

Dynamic Loss For Robust Learning

Shenwang Jiang[†], Jianan Li^{†,*}, Jizhou Zhang, Ying Wang, and Tingfa Xu^{*}

Abstract—Label noise and class imbalance are common challenges encountered in real-world datasets. Existing approaches for robust learning often focus on addressing either label noise or class imbalance individually, resulting in suboptimal performance when both biases are present. To bridge this gap, this work introduces a novel meta-learning-based dynamic loss that adapts the objective functions during the training process to effectively learn a classifier from long-tailed noisy data. Specifically, our dynamic loss consists of two components: a label corrector and a margin generator. The label corrector is responsible for correcting noisy labels, while the margin generator generates per-class classification margins by capturing the underlying data distribution and the learning state of the classifier. In addition, we employ a hierarchical sampling strategy that enriches a small amount of unbiased metadata with diverse and challenging samples. This enables the joint optimization of the two components in the dynamic loss through meta-learning, allowing the classifier to effectively adapt to clean and balanced test data. Extensive experiments conducted on multiple real-world and synthetic datasets with various types of data biases, including CIFAR-10/100, Animal-10N, ImageNet-LT, and Webvision, demonstrate that our method achieves state-of-the-art accuracy.

Index Terms—Robust learning, label noise, class imbalance, meta learning.

1 INTRODUCTION

DEEP neural networks (DNNs) have demonstrated remarkable success attributed to the abundance of labeled data [1], [2], [3]. However, real-world datasets often exhibit long-tailed distributions and inevitably contain noisy labels [4]. The presence of such biased data distributions renders DNNs susceptible to overlooking tail classes [5] and memorizing noisy training labels [6], leading to suboptimal performance on balanced and clean test data. Consequently, addressing the challenge of robust learning from long-tailed data with noisy labels has received growing attention in recent studies [7], [8].

Recent advancements in robust learning [6], [9], [10], [11] have explored various techniques, including correcting noisy labels [12], [13] and adjusting classification margins [5], [14], to address the issues of label noise and class imbalance, respectively. However, these solutions heavily rely on manually-designed rules to distinguish between noisy and clean samples, or pre-set parameters to obtain prior knowledge of class distribution. When label noise meets class imbalance, manual interventions become impractical since tail classes make it difficult to identify label noise, while noisy labels make observed class distribution unreliable. As a result, we aim to pave a new path for robust learning from biased data in a fully self-adaptive manner.

Taking the aforementioned challenges into consideration, we propose a novel meta-learning-based dynamic loss

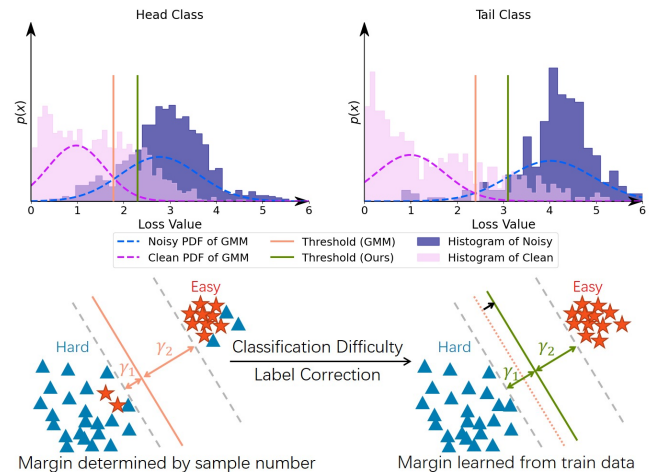


Fig. 1: **Upper:** Conventional class-specific Gaussian Mixture Model (GMM) with fixed hyperparameters tends to misclassify clean samples into the noisy split, whereas our label corrector learns optimal division thresholds for different classes using training data and samples' class-specific loss rank. **Lower:** Previous methods generate classification margin (γ) based solely on sample number, whereas our margin generator takes into account the presence of noisy labels and the distinct classification difficulties of different classes, resulting in a more appropriate margin.

that comprises of a learnable *label corrector* and a *margin generator*. Our approach aims to learn a classifier robustly from long-tailed noisy data. In contrast to the predefined fixed objective functions adopted by priors, our dynamic loss learns to correct per-sample labels and adjust per-class classification margins simultaneously by perceiving the underlying data distribution and the learning state of the classifier. As a result, our approach provides suitable dynamic learning objectives throughout the training process.

The main challenge in correcting noisy labels lies in identifying noisy samples and disregarding or reusing them

This work was financially supported by the National Natural Science Foundation of China (No. 62101032), the Postdoctoral Science Foundation of China (Nos. 2021M690015, 2022T150050), and Beijing Institute of Technology Research Fund Program for Young Scholars (No. 3040011182111).

S. Jiang, J. Li, J. Zhang, Y. Wang, and T. Xu are with Beijing Institute of Technology, Beijing 100081, China. E-mail: jiangwenj02@gmail.com, {lijianan, zhangjizhou, 3120215325, ciom_xtf1}@bit.edu.cn. T. Xu is also with Big Data and Artificial Intelligence Laboratory, Beijing Institute of Technology Chongqing Innovation Center, Chongqing 401135, China.

[†] Shenwang Jiang and Jianan Li contributes equally to this work.

^{*} Jianan Li and Tingfa Xu are co-corresponding authors.

The code will present in https://github.com/jiangwenj02/dynamic_loss.

appropriately. Previous solutions [13], [15] have used a GMM with fixed hyperparameters to differentiate between noisy and clean samples based on their loss values, which is inadequate for long-tailed noisy data. This is because some clean samples of tail classes may have higher loss values than noisy samples of head classes due to the highly skewed class distribution. While a class-specific GMM [16] is a straightforward solution, it is still prone to confusion due to the significant variation in noise rates among different classes (Figure 1). Furthermore, previous methods have relabeled identified noisy samples by combining their assigned labels and the classifier’s predictions in a fixed manner. This approach can lead to incorrect corrections since the classifier is imprecise at the beginning of training and susceptible to overfitting to noisy samples towards the end. To address these issues, we propose a novel label corrector that can learn to jointly identify and relabel noisy samples in a fully learnable fashion. This is achieved by taking into account both the class-specific loss rank of samples and the learning state of the classifier.

In order to adjust per-class classification margins, we are inspired by the fact that classes with fewer samples are associated with larger generalization error bounds, which can be minimized by increasing the classification margin [5]. However, most prior approaches [14] pre-define a fixed margin for each class based solely on its sample number (Figure 1). This approach suffers from two main drawbacks: i) the per-class sample number becomes unreliable in the presence of noisy labels; and ii) the distinct classification difficulties among different classes are simply ignored. To address these limitations, we propose a novel margin generator that can learn to produce a suitable dynamic margin for each class by self-perceiving the true class distribution underlying the noisy data, as well as the classification difficulty of each class.

To enable the learning of the classifier from long-tailed noisy data, we propose a unified dynamic loss that integrates the label corrector and the margin generator, and optimize them through meta-learning. This approach allows the objective function for the classifier to be dynamically adjusted throughout the training process. The label corrector and the margin generator work in tandem to improve the learning from noisy data. Specifically, the label corrector restores the true distribution of the data, which allows the margin generator to produce a more suitable classification margin. In turn, the resulting improved margin boosts the accuracy of the predicted labels and enhances the reliability of the label corrector. Furthermore, the convergence of meta-learning is widely recognized as a challenging issue in the field [17]. In order to tackle this problem, we propose the incorporation of a group optimization strategy and the explicit utilization of known information. These techniques serve to streamline the meta-learning task, mitigate input instability, and ultimately enhance the convergence of the meta-learning process.

To collect a small amount of unbiased meta data for meta-learning, we have developed a new hierarchical sampling strategy that progressively builds a random primary set and then a balanced clean meta set. This approach ensures that the meta set is enriched with diverse and challenging samples that better simulate the distribution

of real test data, thereby circumventing the problem of over-reliance on straightforward samples. By doing so, our dynamic loss can guide the classifier learning robustly on various types of biased data in a fully self-adaptive manner.

We conduct a comprehensive evaluation of our proposed dynamic loss on both synthetic and real-world long-tailed data with label noise, achieving state-of-the-art results on a wide variety of benchmarks featuring various imbalance ratios and noise rates, such as CIFAR-10/100 [18], Animal-10N [19], ImageNet-LT [11], and Webvision [4]. Furthermore, additional tests conducted on purely imbalanced or noisy data further validate the dynamic loss’s exceptional adaptability and robustness.

In summary, our work contributes in the following ways:

- We introduce a straightforward yet powerful dynamic loss, which paves a new way for robust learning on various forms of biased data in a fully self-adaptive manner.
- We devise a novel hierarchical sampling strategy that efficiently generates diverse and unbiased meta data, which allows for better simulation of the true data distribution and improves the generalization ability of the model.
- We achieve state-of-the-art performance on multiple synthetic and real-world datasets, demonstrating the effectiveness and versatility of our proposed approach.

2 RELATED WORKS

Long-Tailed Learning. Previous works on long-tailed learning can be broadly categorized into three main approaches: data re-sampling [20], [21], boundary adjustment [22], and re-weighting [9], [23], [24]. The data re-sampling approach involves balancing the class distribution by over-sampling the tail classes. However, this method is prone to overfitting on the tail classes. Methods belonging to the second category aim to enlarge the classification boundary of the tail classes while narrowing that of the head classes. This is achieved by modifying the classification threshold [25] or by adjusting the weights of the output layer through normalization [10]. The re-weighting approach aims to assign larger loss weights to the tail classes. Conventional approaches of this category [24], [26] impose weights on each training sample directly, which may cause unstable training due to sensitivity to outliers [14]. Recent works [27] modify the predicted scores in the Softmax function to achieve re-weighting, which yields more stable training and promising performance. In this work, we adapt the re-weighting strategy to more challenging long-tailed scenarios with label noise.

Learning under Label Noise. There are two main categories of methods for learning under label noise: sample re-weighting and relabeling. The re-weighting strategy involves assigning lower weights to samples with larger loss values, which are considered to be noisy [28], [29]. MentorNet [6] learns data-driven curriculums for deep convolutional neural networks trained on corrupted labels. Meta-Weight-Net [9] learns an explicit weighting function directly from a small set of clean data. MetaSeg [30] directly generates weights according to the feature and given label of

images. On the other hand, the relabeling strategy leverages noisy samples by refining their labels. Bootstrapping [31] integrates assigned labels and model predictions through interpolation. Some works divide clean and noisy samples based on priors learned from a manually generated noisy set [32] and then take advantage of noisy samples [33].

Long-tailed Learning under Label Noise. Recently, there have been several approaches proposed for addressing long-tailed learning with label noise. For instance, HAR [8] applies a data-dependent regularization technique to regularize different regions of the input space differently. CurveNet [7] learns to assign appropriate weights to different samples based on their loss curves. ROLT [34] combines DivideMix and LDAM to correct noisy labels and improve the performance of tail classes. However, unlike these methods, the approach presented in this work is the first to simultaneously correct noisy labels and adaptively adjust per-class classification margins in a learnable and adaptive manner according to the training data.

3 METHODS

In this section, we will provide a detailed description of our dynamic loss and the process of optimizing it through meta-learning.

3.1 Overview

Our objective is to train a classifier f_ω with learnable parameters ω using a noisy and imbalanced training set $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$, where each training example comprises an image x_i and its corresponding one-hot class label y_i . Although the training set is beset by label noise and class imbalance, our aim is to ensure that the classifier accurately recognizes all classes. To achieve this goal, we use a balanced and clean test set.

We optimize the learnable parameters ω by minimizing the classification loss on the training set, given by

$$\omega^* = \arg \min_{\omega} \sum_{i=1}^N \ell(y_i, f_\omega(x_i)). \quad (1)$$

where $\ell(\cdot)$ represents the cross-entropy loss. However, due to the presence of label noise and class imbalance in the training set, optimizing using this naive cross-entropy loss suffers from two significant drawbacks. Firstly, the labels assigned to noisy samples do not correspond to their ground-truths, leading to high loss values and causing the model to memorize the noisy labels. Secondly, tail classes occur much less frequently than head classes, but have the same classification margins, making them susceptible to poor generalization.

To address the aforementioned issues, we introduce a novel dynamic loss function that corrects noisy labels and adjusts the classification margins for different classes in an adaptive and learnable manner. The dynamic loss is given by:

$$\mathcal{DL} = \ell(y_i^*, f_\omega(x_i) + q), \quad (2)$$

where y_i^* and $q \in \mathbb{R}^C$ represent the reassigned label and the additive classification margin for x_i , respectively.

As depicted in Figure 2, the dynamic loss incorporates a learnable label corrector \mathcal{G}_{θ_l} , parameterized by θ_l ,

and a margin generator \mathcal{G}_{θ_m} , parameterized by θ_m , which are responsible for correcting per-sample labels and per-class classification margins, respectively. We jointly optimize these components along with the classifier f_ω using meta-learning. In the following, we elaborate on the two components and their optimization in detail.

3.2 Label Corrector

The label corrector operates in a class-wise manner to detect and correct wrongly assigned labels in noisy samples. To identify such samples, it first divides all samples into C groups based on their class and sorts the samples in each group individually according to their loss values. It then divides the sorted samples in each group into R bins of equal size. To determine whether a bin $r \in \{1, \dots, R\}$ is dominated by noisy or clean samples, the label corrector uses a lightweight class-wise meta network. As a result, the bin index r_i of the loss value for sample i can be used as a reliable indicator for identifying label noise.

In terms of label correction, the classifier can learn from the clean samples that dominate the data and transfer this knowledge to the noisy samples, provided that the classifier has not severely over-fitted on biased data. This allows the classifier's predictions on noisy samples to be more accurate and closer to their ground-truth labels, facilitating the correction of wrongly assigned labels.

Drawing from the aforementioned observation, we present our label corrector which reassigns a ground-truth label y_i^* for a given sample x_i by calculating a weighted sum of its assigned label y_i and the prediction y'_i made by the classifier. The weights used for the summation are dependent on the loss bin index r_i , given by:

$$\begin{aligned} y_i^* &= \mathcal{G}_{\theta_l}(y_i, y'_i, r_i) \\ &= y_i * g(r_i|y_i) + y'_i * (1 - g(r_i|y_i)), \end{aligned} \quad (3)$$

The function $g : r_i|y_i \rightarrow [0, 1]$ is a class-dependent weighting function that maps the bin index r_i to a balance weight and is learned by a small meta network. The network is comprised of a one-hot encoder and a two-layer perceptron (MLP) with a Sigmoid activation function.

In instances where sample i is regarded as noisy due to a high loss value resulting in a large bin index r_i , the computed value of $g(r_i|y_i)$ tends to approach zero. Consequently, the label corrector adjusts its label assignment by incorporating the prediction made by the classifier to rectify its initial erroneous label assignment. The reverse is also true: if the sample has a low loss value and a small bin index, the weight assigned to the classifier's prediction is close to zero, allowing the assigned label to remain dominant in the label correction process.

3.3 Margin Generator

To design the margin generator, we revisit the Label-Distribution-Aware Margin Loss (LDAM) [5] from the perspective of generalization error bound. Given that tail classes often have fewer training samples, they typically exhibit larger generalization error bounds when compared with head classes. Since the generalization error bound is often negatively correlated with the magnitude of the

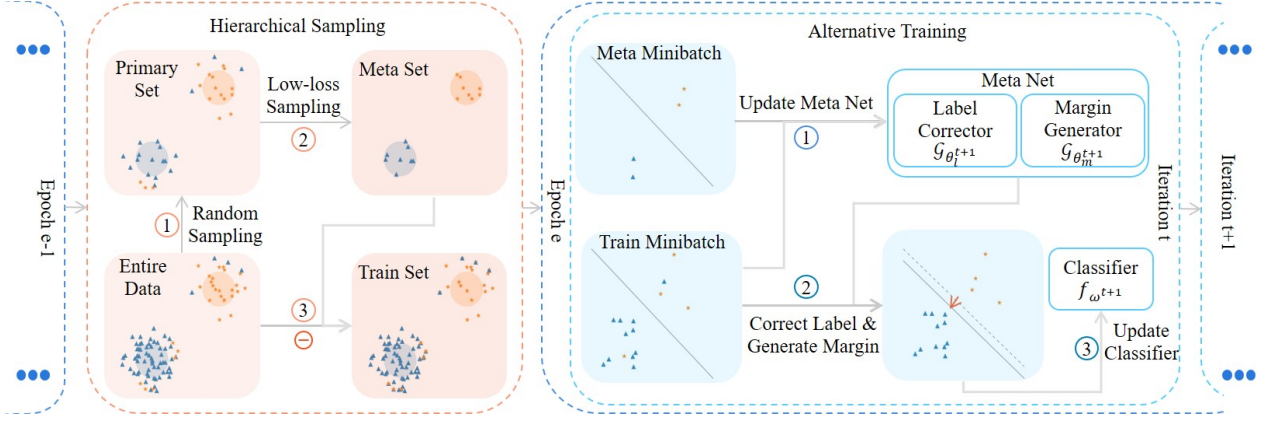


Fig. 2: Overview of our learning paradigm using the dynamic loss. Each training epoch involves splitting the entire data into an unbiased meta set and a biased training set. In each iteration, we jointly update the label corrector and margin generator using meta-learning on mini-batch meta and training data, followed by updating the classifier on mini-batch training data by minimizing the dynamic loss using corrected per-sample labels and generated per-class margins.

classification margin, increasing the classification margins for the tail classes can effectively reduce their generalization error bounds.

In this regard, Balanced Meta-Softmax [14] presents an unbiased extension of standard Softmax by adjusting the classification margin for class j based on its sample number n_j , and adding the margin $-\log(n_j)$ to the confidence score p_j predicted by the classifier:

$$\ell'(\omega|y=j) = -\log\left(\frac{e^{p_j+\log(n_j)}}{\sum_{i=1}^C e^{p_i+\log(n_i)}}\right). \quad (4)$$

However, when dealing with long-tailed data with noisy labels, the sample number n_j may not accurately reflect the true number of samples belonging to class j due to the existence of label noise. Furthermore, manually defining the margin solely based on the sample number may not account for the distinct classification difficulties among different classes.

We propose a learnable margin generator \mathcal{G}_{θ_m} , which consists of a two-layer multi-layer perceptron (MLP) that can dynamically adjust the margin for each class. During classifier training, the margin generator optimizes a learnable margin vector $\mathbf{q} \in \mathbb{R}^C$, which is initialized with an all-ones vector:

$$\mathbf{q} = \mathcal{G}_{\theta_m}(\mathbf{1}) = [q_1, \dots, q_C]. \quad (5)$$

By integrating the margin vector into the standard Softmax loss, we have the modified loss function:

$$\hat{\ell}(\omega|y=j) = -\log\left(\frac{e^{p_j+q_j}}{\sum_{i=1}^C e^{p_i+q_i}}\right). \quad (6)$$

where p_j and q_j denote the predicted score and the learned margin for class j , respectively. Since the classification margin is $-q_j$ in our formulation, the learned margin q_j for class j should be positively correlated with its sample number.

Hence the margin generator can automatically adjust per-class margins by adapting to the underlying true class distribution in long-tailed noisy data, as well as the classification difficulty of each class. Importantly, the margin generator can accomplish this in a learnable manner, without requiring any manual intervention or prior information.

3.4 Hierarchical Sampling Strategy

We combine the label corrector \mathcal{G}_{θ_l} and margin generator \mathcal{G}_{θ_m} into a unified meta net \mathcal{G}_{θ} , which is a crucial element of our dynamic loss. We employ meta-learning to optimize \mathcal{G}_{θ} and facilitate the learning of the classifier \mathbf{f}_{ω} to better adjust to balanced and clean test data.

To enable meta-learning, a meta set $\mathcal{D}_c = \{(x_i, y_i)\}_{i=1}^M$ containing a small number of balanced and clean data must be constructed. A simple approach is to select M_1 samples with the lowest classification loss values from each class in \mathcal{D} , as samples with lower loss values computed by \mathbf{f}_{ω} are more likely to have correctly assigned labels. However, this approach can be problematic because easier samples often have lower loss values during training, which can result in fixed easy samples being selected at each epoch and increase the risk of overfitting to such samples.

Therefore, we have devised a hierarchical sampling strategy to create \mathcal{D}_c . This strategy involves a two-step process: first, we randomly select M_0 samples from each class in \mathcal{D} to construct a *primary set*; second, we choose M_1 samples with low loss from each class in the primary set to form the final *meta set*. The samples that are not selected for the meta set are included in the counterpart set \mathcal{D}_n . Refer to Fig. 2 for an illustration of this process.

The advantages of incorporating the primary set are two-fold. Firstly, since the samples in the primary set are randomly selected at each epoch, it ensures that the resulting meta set is distinct across different epochs. Secondly, the primary set contains fewer samples than \mathcal{D} , increasing the likelihood of selecting hard yet clean samples located near the decision boundary into the meta set. This hierarchical sampling strategy guarantees the dynamism and diversity of the meta set, preventing the model from overfitting to biased data.

3.5 Optimization

The convergence of meta-learning poses significant challenges [17]. Previous approaches [7], [9], [15] employ a strategy of inputting all known information, particularly given labels, predicted labels and loss value, into the meta net,

TABLE 1: Detailed settings for training on different datasets.

Settings		CIFAR-10	CIFAR-100	Webvision	Animal-10N	ImageNet-LT
Classifier	Optimizer	SGD				
	Momentum	0.9				
	Weight Decay	5e-4	5e-4	1e-4	5e-4	1e-4
	Learning Rate	0.1	0.1	0.02	0.02	0.1
	Learning Scheduler	Cosine Annealing				
MetaNet	Optimizer	Adam				
	Weight Decay	0				
	Learning Rate	3e-3				
	Learning Scheduler	Fixed				
Others	M0	0.5	0.5	0.5	0.5	-
	M1	0.25	0.25	0.25	0.25	-
	Epoch	300	300	150	100	90 / 400
	Warmup Epoch	5	5	1	5	0
	Batch Size	512	512	64	128	128
	Rank bins	100	50	100	100	-

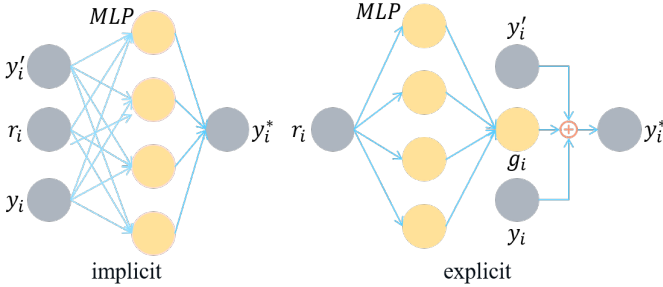


Fig. 3: Visualization of implicit and explicit inputs. y_i , y'_i , y_i^* and r_i represent the assigned label of the training set, predicted label of the classifier, reassigned label of our meta net, and the bin index of the loss rank, respectively.

generating customized weights for each sample. However, due to computational complexity, meta nets often exhibit a relatively simple network structure consisting of a few fully connected layers, limiting their capability to handle complex tasks. Notably, in scenarios involving biased data, a substantial portion of the known information, such as predicted and given labels, can be omitted from the meta net’s input. This is because only one of these labels needs to be selected as the sample’s label, and a weight can be learned to combine the two labels and obtain a new label. Directly incorporating these known information into the meta net not only increases the optimization difficulty but also diminishes the utilization of the information. This is due to the implicit nature of using these information, requiring the meta net to extract it from the input, as expressed by:

$$y_i^* = \mathcal{G}_{\theta_i}(y_i, y'_i, r_i) = g(y_i, y'_i, r_i). \quad (7)$$

To enhance the utilization of known information, we opt for explicit utilization through our Label Corrector, as illustrated in Figure 3 and governed by Equation 3.

The customization of weights for individual samples diminishes the impact of diverse input information on the meta net. However, input information often encompasses errors, such as varying sample loss values during training and

overlapping loss values between noisy and clean samples. These factors make it challenging to discern between them solely based on loss values. To overcome this challenge, we employ a grouping strategy that involves sorting by loss value for noisy data and given labels for imbalanced data. This strategy brings stability to the input information, reducing its instability and enhancing the influence of each input on the meta net. By explicitly utilizing known information and implementing the grouping strategy, we simplify the task of the meta net, promoting easier convergence and the assignment of accurate weights and margins to each sample group.

4 EXPERIMENTS ON LONG-TAILED NOISY DATA

4.1 Experiments on CIFAR-N-LT

Dataset and Implement Details. We assess the performance of our method on CIFAR-N-LT dataset, which comprises CIFAR-10 and CIFAR-100 [18]. These datasets contain 60,000 RGB images, out of which 50,000 images are used for training and 10,000 images for testing. The images are evenly distributed across 10 and 100 categories, respectively. Additionally, the datasets are subjected to simulated label noise and class imbalance.

To simulate a long-tailed dataset, we follow the exponential profile proposed in [5], where the imbalance ratio ρ leads to an exponential decay in the sample number across different classes. We then inject label noise into the long-tailed dataset to create the training set, with each sample’s label independently changed to class j with probability $\frac{N_j}{N} \lambda$, where N is the total number of training samples, N_j is the frequency of class j , and λ represents the noise rate.

Following ROLT [34], we consider imbalance ratios of $\rho \in \{10, 100\}$ and noise rates of $\lambda \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$. The ResNet-32 [2] is adopted as the classifier, and we train all classifiers using balanced-softmax for 300 epochs with a batch size of 512. The learning rate is initialized as 0.1 and controlled by a cosine annealing learning scheduler [37]. We train all classifiers using the same SGD optimizer with a momentum of 0.9 and a weight decay of 5e-4. Further training details can be found in Table 1.

TABLE 2: Accuracy (%) on CIFAR-10-N-LT with varying imbalance ratio and noise rate.

Imbalance Ratio		10					100					Avg.
Noise Rate		0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5	
Cross Entropy	Last	76.63	68.89	61.21	54.01	44.09	59.99	50.56	44.81	37.13	30.02	52.73
	Best	77.14	74.35	72.22	68.28	59.65	63.45	53.45	49.49	45.53	40.76	60.43
DivideMix [13]	Last	88.11	86.95	66.42	73.78	73.85	63.08	61.14	24.31	24.13	20.09	58.19
	Best	88.50	87.13	67.67	74.26	74.15	63.08	61.61	33.79	31.06	26.10	60.74
ELR+ [35]	Last	85.21	85.71	84.13	81.16	78.91	64.76	61.41	56.42	51.31	44.05	69.31
	Best	87.04	86.14	84.38	83.30	80.23	65.73	62.12	56.42	52.26	45.73	70.34
MOIT+ [36]	Last	86.67	84.60	84.85	81.34	79.02	66.21	61.63	60.01	37.71	39.01	68.11
	Best	87.00	84.84	84.85	81.95	79.52	66.92	62.16	60.49	37.85	39.19	68.48
Balanced-Softmax [14]	Last	87.52	84.00	79.49	73.53	65.27	76.14	69.37	62.03	52.79	46.44	69.66
	Best	87.62	84.50	82.84	79.09	75.90	78.63	73.29	70.89	67.57	63.40	76.37
DivideMix +Balanced-Softmax	Last	88.23	86.96	78.02	79.17	76.89	76.63	75.72	46.82	53.25	54.33	71.60
	Best	88.93	86.06	79.19	79.81	78.20	79.36	76.55	48.60	54.99	55.84	72.75
FaMUS [15]	Last	77.31	78.68	79.15	81.17	77.57	55.46	51.67	49.40	38.21	26.69	61.53
	Best	83.12	84.34	84.85	85.44	84.12	57.50	54.61	54.30	45.35	32.93	66.66
CurveNet [7]	Last	84.10	81.70	78.47	78.73	75.65	65.77	66.21	62.37	48.71	51.85	69.36
	Best	84.87	84.62	79.98	81.33	78.37	67.55	68.72	63.71	51.63	52.84	71.36
HAR [8]	Last	86.46	84.27	81.78	79.55	78.07	78.60	75.05	72.08	65.48	63.90	76.52
	Best	87.03	84.47	81.94	79.87	78.25	79.02	76.14	72.74	67.22	65.00	77.17
Dynamic Loss (Ours)	Last	89.23	88.39	86.58	84.43	83.34	77.80	76.31	74.10	69.64	67.45	79.73
	Best	89.44	88.46	86.72	84.73	83.71	78.96	76.64	76.17	70.37	70.26	80.55

Main Results. Table 2 and 3 present the average accuracy on CIFAR-10-N-LT and CIFAR-100-N-LT with varying imbalance ratios and noise rates. Our method exhibits consistently high accuracy across a wide range of biases, whereas previous methods suffer rapid degradation. Specifically, our dynamic loss improves the average last accuracy by 3.21% and 5.79% compared to HAR on CIFAR-10-N-LT and CIFAR-100-N-LT, respectively. Additionally, our method significantly outperforms the baseline model that simply combines strategies from DivideMix and Balanced-Softmax.

Furthermore, the performance of our last model is generally very close to that of our best model despite varying bias settings. In contrast, the last models of DivideMix and Balanced-Softmax degrade significantly compared to their corresponding best models, especially on CIFAR-10-N-LT with severe imbalance and noise (e.g., $\rho = 100$ and $\lambda = 0.5$). These results suggest that our method is much more resistant to overfitting on biased data than the aforementioned priors.

It is worth noting that previous methods require carefully tuned hyperparameters based on unobservable noise rate [13] or perform two-stage training to obtain prior information on class distribution [8]. In comparison, our dynamic loss employs a fixed set of hyperparameters and requires only one-round end-to-end training without manual interventions on the same dataset.

4.2 Experiments on Webvision.

Dataset and Implement Details. We also evaluate our dynamic loss on the WebVision dataset [4], which is a large-scale real-world dataset that suffers from label noise and class imbalance. It comprises 2.4 million images, of which approximately 20% are mislabeled [41]. To construct the

miniWebVision dataset, we follow the methodology proposed in MentorNet [6] by selecting the top 50 classes, resulting in an observed imbalance ratio of around 6.78. Following priors [13], we train the Inception-ResNet V2 [42] for 150 epochs using the SGD optimizer with momentum 0.9 and weight decay of $1e-4$. For ResNet-50, the training details is following INOLML [40]. During the warm-up stage, which lasts for one epoch, we use a batch size of 64 and an initial learning rate of 0.02. The learning rate is then adjusted using a cosine annealing learning scheduler. We evaluate the model’s performance on the validation sets of both WebVision and ImageNet [43]. Additional training details are available in Table 1.

Main Results. Table 4 displays the performance of our method on the WebVision and ImageNet validation sets. Our approach outperforms other state-of-the-art methods even though most of them utilize additional model co-training and ensembling techniques. Notably, our method surpasses HAR and ROLT+, which are specifically designed to handle long-tailed noisy data, by at least 2.48% in terms of accuracy on WebVision, demonstrating its superiority.

The rapid advancements in the field of self-supervised learning have highlighted the increasing importance of adaptive methods for learning from biased data. Particularly, methods that can effectively adapt to models with varying initial parameters are of great significance. In this study, we performed experiments on ResNet-50 that was pre-trained using self-supervised techniques, such as CLIP [44]. The results, presented in the last row of Table 4, demonstrate a significant improvement of 3.92% in accuracy, achieved through our proposed approach, compared to the model trained using cross entropy with the same initial parameters pretrained by CLIP. These findings provide evidence of the dynamic adaptability of our method to classifiers with different initial parameters.

TABLE 3: Accuracy (%) on CIFAR-100-N-LT with varying imbalance ratio and noise rate. NC: not converging, NA: not available.

Imbalance Ratio		10					100					Average
Noise Rate		0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5	
Cross Entropy	Last	43.48	37.25	31.34	25.53	19.45	29.92	21.85	19.32	13.71	12.21	25.41
	Best	43.95	37.64	32.18	29.44	23.87	30.52	22.13	19.58	14.59	12.81	26.67
DivideMix [13]	Last	54.17	51.92	50.44	45.02	43.43	36.31	35.68	34.10	33.19	27.22	41.15
	Best	54.94	53.35	50.93	45.36	43.44	36.99	36.24	34.87	33.64	27.74	41.75
ELR+ [35]	Last	52.48	51.30	46.24	39.98	34.91	33.01	28.10	24.92	22.11	16.54	34.96
	Best	53.91	51.90	47.88	42.61	37.35	33.81	28.94	26.10	22.11	17.39	36.20
MOIT+ [36]	Last	44.66	40.12	NC	NC	NC	NC	NC	NC	NC	NC	NA
	Best	47.70	44.94	42.10	39.12	35.50	32.66	30.35	28.99	25.99	22.82	35.02
Balanced-Softmax [14]	Last	58.38	54.59	50.49	44.83	40.45	43.17	38.67	33.27	27.08	22.10	41.30
	Best	58.62	54.73	50.66	45.63	40.56	43.50	38.67	33.62	28.05	24.19	41.82
DivideMix +Balanced-Softmax	Last	56.37	54.80	54.83	51.29	48.89	43.26	42.42	40.46	37.83	30.95	46.11
	Best	56.85	56.06	55.64	52.30	50.01	43.67	42.79	40.99	38.50	32.38	46.92
FaMUS [15]	Last	46.07	51.59	46.07	46.93	43.83	29.33	30.22	28.53	27.83	24.57	30.72
	Best	47.03	52.05	46.41	47.88	44.30	29.66	30.31	28.50	27.24	24.85	30.81
CurveNet [7]	Last	50.41	47.14	43.18	41.23	34.85	22.10	20.44	17.80	11.87	9.24	29.83
	Best	52.73	51.93	47.56	44.08	39.74	25.26	21.35	18.72	13.60	12.20	32.72
HAR [8]	Last	58.88	55.43	52.57	46.01	43.96	42.67	39.39	34.43	29.43	24.94	42.77
	Best	59.32	55.80	53.44	46.75	44.61	44.45	40.98	36.09	31.17	27.15	43.98
Dynamic Loss (Ours)	Last	59.24	57.57	56.85	52.07	50.74	47.23	45.74	42.72	39.58	33.87	48.56
	Best	59.52	57.85	57.32	52.66	51.26	47.55	45.82	43.54	39.96	34.30	48.98

TABLE 4: Accuracy (%) on WebVision and ImageNet validation sets. * denotes the use of model cotraining or ensembling. † indicates the backbone is pretrained using self-supervised techniques (CLIP). IRV2: Inception-ResNet V2.

Methods	Backbone	Webvision		ImageNet	
		top 1	top 5	top 1	top 5
Cross Entropy	IRV2	72.48	88.48	65.08	87.88
Cross Entropy	ResNet-50	71.72	87.84	65.96	86.88
Cross Entropy	ResNet-50†	78.04	93.04	72.84	91.68
HAR [8]	IRV2	75.50	90.70	70.30	90.00
DivideMix* [13]	IRV2	77.32	91.64	75.20	90.84
ROIT+* [34]	IRV2	77.64	92.44	74.64	92.48
ELR+* [35]	IRV2	77.78	91.68	70.29	89.76
CMW-Net-SL* [38]	IRV2	78.08	92.96	75.72	92.52
FaMUS* [15]	IRV2	79.40	92.80	77.00	92.76
MOIT+ [36]	IRV2	78.76	-	-	-
NCR [39]	ResNet-50	80.50	-	-	-
INOLML* [40]	ResNet-50	81.70	93.80	78.10	92.90
Dynamic Loss	IRV2	80.12	93.64	74.76	93.08
Dynamic Loss	ResNet-50	78.56	92.52	72.08	91.36
Dynamic Loss	ResNet-50†	81.96	94.48	77.43	93.40

5 EXPERIMENTS ON NOISY DATA

5.1 Experiments on CIFAR-N

Dataset and Implement Details. CIFAR-N is a synthetic noisy dataset derived from CIFAR. It includes two common types of simulated label noise: symmetric and asymmetric. Symmetric noise is introduced by randomly changing the labels with all possible labels based on a fixed probability of λ (also known as noise rate). Asymmetric noise, on the other hand, is designed manually to mimic real-world label

noise, where labels are only altered with those in similar classes, such as deer \rightarrow horse and dog \leftrightarrow cat. We evaluate the performance of our dynamic loss in handling noisy data using CIFAR-10-N and CIFAR-100-N with symmetric noise rates of 0.2, 0.4, 0.6 and asymmetric noise rates of 0.2, 0.4. We utilize PreAct ResNet (PAREs18) [49] as the classifier, following the approach proposed in DivideMix [13]. Additional training details can be found in Table 1.

Results on CIFAR-N. Table 5 illustrates that our proposed method outperforms previous methods specifically designed for learning on noisy data, achieving the highest average accuracy. Unlike DivideMix, which requires manual tuning of hyperparameters under different noise types and rates, our dynamic loss adapts well to various noisy scenarios in a fully self-adaptive manner without any manual intervention.

5.2 Experiments on Animal-10N

Dataset and Implement Details. We also evaluate our method on the real-world noisy dataset ANIMAL-10N [19], which contains a total of 55,000 images (50,000 for training and 5,000 for testing) of 5 pairs of confusing animals. All image categories are equally represented in the dataset. The images are collected from websites using the search keywords as labels, resulting in significant label noise with an estimated rate of 8%. To ensure a fair comparison with prior work [19], we train the VGG10-BN [1] for 100 epochs using the SGD optimizer with a momentum of 0.9 and weight decay of $5e-4$. The warm-up stage lasts for 5 epochs, and the batch size is set to 128. We initialize the learning rate to 0.02 and control it with a cosine annealing learning scheduler. Further training details are presented in Table 1.

Main Results. Table 6 illustrates that our dynamic loss achieves state-of-the-art performance compared to all pre-

TABLE 5: Accuracy (%) on CIFAR-N with varying noise rate. PAREs18 used as the classifier.

Datasets	CIFAR-10-N						CIFAR-100-N			
Noise Rate	20	40	60	20 (Asym.)	40 (Asym.)	Avg.	20	40	60	Avg.
Cross Entropy	86.98	77.52	73.63	83.60	77.85	79.92	60.38	46.92	31.82	46.37
SELFIE [19]	86.39	82.23	74.81	-	-	-	55.71	51.14	43.85	50.23
PLC [45]	86.40	71.72	65.22	90.23	85.40	79.79	59.66	49.24	33.18	47.36
NCT [46]	95.00	87.00	73.22	91.51	93.00	87.95	67.65	57.97	45.01	56.88
Coteaching [47]	93.83	91.74	57.65	93.23	90.78	85.45	70.81	62.65	41.55	58.34
CMW-Net [38]	91.09	86.91	83.33	93.02	92.70	89.41	70.11	65.84	56.93	64.29
DivdeMix [13]	95.63	93.78	94.23	94.18	92.73	94.11	77.20	73.37	70.75	73.77
GJS [48]	94.20	92.80	89.72	91.92	86.07	90.94	73.31	71.33	66.92	70.52
NCR [39]	95.20	94.50	78.45	-	90.70	-	76.60	74.20	38.25	63.02
MOIT+ [36]	94.08	91.95	89.38	94.50	93.27	92.64	75.89	70.88	65.30	70.69
Dynamic Loss (Ours)	95.90	94.69	92.28	95.74	94.51	94.62	78.26	75.28	69.18	74.24

TABLE 6: Accuracy (%) on Animal-10N.

Methods	Cross Entropy	SELFIE [19]	PLC [45]	NCT [46]	Co-teaching [47]	CMW-Net [38]	DivdeMix [13]	GJS [48]	Dynamic Loss
Accuracy	79.40	81.80	83.40	84.10	80.20	84.70	84.00	84.17	86.54

vious methods on the ANIMAL-10N dataset. Specifically, our method outperforms the classifier trained with cross entropy by a large margin of 7.14%, demonstrating its superior ability to handle real-world noisy data.

6 EXPERIMENTS ON LONG-TAILED DATA

6.1 Experiments on CIFAR-LT

Dataset and Implement Details. The CIFAR-LT is a simulated long-tailed dataset that is derived from CIFAR by reducing the number of training samples per class according to an exponential function $n_i = n\mu^i$, where n_i , i , and n denote the number of samples in the i -th class, the class index, and the maximum number of samples across all classes, respectively. We evaluate the efficacy of our dynamic loss approach in dealing with long-tailed data on clean CIFAR datasets that have varying imbalance ratios ($\rho \in \{10, 20, 50, 100\}$). For more detailed information on the training process, please refer to Table 1.

Results Table 7 presents the superior performance of our proposed method compared to previous approaches specifically designed for learning on long-tailed data. Notably, our method surpasses LDAM, which adjusts classification margins based solely on the sample number, by a significant margin of 4.03% and 5.75% on CIFAR-10-LT and CIFAR-100-LT, respectively. This result demonstrates the effectiveness of our dynamic loss in accurately perceiving the classification difficulty of different classes and adaptively adjusting their margins accordingly.

Furthermore, our approach seamlessly integrates with Logit Adjustment and Balanced Softmax techniques by utilizing them as base margins while learning the residual margin through our margin generator. As shown in Table 7, our proposed approach with Balanced-softmax and logit adjustment yields a notable improvement in final accuracy, namely 0.19% and 3.05% on CIFAR-10-LT, compared to the initial Balanced-softmax and logit adjustment techniques.

However, the performance level remains comparable to that of our original dynamic loss. These findings indicate that our original dynamic loss is sufficient for accurately identifying the classification complexity of diverse categories and accounting for the number of samples per class.

6.2 Experiments on ImageNet-LT

Dataset and Implement Details. The ImageNet-LT dataset consists of 115.8K images, categorized into 1,000 classes based on the Pareto distribution. Consequently, the number of images per class varies between 5 and 1280. Following previous studies [11], [55], we employ several architectures for training, including ResNet-10, ResNet-50, ResNet-152, ResNeXt-50 [56], and self-supervised ResNet-50. These models are trained for either 90 or 400 epochs, utilizing the SGD optimizer with a momentum of 0.9 and weight decay of $1e-4$. No warm-up stage is conducted, and a batch size of 128 is employed. The learning rate is initialized at 0.1 and controlled using a cosine annealing learning scheduler. Further details regarding the training procedure are provided in Table 1.

Main Results Table 8 presents the experimental findings, showcasing the superior accuracy achieved by our proposed method when employed in conjunction with various classifiers, namely ResNet-50, ResNet-152, ResNeXt-50, and self-supervised ResNet-50. Remarkably, even when utilizing the ResNet-10 classifier trained from scratch, our method surpasses the performance of the PaCo framework [55] equipped with self-supervised techniques by 2.12%. These results underscore the robustness and efficacy of our proposed method in effectively handling the challenges posed by long-tailed datasets.

Additionally, we conducted supplementary experiments by incorporating our method with the pretrained ResNet50 using self-supervised techniques. The results reveal a significant performance improvement of 2.42% compared to the PaCo framework. This further highlights the compatibility

TABLE 7: Accuracy (%) on CIFAR-LT with varying imbalance ratio. ResNet32 used as the classifier.

Datasets	CIFAR-10-LT					CIFAR-100-LT				
Imbalance Ratio	10	20	50	100	Avg.	10	20	50	100	Avg.
Cross entropy	86.39	82.23	74.81	70.36	78.45	55.71	51.14	43.85	38.32	47.26
Focal Loss [23]	86.66	82.76	76.71	70.38	79.13	55.78	51.95	44.32	38.41	47.62
CB Focal [50]	87.49	84.36	79.27	74.57	81.42	57.99	52.59	45.32	39.6	48.88
LDAM-DRW [5]	87.68	85.51	81.64	78.02	83.21	44.70	52.93	48.22	59.59	51.36
FaMUS [15]	87.9	86.24	83.32	80.96	84.61	59.00	55.95	49.93	46.03	52.73
Balanced-Softmax [14]	91.01	88.85	86.44	82.31	87.15	64.00	59.48	54.36	50.47	57.08
WD [51]	89.80	84.81	79.66	74.84	82.28	61.60	52.75	45.89	40.79	50.26
MiSLAS [52]	90.00	88.52	85.70	82.10	86.58	63.20	59.25	52.30	47.00	55.44
Logit Adjustment [53]	89.64	86.77	82.61	78.38	84.35	62.83	58.81	52.15	48.36	55.54
CMO [54]	83.26	89.27	87.19	85.35	86.27	62.30	60.12	51.40	46.60	55.11
Dynamic Loss	91.24	88.30	86.46	82.95	87.24	63.99	59.79	54.51	50.14	57.11
Dynamic Loss + Balanced-Softmax	90.99	89.66	85.49	83.21	87.34	64.18	60.13	54.60	50.54	57.36
Dynamic Loss + Logit Adjustment	91.10	89.54	85.52	83.43	87.40	64.15	60.08	54.14	49.65	57.01

TABLE 8: Accuracy (%) on Imagenet-LT. * indicates the use of self-supervised techniques. RN: ResNet. RNeXt: ResNeXt.

Methods	Epoch	RN-10	RN-50	RN-152	RNeXt-50
Cross Entropy	90	34.01	44.60	46.20	42.78
Focal Loss [23]	90	32.64	41.61	44.36	41.58
LDAM-DRW [5]	90	36.03	48.80	51.83	51.43
CDB-S [57]	90	37.70	41.80	46.40	45.10
Logit adjustment [53]	90	38.43	48.89	47.86	51.85
Balanced-Softmax [14]	90	38.21	50.96	53.93	51.73
MiSLAS [52]	400	44.36	53.05	48.82	51.88
PaCo* [55]	400	42.89	57.00	58.55	58.20
Dynamic Loss	90	38.87	51.16	54.40	52.37
Dynamic Loss	400	45.01	53.19	56.41	53.48
Dynamic Loss*	400	-	59.42	-	-

of our proposed method with self-supervised approaches, leading to enhanced performance gains.

7 QUALITATIVE ANALYSIS.

7.1 Label corrector

Behavior of Label Corrector. We examine the behavior of the label corrector on the balanced CIFAR-10-N with asymmetric noise, which is designed to emulate the structure of real-world label noise by assigning distinct noise rates to various classes. Figure 4 illustrates the learned weight $g(r|y)$ by the label corrector and the percentage of noisy labels corresponding to an increasing loss bin index r for each class. We observe that for the classes that contain noisy labels, clean samples mainly appear in the top-ranked (low-loss) bins while noisy samples occupy the bottom-ranked (high-loss) bins. This finding supports our hypothesis that the loss bin index r can be used as a reliable input indicator for the label corrector to differentiate between noisy and clean samples. Accordingly, the generated weight $g(r|y)$ remains at 1 and suddenly drops to 0 at around bin 60, indicating that the label corrector preserves the assigned ground-truth label for clean samples and resorts to the predicted label that is more likely to be the ground-truth for noisy samples. On

the other hand, for the classes without noisy labels, $g(r|y)$ remains at 1. As a result, our label corrector consistently outputs the correct labels for both noisy and clean samples across different classes.

We also visualize the learned label weights of the real-world datasets miniWebvision and Animal-10n in Figures 5 and 6, respectively. For miniWebvision, we select 2 categories at intervals of 25 to visualize their learned label weights considering the large number of categories. As shown in Figure 5, the learned label weights vary with different classes, suggesting that the noise rates of different classes are different, which is consistent with real-world datasets. In Animal-10N, two categories were selected at intervals of 5 to visualize the learned label weights. As depicted in Figure 6, the learned label weights remain 1 and drop to 0 at around bin 92 (red dotted line), indicating that the noise rate estimated by the label corrector is about 8%, which is consistent with the well-recognized estimated noise rate on Animal-10N [19].

Corrected Label Accuracy. Figure 7 illustrates the accuracy of the corrected labels on balanced CIFAR-10-N, measured as the proportion of samples with ground-truth labels after the label correction process. The label accuracy gradually increases as the training progresses, and the classifier becomes more reliable. Eventually, the label accuracy reaches over 90% on CIFAR-10-N with noise rates of 0.2 and 0.4. Remarkably, the accuracy has also significantly improved by 35% under a noise rate of 0.6. The high accuracy of the corrected labels validates our design choices from two perspectives. Firstly, this finding supports our assumption that the classifier primarily focuses on fitting in the dominant clean samples, and can transfer the acquired knowledge to the noisy samples for predicting their ground-truth labels. Secondly, the label corrector can precisely identify the noisy samples and rectify their labels with the predicted correct ones.

Learnable Weight Varied Trend in Training. Figure 8 illustrates the evolution of the learnable weight across training epochs. At the beginning, the label corrector mainly relies on the given labels to train the classifier, while gradually shifting towards trusting the predicted labels for samples

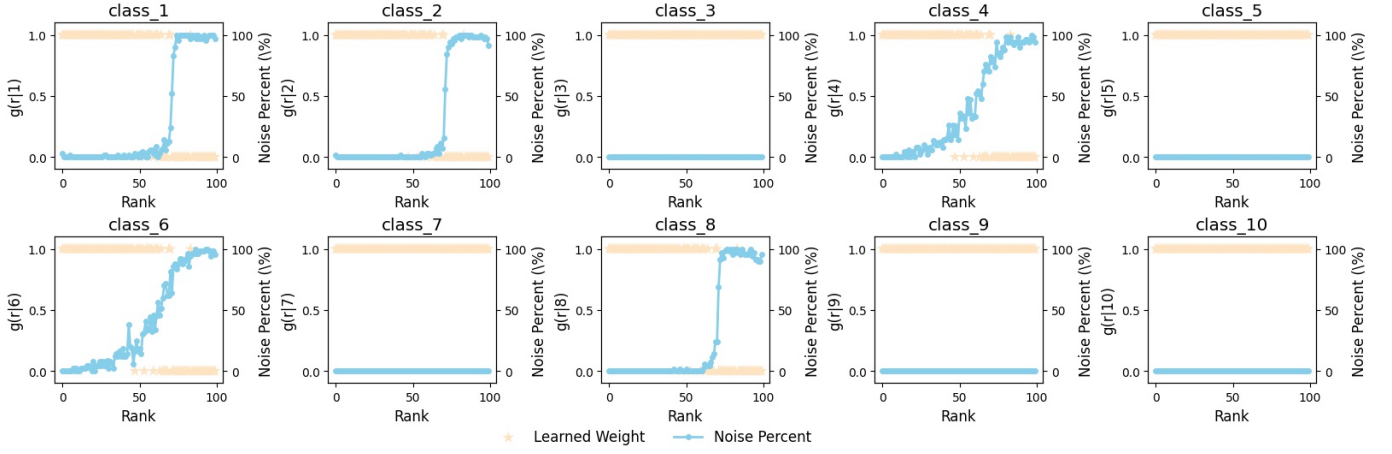


Fig. 4: Visualizing the learned label weights $g(r|y)$ and noise percentage for each class on CIFAR-10-N (Asym. $\lambda = 0.4$).

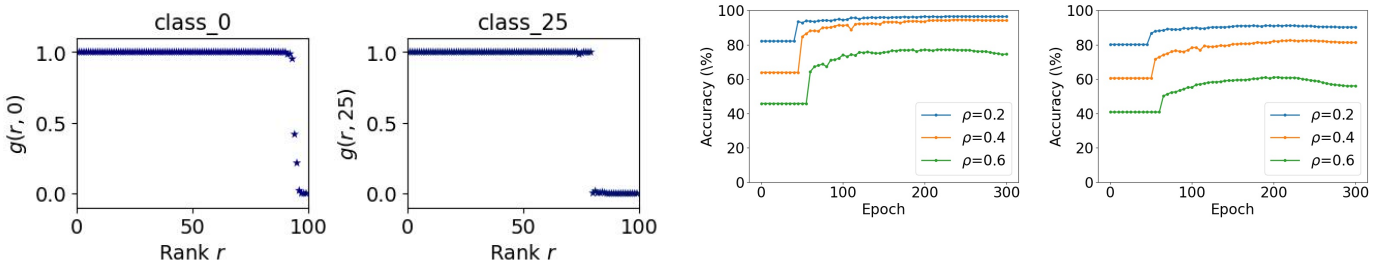


Fig. 5: The visualization of the per-class learned label weights on WebVision.

Fig. 7: The visualization of the accuracy of generated labels varied with epoch of \mathcal{G}_{θ_i} on CIFAR-10-N (left) and CIFAR-100-N (right) with noise rates λ ranging from 0.2 to 0.6.

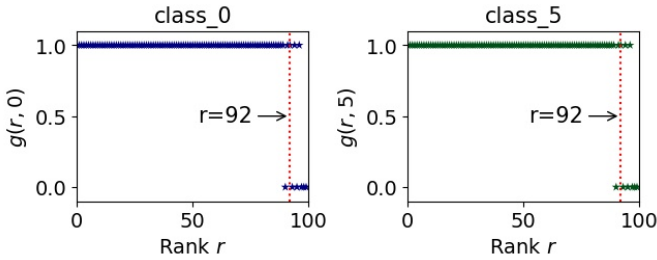


Fig. 6: The visualization of the per-class learned label weights on Animal-10N.

with higher rank bins, as indicated by the decreasing values of the learnable weight at epoch 45. Furthermore, the label corrector accurately estimates the noise rate to be approximately 35% (for a noise rate of 40%, there are actually 35% noisy samples). The plot, together with Figure 7, reveals that the epoch at which the label corrector begins to trust the classifier is delayed as the noise rate increases. This suggests that the label corrector considers the classifier to require more training epochs to produce more reliable predicted labels as the noise rate of the training set increases. Therefore, our results demonstrate that the label corrector is capable of dynamically adjusting the relabelling of noisy labels based on the status of the classifier and the training set.

7.2 Margin generator

Behavior of Margin Generator. We conducted an analysis of the behavior of the margin generator on clean CIFAR-10-LT with imbalance factor $\rho = 20$. The left subfigure of Figure 9 shows the generated margins for different classes. We observed that as the class index increases and the sample number decreases, the learned margin also decreases as expected. This suggests that the margin generator has the ability to automatically discern the sample numbers of different classes and adaptively adjust the margins for each class accordingly.

Interestingly, we also observed some irregularly larger margins on class 9 and 10. To explain this observation, we visualized the feature distribution of the meta set using T-SNE [58] as shown in the right subfigure of Figure 9. We found that the feature distribution of these two classes correspond to the two rightmost clusters, indicating that they are easier to be distinguished from the other classes. This evidence supports the claim that the margin generator takes into account not only the sample number but also the classification difficulty of each class when generating comprehensively adaptive margins during classifier training.

We present visualizations of the learned label weights for the real-world miniWebvision and simulated ImageNet-LT datasets in Figures 10 and 11, respectively. In Figure 10, the generated margins over different classes in miniWebvision accord with the complex variation of sample size, which demonstrates the adaptability of our method to handle com-

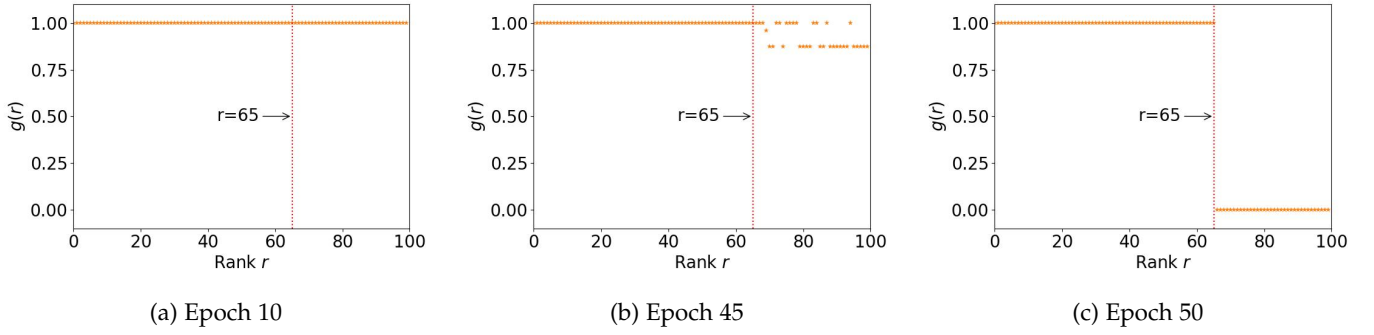


Fig. 8: The visualization of the learnable weight $g(r)$ of class 1 varied with training epochs on CIFAR-10-N with noise rate of $\lambda=0.4$.

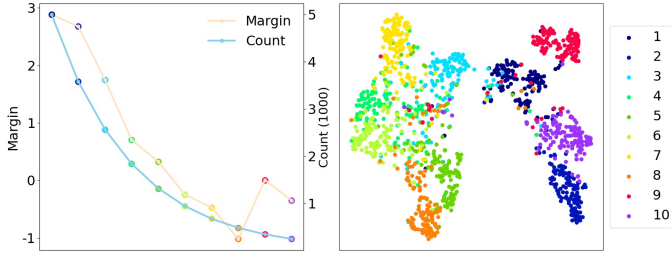


Fig. 9: The learned per-class margins and the number of samples (left), and the feature distribution of samples in the meta set (right) on CIFAR-10-LT with $\rho=20$, are presented.

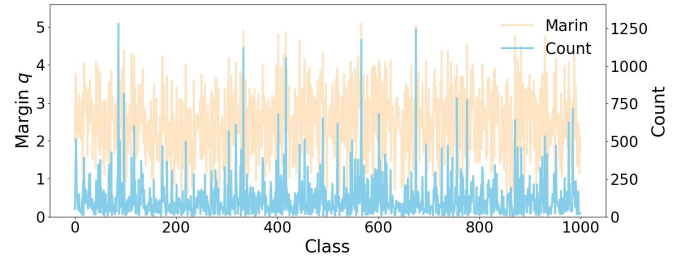


Fig. 11: Visualization of the learned class margins on ImageNet-LT.

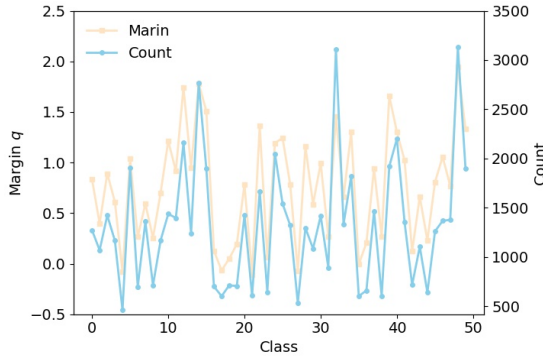


Fig. 10: Visualization of the learned per-class margins on Webvision.

plex real-world biased data. For the simulated ImageNet-LT dataset, we visualize the learned margins for 333 categories in intervals of 3 in Figure 11. As shown in the figure, the learned margins consistently vary with the sample number of different classes, suggesting that the margin generator can generate proper margins for different classes. This finding provides evidence that our method can handle datasets with numerous categories.

Learnable Weight Varied with Imbalance Ratios. Figure 12 depicts the learned margins generated by the margin generator under different imbalance ratios. We can observe that as the class index increases, corresponding to decreasing sample size, the generated margin consistently decreases, irrespective of the varying imbalance ratios. Additionally, the

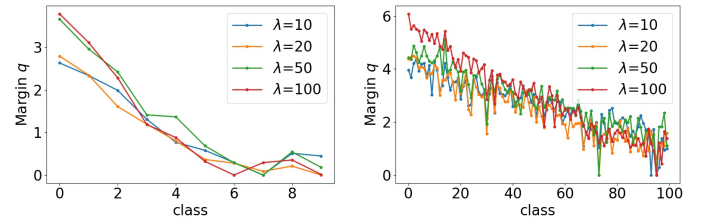


Fig. 12: The visualization of the learned class-aware margin of $G_l(\theta_l)$ on CIFAR-10-LT (left) and CIFAR-100-LT (right) with imbalance factor ρ ranging from 10 to 100.

variation in learned margins across different classes tends to increase as the imbalance ratio becomes more severe. Both quantitative and qualitative analyses provide evidence that the margin generator effectively respects and adapts to various class distributions by automatically learning to assign appropriate margins.

Learnable Weight Varied in Training. In Figure 13, the variations of classification margins and feature distributions of meta set are depicted across different training epochs. It can be inferred that the margin generator continually adapts the classification margins based on the feature distributions of meta set during the training process. For instance, in the case of class 10, a relatively small classification margin is assigned to it since it is challenging to recognize at epoch 50. However, as it becomes easier to recognize at epoch 250, the margin generator assigns a larger classification margin to it. This observation supports the claim that the margin generator can dynamically regulate the classification margin based on the classification difficulty.

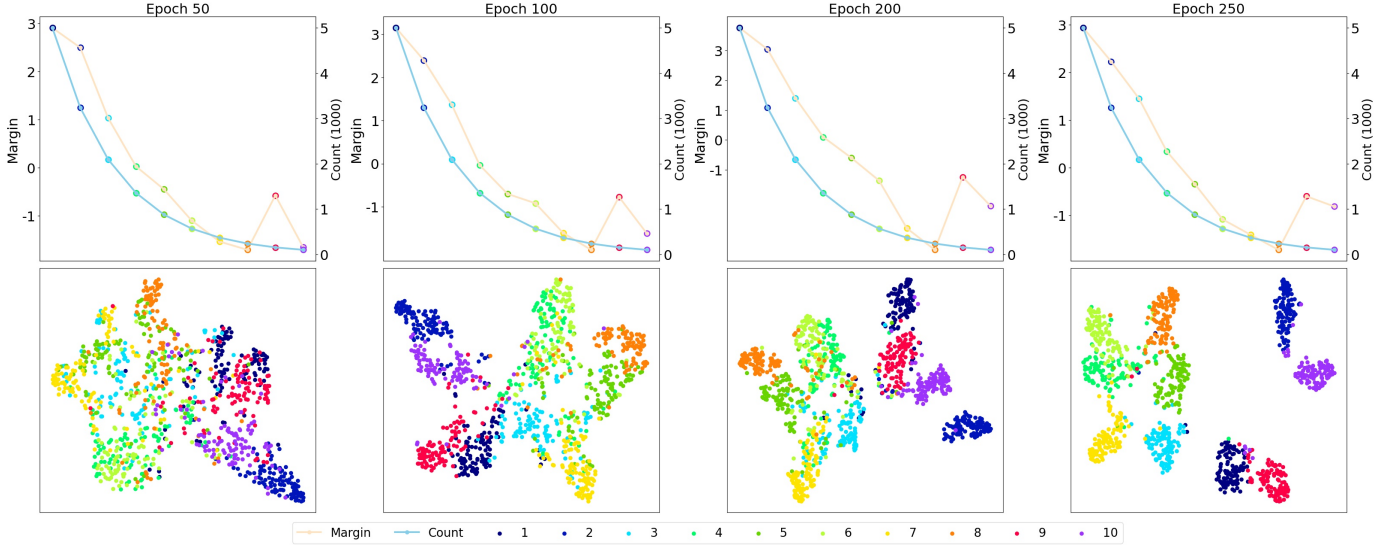


Fig. 13: Visualization of learned classification margins and feature distributions of meta set varied with training epochs.

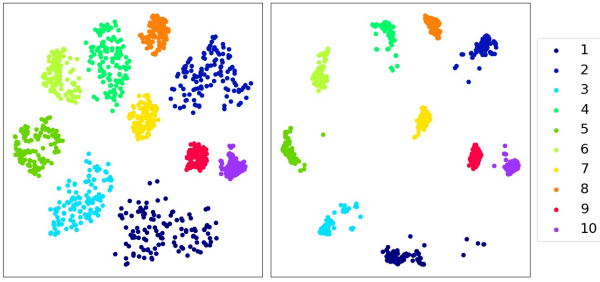


Fig. 14: The feature distribution of samples in the meta set, constructed using hierarchical sampling (left) and naive sampling (right) on CIFAR-10-N-LT with $\rho=10$ and $\lambda=0.2$.

7.3 Behavior of hierarchical sampling.

We examine the efficacy of our hierarchical sampling strategy in enhancing the construction of meta sets. Figure 14 visually presents the feature distribution of the meta set generated using our hierarchical sampling method (left) and the meta set generated using a naive sampling method (right). The feature distribution of samples selected using hierarchical sampling demonstrates greater dispersion within distinct clusters compared to naive sampling. This indicates that the creation of a primary set prior to random sampling enables a more diverse meta data selection, encompassing both easy and challenging samples. As a result, this approach helps to mitigate biased learning on easier samples.

8 ABLATION STUDIES.

Effect of Label Corrector. We construct a model variant wherein the label corrector component is removed to assess its efficacy. Table 9a illustrates that the average final accuracy of this variant decreases by 8.22% compared to the complete dynamic loss configuration. This notable decrease in performance provides compelling evidence of the effectiveness of the label corrector.

Effect of Margin Generator. As indicated in Table 9a, to assess the efficacy of the margin generator, we initially

utilize only the label corrector, which yields an average final accuracy of 71.08%. Subsequently, we incorporate Balanced-Softmax to address class imbalance, which only results in a marginal improvement of 1.56% in accuracy. Finally, by implementing our margin generator, we are able to further improve the accuracy to 79.73%. These results serve as evidence supporting the importance of utilizing a dynamic margin in managing long-tailed noisy data.

Effect of Hierarchical Sampling. Table 9b demonstrates that the substitution of hierarchical sampling with naive random sampling leads to a reduction of up to 1.20% in average final accuracy. This observation suggests that the meta set created using hierarchical sampling possesses a more comparable distribution to the test set.

Effect of Class-specific Label Corrector. To validate the class-specific design of our label corrector, we constructed a class-agnostic variant and evaluated its performance on CIFAR-10-N, a dataset containing 40% asymmetric noise with varying noise rates across different classes. Our results demonstrate that the class-specific label corrector significantly outperforms its class-agnostic counterpart by a large margin of 3.95% in terms of accuracy (94.51% vs. 90.56%), thus providing clear evidence for the effectiveness of our class-specific design.

Effect of Explicit Utilizing Given and Predicted labels. To assess the impact of explicitly utilizing known information, particularly given labels and predicted labels, we conducted a comparative analysis between the implicit and explicit approaches in utilizing such information. The results, as presented in Table 9c, demonstrate that the explicit utilization of known information resulted in a remarkable 22.52% increase in final accuracy and a substantial 6.13% improvement in best accuracy compared to the implicit approach. These findings indicate that simplifying the task of the meta net through the explicit utilization of known information can significantly enhance the performance of the classifier. Furthermore, the notable discrepancy between the final accuracy and the best accuracy underscores the contribution of explicit utilization of known information in

TABLE 9: Ablation studies conducted on CIFAR-10-N-LT, considering varying noise rates (0.1-0.5) and imbalance ratios (10, 100). The average accuracy is reported to evaluate model performance. BS represents Balanced Softmax.

\mathcal{G}_{θ_l}	\mathcal{G}_{θ_m}	BS	Last	Best
\times	\times	\times	52.73	60.43
\times	\checkmark	\times	71.51	76.84
\checkmark	\times	\times	71.08	77.11
\checkmark	\times	\checkmark	72.64	74.96
\checkmark	\checkmark	\times	79.73	80.55

Sampling	Last	Best
Naïve	78.53	79.58
Hierarchical	79.73	80.55

Methods	Last	Best
Cross Entropy	52.73	60.43
Implicit	57.21	74.42
Explicit	79.73	80.55

\mathcal{G}_{θ_l}	\mathcal{G}_{θ_m}	Last	Best
-	-	52.73	60.43
Sample-wise	Sample-wise	44.63	74.42
Sample-wise	Group-wise	40.76	68.08
Group-wise	Group-wise	79.73	80.55

Loss rank	Uncertainty	Last	Best
\times	\times	52.73	60.43
\times	\checkmark	72.91	77.18
\checkmark	\times	79.73	80.55
\checkmark	\checkmark	80.62	81.23

Architecture	Last	Best
Vector	79.50	80.35
3 layer MLPs	79.73	80.55

TABLE 10: Accuracy (%) on CIFAR10-N-LT with varying imbalance ratio. WRN-28-10 used as the classifier.

Noise Rate	0.2			0.4			Avg.
Imbalance Ratio	10	50	200	10	50	200	
Cross Entropy	78.03	65.53	42.06	63.04	47.56	29.09	54.22
HAR-DRW [8]	88.81	82.74	73.98	84.03	75.36	63.95	78.15
FSR [59]	85.70	77.40	65.50	81.60	69.80	49.50	71.58
INOLML [40]	90.10	80.10	66.60	89.10	78.10	61.60	77.60
Dynamic Loss	92.40	83.56	76.72	90.46	78.54	66.47	81.36

facilitating the convergence of the meta net.

Effect of Group-wise Optimization. To validate the efficacy of the group optimization approach employed in our method, we conducted experiments comparing it with sample-wise optimization. The results, presented in Table 9d, reveal that utilizing a sample-wise approach with both the label corrector and margin generator hindered classifier convergence, leading to significantly lower final accuracy compared to the best achieved accuracy. Similarly, employing the margin generator in a group-wise manner alone failed to yield convergence. However, when extending the group-wise approach to the label corrector, a final accuracy of 79.73% was attained. These findings highlight the significance of optimizing the meta net task, as it enhances learning efficiency and improves the classifier’s robustness.

Test with MC-Dropout Uncertainty. To further evaluate the generalizability of our proposed approach, we conducted an assessment by substituting the loss rank metric with MC-Dropout uncertainty, renowned for its simplicity and effectiveness in uncertainty estimation. The results of this evaluation are presented in Table 9e. Initially, we exclusively fed the MC-Dropout uncertainty to the meta net, which yielded a final accuracy of 72.91% and a best accuracy of 77.18%. However, this performance was evidently inferior compared to when the loss rank was provided to the MetaNet. Subsequently, we incorporated both the loss rank and uncertainty as inputs to the meta net, resulting in the highest performance. Specifically, this configuration achieved a final accuracy of 80.62% and a best accuracy of

81.23%. These findings suggest that uncertainty estimation can serve as a significant supplement to loss rank in the context of learning with noisy labels.

Effect of Meta Net Architecture. In order to validate the architecture design of the meta net, we simplified the label corrector and margin generator by replacing them with a R -length and C -length learnable vector, respectively. As shown in Table 9f, this modification resulted in a noticeable performance drop of 0.23%. Our experimental observation suggests that this is due to the fact that MLPs are able to quickly learn appropriate label weights and per-class margins, whereas the learnable vectors suffer from slow convergence.

Test on More Classifiers. To demonstrate the broad applicability of our method, we conducted additional evaluations using the Wide-ResNet-28-10 (WRN-28-10) architecture [60]. Specifically, we set the imbalance ratios to 10, 50, and 200, and the noise rates to 0.2 and 0.4. The mean accuracy results are presented in Table 10. Our findings indicate that our method surpasses the performance of HAR, FSR [59], and INOLML [40] on both CIFAR-10-N-LT datasets. These results provide compelling evidence supporting the effectiveness of our dynamic loss across different classifier architectures.

Training Time Analysis As training time is a critical concern in meta learning, we evaluated the total training time of our methods following the approach taken in DivideMix. Thanks to the use of FaMUS [15] and CurveNet [7] to accelerate the training speed of meta learning, we were able to train a model in approximately 7.2 hours using an NVIDIA GTX 1080 Ti. This is slightly slower than DivideMix with Nvidia V100 GPU (5.2 hours), but it provides evidence of the efficiency of our method.

9 CONCLUSIONS

This work introduces a novel dynamic loss for robust learning from long-tailed data with noisy labels. The dynamic loss consists of a learnable label corrector and margin generator, which jointly correct noisy labels and adjust classification margins to guide classifier learning. The meta net and classifier are co-optimized through meta-learning using

a hierarchical sampling strategy that provides unbiased yet diverse meta data. Extensive evaluations on both synthetic and real-world data demonstrate the effectiveness of our dynamic loss, which exhibits high adaptability and robustness to various types of data biases.

REFERENCES

- [1] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [3] X. Ding, X. Zhang, N. Ma, J. Han, G. Ding, and J. Sun, "Repvgg: Making vgg-style convnets great again," in *CVPR*, 2021.
- [4] W. Li, L. Wang, W. Li, E. Agustsson, and L. Van Gool, "Webvision database: Visual learning and understanding from web data," 2017.
- [5] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma, "Learning imbalanced datasets with label-distribution-aware margin loss," in *NeurIPS*, 2019.
- [6] L. Jiang, Z. Zhou, T. Leung, L.-J. Li, and L. Fei-Fei, "Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels," in *ICML*, 2018.
- [7] S. Jiang, J. Li, Y. Wang, B. Huang, Z. Zhang, and T. Xu, "Delving into sample loss curve to embrace noisy and imbalanced data," in *AAAI*, 2022.
- [8] K. Cao, Y. Chen, J. Lu, N. Arechiga, A. Gaidon, and T. Ma, "Heteroskedastic and imbalanced deep learning with adaptive regularization," in *ICLR*, 2020.
- [9] J. Shu, Q. Xie, L. Yi, Q. Zhao, S. Zhou, Z. Xu, and D. Meng, "Meta-weight-net: Learning an explicit mapping for sample weighting," in *NeurIPS*, 2019.
- [10] B. Kang, S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and Y. Kalantidis, "Decoupling representation and classifier for long-tailed recognition," in *ICLR*, 2019.
- [11] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, and S. X. Yu, "Large-scale long-tailed recognition in an open world," in *CVPR*, 2019.
- [12] G. Zheng, A. H. Awadallah, and S. Dumais, "Meta label correction for noisy label learning," in *AAAI*, 2021.
- [13] J. Li, R. Socher, and S. C. Hoi, "Dividemix: Learning with noisy labels as semi-supervised learning," in *ICLR*, 2020.
- [14] J. Ren, C. Yu, S. Sheng, X. Ma, H. Zhao, S. Yi, and H. Li, "Balanced meta-softmax for long-tailed visual recognition," in *NeurIPS*, 2020.
- [15] Y. Xu, L. Zhu, L. Jiang, and Y. Yang, "Faster meta update strategy for noise-robust deep learning," in *CVPR*, 2021.
- [16] Y. Huang, B. Bai, S. Zhao, K. Bai, and F. Wang, "Uncertainty-aware learning against label noise on imbalanced datasets," in *AAAI*, 2022.
- [17] A. Antoniou, H. Edwards, and A. Storkey, "How to train your maml," in *ICLR*, 2019.
- [18] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [19] H. Song, M. Kim, and J.-G. Lee, "Selfie: Refurbishing unclean samples for robust deep learning," in *ICML*, 2019.
- [20] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "Smote: synthetic minority over-sampling technique," *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
- [21] M. Buda, A. Maki, and M. A. Mazurowski, "A systematic study of the class imbalance problem in convolutional neural networks," *Neural Networks*, vol. 106, pp. 249–259, 2018.
- [22] K. Tang, J. Huang, and H. Zhang, "Long-tailed classification by keeping the good and removing the bad momentum causal effect," in *NeurIPS*, 2020.
- [23] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *ICCV*, 2017.
- [24] C. Huang, Y. Li, C. C. Loy, and X. Tang, "Learning deep representation for imbalanced classification," in *CVPR*, 2016.
- [25] A. K. Menon, S. Jayasumana, A. S. Rawat, H. Jain, A. Veit, and S. Kumar, "Long-tail learning via logit adjustment," 2020.
- [26] C. Huang, Y. Li, C. C. Loy, and X. Tang, "Deep imbalanced learning for face recognition and attribute prediction," *TPAMI*, vol. 42, no. 11, pp. 2781–2794, 2019.
- [27] J. Tan, C. Wang, B. Li, Q. Li, W. Ouyang, C. Yin, and J. Yan, "Equalization loss for long-tailed object recognition," in *CVPR*, 2020.
- [28] M. P. Kumar, B. Packer, and D. Koller, "Self-paced learning for latent variable models," in *NeurIPS*, 2010.
- [29] J. Huang, L. Qu, R. Jia, and B. Zhao, "O2u-net: A simple noisy label detection approach for deep neural networks," in *CVPR*, 2019.
- [30] S. Jiang, J. Li, Y. Wang, W. Wu, J. Zhang, B. Huang, and T. Xu, "Metaseg: Content-aware meta-net for omni-supervised semantic segmentation," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [31] S. Reed, H. Lee, D. Anguelov, C. Szegedy, D. Erhan, and A. Rabinovich, "Training deep neural networks on noisy labels with bootstrapping," 2014.
- [32] W. Chen, C. Zhu, and Y. Chen, "Sample prior guided robust model learning to suppress noisy labels," 2021.
- [33] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. Raffel, "Mixmatch: A holistic approach to semi-supervised learning," in *NeurIPS*, 2019.
- [34] T. Wei, J.-X. Shi, W.-W. Tu, and Y.-F. Li, "Robust long-tailed learning under label noise," 2021.
- [35] S. Liu, J. Niles-Weed, N. Razavian, and C. Fernandez-Granda, "Early-learning regularization prevents memorization of noisy labels," in *NeurIPS*, 2020.
- [36] D. Ortego, E. Arazo, P. Albert, N. E. O'Connor, and K. McGuinness, "Multi-objective interpolation training for robustness to label noise," in *CVPR*, 2021.
- [37] I. Loshchilov and F. Hutter, "Sgdr: Stochastic gradient descent with warm restarts," in *ICLR*, 2017.
- [38] J. Shu, X. Yuan, D. Meng, and Z. Xu, "Cmw-net: Learning a class-aware sample weighting mapping for robust deep learning," 2022.
- [39] A. Iscen, J. Valmadre, A. Arnab, and C. Schmid, "Learning with neighbor consistency for noisy labels," in *CVPR*, 2022.
- [40] D. A. Hoang, G. Carneiro *et al.*, "Maximising the utility of validation sets for imbalanced noisy-label meta-learning," *arXiv preprint arXiv:2208.08132*, 2022.
- [41] M. Chen, H. Cheng, Y. Du, M. Xu, W. Jiang, and C. Wang, "Two wrongs don't make a right: Combating confirmation bias in learning with label noise," 2021.
- [42] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *AAAI*, 2017.
- [43] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009.
- [44] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.
- [45] Y. Zhang, S. Zheng, P. Wu, M. Goswami, and C. Chen, "Learning with feature-dependent label noise: A progressive approach," in *ICLR*, 2021.
- [46] Y. Chen, S. X. Hu, X. Shen, C. Ai, and J. A. K. Suykens, "Compressing features for learning with noisy labels," *TNNLS*, 2022.
- [47] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. Tsang, and M. Sugiyama, "Co-teaching: Robust training of deep neural networks with extremely noisy labels," in *NeurIPS*, 2018.
- [48] E. Englesson and H. Azizpour, "Generalized jensen-shannon divergence loss for learning with noisy labels," in *NeurIPS*, 2021.
- [49] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *ECCV*, 2016.
- [50] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-balanced loss based on effective number of samples," in *CVPR*, 2019.
- [51] S. Alshammari, Y.-X. Wang, D. Ramanan, and S. Kong, "Long-tailed recognition via weight balancing," in *CVPR*, 2022.
- [52] Z. Zhong, J. Cui, S. Liu, and J. Jia, "Improving calibration for long-tailed recognition," in *CVPR*, 2021.
- [53] A. K. Menon, S. Jayasumana, A. S. Rawat, H. Jain, A. Veit, and S. Kumar, "Long-tail learning via logit adjustment," in *ICLR*, 2021.
- [54] S. Park, Y. Hong, B. Heo, S. Yun, and J. Y. Choi, "The majority can help the minority: Context-rich minority oversampling for long-tailed classification," in *CVPR*, 2022, pp. 6887–6896.
- [55] J. Cui, Z. Zhong, S. Liu, B. Yu, and J. Jia, "Parametric contrastive learning," in *ICCV*, 2021.
- [56] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *CVPR*, 2017.
- [57] S. Sinha, H. Ohashi, and K. Nakamura, "Class-difficulty based methods for long-tailed visual recognition," *International Journal of Computer Vision*, vol. 130, no. 10, pp. 2517–2531, 2022.

- [58] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [59] Z. Zhang and T. Pfister, "Learning fast sample re-weighting without reward data," in *ICCV*, 2021.
- [60] S. Zagoruyko and N. Komodakis, "Wide residual networks," in *BMVC*, 2016.