALA and Meta-CALA. Extensive experiments validate the efficacy of our methodologies in addressing both class- and attribute-wise imbalances, achieving SOTA performance across various learning scenarios.

Acknowledgments

This work was mainly conducted by the authors during their tenure at Tianjin University and was partially supported by the NSFC under Grants 6207617 and 62476191.

References

- Agarwal, V.; Shetty, R.; and Fritz, M. 2020. Towards causal vqa: Revealing and reducing spurious correlations by invariant and covariant semantic editing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9687–9695.
- Ahuja, K.; Caballero, E.; Zhang, D.; Gagnon-Audet, J.-C.; Bengio, Y.; Mitliagkas, I.; and Rish, I. 2021. Invariance principle meets information bottleneck for out-of-distribution generalization. In *Proceedings of the Advances in Neural Information Processing Systems*, 3438–3450.
- Arjovsky, M.; Bottou, L.; Gulrajani, I.; and Lopez-Paz, D. 2019. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*.
- Cao, K.; Wei, C.; Gaidon, A.; Arechiga, N.; and Ma, T. 2019. Learning imbalanced datasets with label-distribution-aware margin loss. In *Proceedings of the Advances in Neural Information Processing Systems*, 1567–1578.
- Chen, A. S.; Lee, Y.; Setlur, A.; Levine, S.; and Finn, C. 2023. Confidence-based model selection: When to take shortcuts for subpopulation shifts. *arXiv preprint arXiv:2306.11120*.
- Chen, X.; Zhou, Y.; Wu, D.; Zhang, W.; Zhou, Y.; Li, B.; and Wang, W. 2022. Imagine by reasoning: A reasoning-based implicit semantic data augmentation for long-tailed classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 356–364.
- Cubuk, E. D.; Zoph, B.; Shlens, J.; and Le, Q. V. 2020. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 3008–3017.
- Cui, J.; Zhong, Z.; Liu, S.; Yu, B.; and Jia, J. 2021. Parametric contrastive learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 715–724.
- Cui, Y.; Jia, M.; Lin, T.-Y.; Song, Y.; and Belongie, S. 2019. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9260–9269.
- De Alvis, C.; and Seneviratne, S. 2024. A survey of deep long-tail classification advancements. *arXiv preprint arXiv:2404.15593*.
- Deng, Y.; Yang, Y.; Mirzasoleiman, B.; and Gu, Q. 2024. Robust learning with progressive data expansion against spurious correlation. In *Proceedings of the Advances in Neural Information Processing Systems*, 1390–1402.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778.
- Hong, Y.; Han, S.; Choi, K.; Seo, S.; Kim, B.; and Chang, B. 2021. Disentangling label distribution for long-tailed visual recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 6626–6636.
- Hong, Y.; Zhang, J.; Sun, Z.; and Yan, K. 2022. Safa: Sample-adaptive feature augmentation for long-tailed image classification. In *Proceedings of the European Conference on Computer Vision*, 587–603.
- Kang, B.; Xie, S.; Rohrbach, M.; Yan, Z.; Gordo, A.; Feng, J.; and Kalantidis, Y. 2020. Decoupling representation and classifier for long-tailed recognition. In *Proceedings of the International Conference on Learning Representations*.
- Koh, P. W.; Sagawa, S.; Marklund, H.; Xie, S. M.; Zhang, M.; Balsubramani, A.; Hu, W.; Yasunaga, M.; Phillips, R. L.; Gao, I.; Lee, T.; David, E.; Stavness, I.; Guo, W.; Earnshaw, B.; Haque, I.; Beery, S. M.; Leskovec, J.; Kundaje, A.; Pierson, E.; Levine, S.; Finn, C.; and Liang, P. 2021. Wilds: A benchmark of in-the-wild distribution shifts. In *Proceedings of the International Conference on Machine Learning*, 5637–5664.
- Krueger, D.; Caballero, E.; Jacobsen, J.-H.; Zhang, A.; Binias, J.; Zhang, D.; Le Priol, R.; and Courville, A. 2021. Out-of-distribution generalization via risk extrapolation (rex). In *Proceedings of the International Conference on Machine Learning*, 5815–5826.
- Li, M.; Cheung, Y.-m.; and Lu, Y. 2022. Long-tailed visual recognition via gaussian clouded logit adjustment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6929–6938.
- Li, S.; Gong, K.; Liu, C. H.; Wang, Y.; Qiao, F.; and Cheng, X. 2021. Metasaug: Meta semantic augmentation for long-tailed visual recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5212–5221.
- Liang, W.; and Zou, J. 2022. Metashift: A dataset of datasets for evaluating contextual distribution shifts and training conflicts. In *Proceedings of the International Conference on Learning Representations*.
- Lin, C.; Tsai, C.-F.; and Lin, W.-C. 2023. Towards hybrid over- and under-sampling combination methods for class imbalanced datasets: an experimental study. *Artificial Intelligence Review*, 56(2): 845–863.
- Liu, Z.; Luo, P.; Wang, X.; and Tang, X. 2015. Deep learning face attributes in the wild. In *Proceedings of the IEEE International Conference on Computer Vision*, 3730–3738.
- Liu, Z.; Miao, Z.; Zhan, X.; Wang, J.; Gong, B.; and Yu, S. X. 2019. Large-scale long-tailed recognition in an open world. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2537–2546.
- Loshchilov, I.; and Hutter, F. 2016. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*.
- Mao, R.; Fan, W.; and Li, Q. 2023. GCARE: Mitigating subgroup unfairness in graph condensation through adversarial regularization. *Applied Sciences*, 13(16): 9166.

- Menon, A. K.; Jayasumana, S.; Rawat, A. S.; Jain, H.; Veit, A.; and Kumar, S. 2021. Long-tail learning via logit adjustment. In *Proceedings of the International Conference on Learning Representations*.
- Ren, J.; Yu, C.; Sheng, S.; Ma, X.; Zhao, H.; Yi, S.; and Li, H. 2020. Balanced meta-softmax for long-tailed visual recognition. In *Proceedings of the Advances in Neural Information Processing Systems*, 4175–4186.
- Sagawa, S.; Koh, P. W.; Hashimoto, T. B.; and Liang, P. 2020. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *Proceedings of the International Conference on Learning Representations*.
- Shi, J.-X.; Wei, T.; Xiang, Y.; and Li, Y.-F. 2023. How re-sampling helps for long-tail learning? In *Proceedings of the Advances in Neural Information Processing Systems*, 75669–75687.
- Shi, Y.; Seely, J.; Torr, P. H.; Siddharth, N.; Hannun, A.; Usunier, N.; and Synnaeve, G. 2022. Gradient matching for domain generalization. In *Proceedings of the International Conference on Learning Representations*.
- Shu, J.; Xie, Q.; Yi, L.; Zhao, Q.; Zhou, S.; Xu, Z.; and Meng, D. 2019. Meta-weight-net: Learning an explicit mapping for sample weighting. In *Proceedings of the Advances in Neural Information Processing Systems*, 1919–1930.
- Srivastava, M.; Hashimoto, T.; and Liang, P. 2020. Robustness to spurious correlations via human annotations. In *Proceedings of the International Conference on Machine Learning*, 9046–9056.
- Sun, B.; and Saenko, K. 2016. Deep coral: Correlation alignment for deep domain adaptation. In *Proceedings of the European Conference on Computer Vision Workshops*, 443–450.
- Tang, K.; Tao, M.; Qi, J.; Liu, Z.; and Zhang, H. 2022. Invariant feature learning for generalized long-tailed classification. In *Proceedings of the European Conference on Computer Vision*, 709–726.
- Tao, Y.; Sun, J.; Yang, H.; Chen, L.; Wang, X.; Yang, W.; Du, D.; and Zheng, M. 2023. Local and global logit adjustments for long-tailed learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 11783–11792.
- Tripathi, A.; Chakraborty, R.; and Kopparapu, S. K. 2021. A novel adaptive minority oversampling technique for improved classification in data imbalanced scenarios. In *Proceedings of the International Conference on Pattern Recognition*, 10650–10657.
- Wan, M.; Zha, D.; Liu, N.; and Zou, N. 2023. In-processing modeling techniques for machine learning fairness: A survey. *ACM Transactions on Knowledge Discovery from Data*, 17(3): 1–27.
- Wang, X.; Lian, L.; Miao, Z.; Liu, Z.; and Yu, S. X. 2020. Long-tailed recognition by routing diverse distribution-aware experts. *arXiv preprint arXiv:2010.01809*.
- Wang, Z.; Xu, Q.; Yang, Z.; He, Y.; Cao, X.; and Huang, Q. 2024. A unified generalization analysis of re-weighting and logit-adjustment for imbalanced learning. In *Proceedings of the Advances in Neural Information Processing Systems*, 48417–48430.
- Xiang, L.; Ding, G.; and Han, J. 2020. Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification. In *Proceedings of the European Conference on Computer Vision*, 247–263.
- Xie, S.; Girshick, R.; Dollár, P.; Tu, Z.; and He, K. 2017. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1492–1500.
- Xu, M.; Zhang, J.; Ni, B.; Li, T.; Wang, C.; Tian, Q.; and Zhang, W. 2020. Adversarial domain adaptation with domain mixup. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 6502–6509.
- Yan, Y. T.; Wu, Z. B.; Du, X. Q.; Chen, J.; Zhao, S.; and Zhang, Y. P. 2019. A three-way decision ensemble method for imbalanced data oversampling. *International Journal of Approximate Reasoning*, 107: 1–16.
- Yao, H.; Wang, Y.; Li, S.; Zhang, L.; Liang, W.; Zou, J.; and Finn, C. 2022. Improving out-of-distribution robustness via selective augmentation. In *Proceedings of the International Conference on Machine Learning*, 25407–25437.
- Zhang, H.; Cisse, M.; Dauphin, Y. N.; and Lopez-Paz, D. 2018. Mixup: Beyond empirical risk minimization. In *Proceedings of the International Conference on Learning Representations*.
- Zhang, S.; Chen, C.; Hu, X.; and Peng, S. 2023. Balanced knowledge distillation for long-tailed learning. *Neurocomputing*, 527: 36–46.
- Zhang, S.; Li, Z.; Yan, S.; He, X.; and Sun, J. 2021. Distribution alignment: A unified framework for long-tail visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2361–2370.
- Zhao, Y.; Chen, W.; Tan, X.; Huang, K.; and Zhu, J. 2022. Adaptive logit adjustment loss for long-tailed visual recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 3472–3480.
- Zhong, Z.; Cui, J.; Liu, S.; and Jia, J. 2021. Improving calibration for long-tailed recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 16489–16498.
- Zhou, B.; Cui, Q.; Wei, X.; and Chen, Z. 2020. Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9719–9728.
- Zhou, X.; and Wu, O. 2023. Implicit counterfactual data augmentation for deep neural networks. *arXiv preprint arXiv:2304.13431*.
- Zhou, X.; Yang, N.; and Wu, O. 2023. Combining adversaries with anti-adversaries in training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 11435–11442.
- Zhou, X.; Ye, W.; Lee, Z.; Xie, R.; and Zhang, S. 2024. Boosting model resilience via implicit adversarial data augmentation. In *Proceedings of the International Joint Conference on Artificial Intelligence*, 5653–5661.