

Learning with Imbalanced Noisy Data by Preventing Bias in Sample Selection

Huafeng Liu*, Mengmeng Sheng*, Zeren Sun, Yazhou Yao, Xian-Sheng Hua, and Heng-Tao Shen

Abstract—Learning with noisy labels has gained increasing attention because the inevitable imperfect labels in real-world scenarios can substantially hurt the deep model performance. Recent studies tend to regard low-loss samples as clean ones and discard high-loss ones to alleviate the negative impact of noisy labels. However, real-world datasets contain not only noisy labels but also class imbalance. The imbalance issue is prone to causing failure in the loss-based sample selection since the under-learning of tail classes also leans to produce high losses. To this end, we propose a simple yet effective method to address noisy labels in imbalanced datasets. Specifically, we propose Class-Balance-based sample Selection (CBS) to prevent the tail class samples from being neglected during training. We propose Confidence-based Sample Augmentation (CSA) for the chosen clean samples to enhance their reliability in the training process. To exploit selected noisy samples, we resort to prediction history to rectify labels of noisy samples. Moreover, we introduce the Average Confidence Margin (ACM) metric to measure the quality of corrected labels by leveraging the model’s evolving training dynamics, thereby ensuring that low-quality corrected noisy samples are appropriately masked out. Lastly, consistency regularization is imposed on filtered label-corrected noisy samples to boost model performance. Comprehensive experimental results on synthetic and real-world datasets demonstrate the effectiveness and superiority of our proposed method, especially in imbalanced scenarios. The source code has been made available at <https://github.com/NUST-Machine-Intelligence-Laboratory/CBS>.

Index Terms—Imbalanced label noise, class-balance-based sample selection, confidence-based sample augmentation, consistency regularization, average confidence margin.

I. INTRODUCTION

DEEP neural networks (DNNs) have obtained remarkable achievements in various tasks (*e.g.*, image classification [1], [2], object detection [3], [4], face recognition [5], [6], instance segmentation [7]–[10], natural language processing [11]) in recent years. These successes are highly attributed to large-scale accurately-labeled training datasets (*e.g.*, ImageNet [12]). Nevertheless, acquiring high-quality manual annotations is expensive and time-consuming, especially for tasks requiring expert knowledge for annotating (*e.g.*, medical images [13]). To obtain large-scale annotated data under a limited budget, recent researchers have started to pay attention to using crowd-sourcing platforms [14] or web image search engines [15] for dataset construction. Despite reducing the cost of

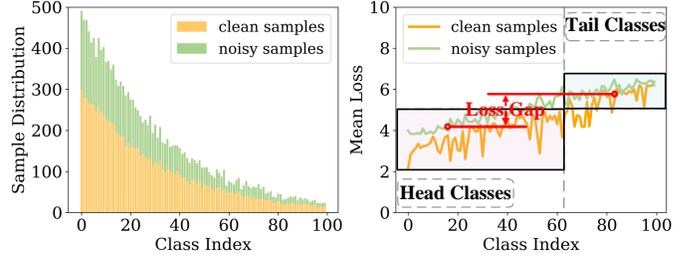


Fig. 1. The sample distribution (left) and the mean loss variation (right) on noisy and imbalanced CIFAR100 (noise rate is 0.4 and imbalance factor is 20). We can find: (1) both tail class samples and noisy samples exhibit large losses; (2) losses of some clean samples belonging to tail classes are even larger than losses of some noisy ones from head classes. Accordingly, existing low-loss-based sample selection methods tend to fail when distinguishing clean and noisy samples. This inspires us to develop a class-balanced sample selection method to combat noisy and imbalanced labels.

data collection, these methods inevitably introduce low-quality samples that are associated with noisy labels. Noisy labels tend to result in inferior model performance due to the strong learning ability of DNNs [16]. Therefore, it is significant to develop robust methods for alleviating noisy labels.

Recently, some methods have been proposed to address the label noise problem [17]–[31]. Existing approaches mainly employ two kinds of strategies for tackling noisy labels: loss/label correction [21], [32], [33] and sample selection [17], [19], [34]. Loss/label correction methods typically attempt to rectify labels by using the noise transition matrix [35]–[37] or model predictions [18], [21], [33]. For example, methods such as loss correction [35] attempt to first estimate the noise transition matrix and then utilize forward and backward correction to mitigate the impact of label noise. [37] proposes a transition-revision (T-Revision) method to effectively learn transition matrices, leading to better classifiers. JoSRC [18] uses the temporally averaged model (*i.e.*, mean-teacher model) to generate reliable pseudo-label distributions for training. PENCIL [33] proposes to directly learn label distributions for corrupted samples in an end-to-end manner. However, loss/label correction methods usually suffer from error accumulation due to the imperfectness and unreliability of the estimated noise transition matrix and model predictions. Contrarily, sample selection methods primarily seek to divide training samples into a “noisy” subset and a “clean” subset, and then use the “clean” one for training [17], [18], [38], [39]. The effectiveness of recent sample-selection-based approaches is mainly attributed to the *Memorization Effect*: DNNs first fit clean samples and then gradually memorize noisy ones. Accordingly, existing methods usually regard samples with

Huafeng Liu, Mengmeng Sheng, Zeren Sun, and Yazhou Yao are with the School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing 210094, China.

Xian-Sheng Hua is with the Terminus Group, Beijing 100027, China.

Heng-Tao Shen is with the School of Electronic and Information Engineering, Tongji University, Shanghai 201804, China.

*Equal Contribution.

small losses as clean ones. For example, Co-teaching [17] cross-updates two networks using small-loss samples selected by its peer networks. Jo-SRC [18] proposes to employ Jensen-Shannon Divergence for selecting clean samples globally. DivideMix [21] extracts the clean subset by fitting the loss distribution with the Gaussian Mixture Model.

However, real-world scenarios contain not only noisy labels but also class imbalance [12], [40], [41]. Most training data tends to belong to the majority classes (*i.e.*, head classes), while some other classes (*i.e.*, tail classes) may possess only a few training samples. Class imbalance leans to mislead the optimization of DNNs to sub-optimal solutions, in which models will predict most samples as head classes. Samples from tail classes will be under-learned. Consequently, the generalization performance and robustness of DNNs are inevitably degraded. Existing approaches designed for noisy labels usually implicitly hypothesize that training samples are class-balanced and thus tend to fail when noisy samples and class imbalance exist simultaneously. Label correction methods cannot guarantee the reliability of corrected labels since tail classes have far fewer samples than head classes. Sample selection methods are prone to suffering from learning bias. These methods mostly rely on the low-loss criterion. Nevertheless, as shown in Fig. 1, samples from tail classes will also have high losses due to under-learning. Nevertheless, as shown in Fig. 1, (1) both tail class samples and noisy samples exhibit large losses; (2) losses of some clean samples belonging to tail classes are even larger than losses of some noisy ones from head classes due to under-learning.

To alleviate the aforementioned issues, we propose a simple yet effective method to learn with noisy labels by balanced sample selection. Our method ensures that tail class samples are learned sufficiently by preventing head classes from prevailing in the selected clean samples. Specifically, we propose **Class-Balance-based sample Selection (CBS)** to divide training samples into a “clean” subset and a “noisy” subset in a class-balanced manner. Subsequently, we propose **Confidence-based Sample Augmentation (CSA)** to minimize the negative effect caused by noisy tail class samples being grouped into the “clean” subset. By fusing selected clean samples based on confidence, CSA promotes the stability of model training by assuring the reliability of samples fed into the model. Moreover, to exploit selected noisy samples and avoid the waste of data, we resort to prediction history to rectify labels of noisy samples and feed them into the model afterward. In order to alleviate the potential harm induced by low-confidence corrected samples (*i.e.*, presumably erroneous corrections), we introduce the **Average Confidence Margin (ACM)** metric to assess the quality of corrected labels. ACM estimates the contribution of a corrected sample to the model by investigating the gap between its confidence scores of the top-2 candidate corrected labels. Additionally, ACM leverages the evolving training dynamic, ensuring that low-quality corrected labels are effectively masked out. Lastly, we design a consistency regularization term to encourage sample-view-wise and epoch-wise prediction consistency, maximizing data exploitation and boosting model performance further. Comprehensive experimental results have been provided to verify the effectiveness

and superiority of our proposed method.

Our main contributions are summarized as follows :

- We propose a simple yet effective approach to address noisy and imbalanced labels. Our proposed class-balanced sample selection assures class balance during the sample selection process to alleviate the learning bias induced by the data imbalance.
- We propose to employ confidence-based sample augmentation to enhance the reliability of selected clean samples. The exponential moving average (EMA) is leveraged to correct labels for noisy samples by resorting to prediction history. Moreover, consistency regularization is adopted to achieve further model enhancement.
- We propose the average confidence margin metric to measure the quality of corrected labels during training. It quantifies the gap between the confidence scores corresponding to the top-2 candidate corrected labels, thereby ensuring that low-quality corrected noisy samples are appropriately discarded from training.
- We provide comprehensive experimental results on synthetic and real-world datasets to illustrate the superiority of our approach. Extensive ablation studies are conducted to verify the effectiveness of each proposed component.

II. RELATED WORK

A. Learning with Noisy Labels

Label noise in training data has been evidenced to have a detrimental impact on the training of deep neural networks [18], [38], [42]–[45]. Existing methods designed for noisy labels can be primarily categorized into the following three groups: label correction [33], [35], [36], sample selection [17]–[19], and other methods [46]–[52].

1) *Label or Loss Correction*: To cope with label noise, one intuitive idea is to correct sample losses or corrupted labels. Methods such as loss correction [35] attempt to first estimate the noise transition matrix and then utilize forward and backward correction to mitigate the impact of label noise. [37] proposes a transition-revision (T-Revision) method to effectively learn transition matrices, leading to better classifiers. Goldberger *et al.* [36] proposes to use an additional layer to estimate the noise transition matrix. Some other researchers focus on correcting labels based on model predictions. For instance, PENCIL [33] proposes to learn label distributions according to model predictions. Tanaka *et al.* [53] proposes to relabel samples by directly using pseudo-labels in an iterative manner. However, the noise transition matrix is difficult to estimate accurately, while prediction-based label correction tends to suffer from error accumulation. Consequently, these methods are prone to struggling with significant performance drops under high noise settings due to the low quality of corrected labels.

2) *Sample Selection*: Another straightforward idea for addressing noisy labels is to select clean samples and discard selected noisy ones from training. For example, Co-teaching [17] maintains two networks and lets each network select small-loss samples as clean ones for its peer network. Co-teaching+ [20] integrates Co-teaching and model disagreement

to identify clean samples. JoCoR [19] exploits a joint loss to select small-loss samples to encourage agreement between models. Besides the popular low-loss-based sample selection, some recent methods propose new selection criteria for finding clean samples. For instance, PNP [22] simultaneously trains two networks, in which one predicts the category label and the other predicts the noise type. NCE [54] resorts to neighbor data to identify clean and noisy samples. BARE [55] proposes a data-dependent, adaptive sample selection strategy that relies only on batch statistics of a given mini-batch to promote the model robustness against label noise. Nevertheless, these methods usually rely on the class-balanced hypothesis, rendering them inadequate for addressing noisy and imbalanced datasets in real-world scenarios. In this paper, we introduce the class-balance-based sample selection strategy to simultaneously tackle label noise and class imbalance issues. Our method is applied per class, mitigating the loss gap between different classes.

3) *Other Methods*: Apart from the two types of methods mentioned above, there are other attempts that have been established to address noisy labels [46]–[51]. For example, AGCE [47] proposes asymmetric loss functions to address discrete and continuous noisy labels. ELR [49] aims to mitigate the impact of noisy data by applying loss gradient regularization. SR [56] introduces a sparse regularization approach that constrains the network output to a permutation set within a one-hot vector framework. Recently, some researchers have strived to take advantage of contrastive learning methods. TCL [57] proposes to focus on learning discriminative representations aligned with estimated labels through mixup and contrastive learning. Sel-CL [51] introduces selective-supervised contrastive learning to learn robust representations and handle noisy labels.

B. Class Imbalance

Real-world scenarios contain not only noisy labels but also class imbalance, posing a more challenging problem. Prior works mainly resort to the sample re-weighting strategy for addressing class imbalance [58]–[61]. These methods usually assign larger weights to tail classes while smaller weights to head classes. For example, [58] proposes to assign different weights to training samples based on gradient directions. [59] proposes a sample weighting function based on meta-learning. However, existing approaches are usually vulnerable when training with noisy and imbalanced data. It should be noted that noisy and tail class samples exhibit high losses. Noisy samples require smaller weights, while tail class samples require larger weights. CNLCU [62], CoDis [63], CurveNet [60] and ULC [61] propose initial attempts to address noisy labels and class imbalance simultaneously. CNLCU [62] extends time intervals and utilizes the mean of training losses at different training iterations to reduce the uncertainty of small-loss examples. CoDis [63] measures the discrepancy by using the distance of prediction probabilities between two networks. CurveNet [60] proposes to learn valuable priors for sample weight assignment based on the loss curves. ULC [61] performs epistemic uncertainty-aware class-specific noise modeling to identify trustworthy clean samples and

refine/discard highly confident true/corrupted labels. In this work, instead of following the re-weighting paradigm, we propose a class-balanced sample selection method to ensure that tail classes are sufficiently learned during training.

III. METHODS

To effectively mitigate the performance degradation caused by noisy labels and class imbalance, we propose to learn from noisy labels by employing balanced sample selection. Initially, we partition the training dataset into two subsets (*i.e.*, the clean and noisy subsets) based on our proposed class-balance-based sample selection (CBS) method. For samples in the clean subset, we further enhance their reliability using the proposed confidence-based sample augmentation (CSA). For samples inside the noisy subset, we correct their given labels based on the exponential moving average (EMA). Besides, we introduce the average confidence margin (ACM) metric to enhance the quality of corrected labels as the training progresses. Lastly, we incorporate a consistency regularization term to further boost the model performance by encouraging sample-view-wise and epoch-wise prediction consistency. The overall framework of our approach is shown in Figure 2.

A. Preliminaries

Let $D_{train} = \{(x_i, y_i) | i = 1, \dots, N\}$ be a noisy C -class dataset containing N training samples, where x_i denotes the i -th image and $y_i \in \{0, 1\}^C$ is its associated label (potentially noisy). y_i^* is the ground-truth label of x_i . We denote $\mathcal{F}(\cdot, \theta)$ as the neural network model parameterized by θ . Given an image-label pair (x, y) , we optimize the network by employing the loss function $\mathcal{L}(\mathcal{F}(x, \theta), y)$ (*e.g.*, cross-entropy loss) during the training process. In the conventional training process, we implicitly assume that the annotated labels of all training samples are accurate (*i.e.*, $y_i = y_i^*$), and use the following cross-entropy loss to optimize the model parameters.

$$\begin{aligned} \mathcal{L}(\mathcal{F}(x, \theta), y) &= \frac{1}{N} \sum_{i=1}^N l_{ce}(x_i, y_i) \\ &= -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_i^c \log(p^c(x_i, \theta)), \end{aligned} \quad (1)$$

in which $p^c(x_i, \theta)$ denotes the predicted softmax probability of the i -th training sample x_i over its c -th class.

Due to the existence of label noise, the empirical risk minimization based on the above loss \mathcal{L} leads to an ill-suited solution. Recent researchers [21], [54], [64] have attempted to employ the semi-supervised learning (SSL) framework by combining the sample selection and label correction methods. First, the sample selection method is adopted to divide the training set D_{train} into a clean subset D_c and a noisy subset D_n . Then, the SSL-based method performs label correction on the subset D_n and subsequently uses the corrected labels for model training. This work also follows the SSL-based paradigm for addressing noisy and imbalanced labels.

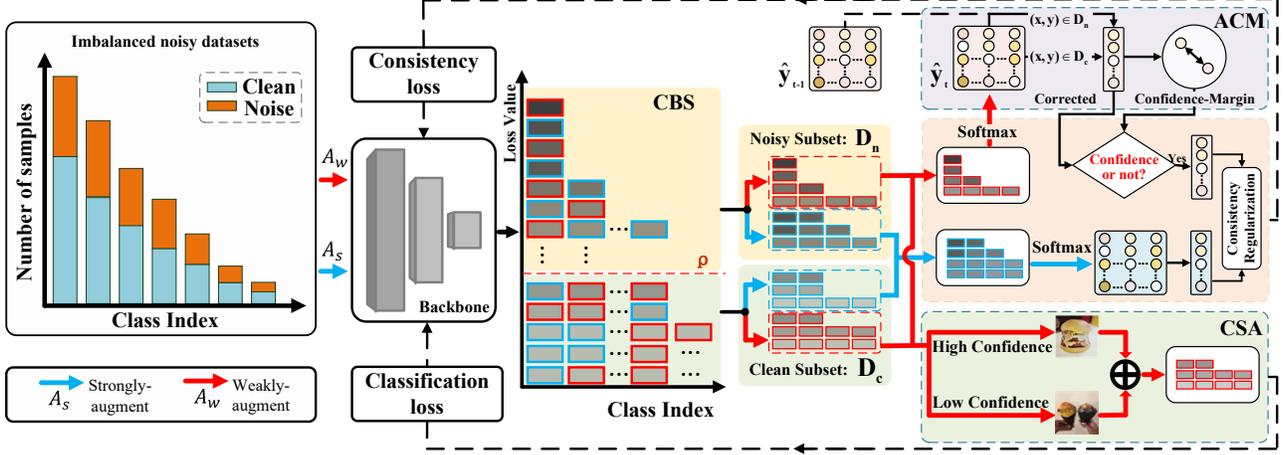


Fig. 2. The overall framework of our proposed approach. We first divide the noisy training set into clean and noisy subsets in a class-balanced manner based on the proposed class-balance-based sample selection (CBS) method. Then, for samples in the clean subset, we propose a confidence-based sample augmentation (CSA) method to enhance the reliability of the selected clean samples. Subsequently, the exponential moving average (EMA) is adopted for correcting labels of noisy samples. Thus, noisy samples are also used for model training. Besides, the average confidence margin (ACM) is proposed to measure the quality of corrected labels as the training progresses. Finally, we employ consistency regularization to boost the model performance further. This regularization term can not only enhance the extracted features but also stabilize training by encouraging epoch-wise prediction consistency.

B. Class-balance-based Sample Selection

Due to the memorization effect of deep neural networks, previous studies [17]–[19], [21] usually select low-loss samples as clean ones in each mini-batch B . The conventional low-loss sample selection strategy empirically performs well in class-balanced datasets. However, these methods tend to suffer from learning bias since the losses of simple class samples tend to be lower than those of hard class samples. This learning bias issue is amplified when the dataset is not only label-noisy but also class-imbalanced. As mentioned above, losses of tail classes tend to be higher than those of head classes. Clean samples belonging to tail classes may even have higher losses than noisy samples from head classes. The class imbalance substantially prevents current sample selection methods from distinguishing between clean samples and noisy ones.

To this end, we propose a *class-balance-based sample selection* (CBS) strategy to address the concurrent label noise and class imbalance issues. Specifically, we first normalize the losses of all samples. The normalization is done on all samples to map all losses to $[0, 1]$, bringing them into a common scale and making loss comparison more meaningful.

$$l(\mathcal{F}(x, \theta), y) = \frac{l_{ce} - \min\{l_{ce}\}}{\max\{l_{ce}\} - \min\{l_{ce}\}}, \quad (2)$$

in which

$$l_{ce} = l_{ce}(\mathcal{F}(x, \theta), y), (x, y) \in D_{train}. \quad (3)$$

It should be noted that the loss-based sample selection is done per class after normalization, thus mitigating the loss gap between different classes. Let D_{sub_i} denote the sample set of the i -th class. We select the top $\lfloor \rho \frac{|D_{train}|}{C} \rfloor$ small-loss samples from D_{sub_i} as “clean” samples belonging to the i -th class. ρ indicates the sample selection ratio, which is designed to control the number of the selected “clean” samples. In practice, we set the value of ρ based on the estimated noise rate η (i.e., $\rho = 1 - \eta$). If $|D_{sub_i}| < \lfloor \rho \frac{|D_{train}|}{C} \rfloor$, all samples

in D_{sub_i} will be selected as “clean” ones. Thus, we can get D_c and D_n as follows:

$$D_c = \bigcup_{i \in \{1, \dots, C\}} D_{c_i}, \quad (4)$$

$$D_n = D_{train} - D_c, \quad (5)$$

in which

$$D_{c_i} = \arg \min_{D'_{c_i} \subseteq D_{sub_i}: |D'_{c_i}| = \delta, (x_j, y_j) \in D_{sub_i}} l_{ce}(\mathcal{F}(x_j, \theta), y_j), \quad (6)$$

$$\delta = \min(\lfloor \rho \frac{|D_{train}|}{C} \rfloor, |D_{sub_i}|). \quad (7)$$

Our proposed class-balanced-based sample selection prevents samples of tail classes from being neglected in the selection procedure, ensuring their adequate participation in the training process. Consequently, the network can sufficiently learn from the tail class samples and correctly produce label predictions.

Discussion. Deep networks usually learn categories that have more samples better than those having fewer ones. Given that the cross-entropy loss is unbounded, losses of samples from different categories tend to have different scales, resulting in the class-wise loss gap. This loss gap is prone to hampering sample selection and thereby downgrading model performance. To address this issue, our method performs sample selection per class using normalized losses. The loss normalization brings selection metrics to a common scale, making the comparison easier and more meaningful. Meanwhile, the class-wise sample selection strategy effectively mitigates the negative impact caused by the loss gap between different categories.

It is also worth noting that although our proposed class-balanced-based sample selection is designed for imbalanced noisy datasets, it is also beneficial for balanced noisy ones. In balanced noisy datasets, the learning difficulties of various categories are inconsistent. Samples from simple categories

tend to yield smaller losses since they are better learned by the network. Contrarily, samples from hard classes usually result in larger losses. This issue leans to make the trained network have biased and inferior recognition performance. By using our proposed class-balanced-based sample selection method, we can alleviate the imbalanced selection results caused by the biased learning ability of the model, thus achieving better model performance.

C. Confidence-based Sample Augmentation

Resorting to the proposed class-balance-based sample selection strategy, we can effectively distinguish between clean and noisy samples while ensuring that tail class samples are selected sufficiently for the subsequent training. However, this selection process will inevitably result in some noisy samples from tail classes being misselected into the clean subset. This issue may lead to a decrease in the model performance.

Therefore, we propose a *confidence-based sample augmentation* (CSA) method for enhancing the reliability of selected clean samples. To be specific, for each sample $(x_i, y_i) \in D_c$ selected by CBS, we randomly choose another sample $(x_j, y_j) \in D_c$ and integrate them to obtain $(\tilde{x}_i, \tilde{y}_i)$ for sample enhancement. Samples with higher prediction confidence are more likely to be truly clean. When integrating the selected two samples, we assign a larger coefficient for the sample whose prediction confidence is higher and a lower coefficient for the sample with lower prediction confidence. Here, we use the max predicted softmax probability to measure the prediction confidence. Thus, the generated $(\tilde{x}_i, \tilde{y}_i)$ is as follows:

$$\tilde{x}_i = \begin{cases} lx_i + (1-l)x_j, p(x_i)^{max} \geq p(x_j)^{max}, \\ (1-l)x_i + lx_j, p(x_i)^{max} < p(x_j)^{max}, \end{cases} \quad (8)$$

$$\tilde{y}_i = \begin{cases} ly_i + (1-l)y_j, p(x_i)^{max} \geq p(x_j)^{max}, \\ (1-l)y_i + ly_j, p(x_i)^{max} < p(x_j)^{max}. \end{cases} \quad (9)$$

$l = \max(l', 1-l')$, in which l' is sampled from a Beta distribution $B(\Phi, \Phi)$ (In our implementation, Φ is empirically set to 4). $p(x_i)^{max}$ and $p(x_j)^{max}$ denote the max predicted softmax probabilities of x_i and x_j , respectively.

By adopting the proposed confidence-based sample augmentation, we reconstruct the selected clean subset as

$$\widetilde{D}_c = \{(\tilde{x}, \tilde{y}) | (x, y) \in D_c\}. \quad (10)$$

We accordingly enhance the reliability of selected clean samples. Then, based on Eq. (1), we calculate the loss on the obtained clean subset \widetilde{D}_c as follows:

$$\mathcal{L}_{D_c} = -\frac{1}{|\widetilde{D}_c|} \sum_{(\tilde{x}, \tilde{y}) \in \widetilde{D}_c} \tilde{y} \log p(\tilde{x}, \theta). \quad (11)$$

Discussion. It is worth noting that our CSA, inspired by Mixup [48], is designed to minimize the negative effect of selecting noisy samples as “clean”. However, unlike Mixup, CSA integrates selected “clean” samples from the clean subset and assigns larger coefficients to samples with higher prediction confidence. High-confidence samples are more likely to

be truly clean. Accordingly, CSA maximizes data reliability by ensuring augmented data contains at least some clean knowledge, thus promoting generalization performance. Although this is a rare occurrence, it is still possible that these two samples are both noisy. When tackling this kind of extreme case where two samples are both noisy, the combination of y (Eq. (9)) smooths label distributions and thus slows down the fitting on label noise, thereby effectively enhancing the model’s generalization performance.

D. Label Correction & Average Confidence Margin

Discarding selected noisy samples directly leads to a waste of data. Meanwhile, our proposed sample selection method may introduce another issue: some clean samples belonging to head classes may be mistakenly identified as noisy samples. Consequently, we follow semi-supervised learning and conduct label correction for selected noisy samples before feeding them to the network. Moreover, we propose to impose the metric of *average confidence margin* (ACM) to measure the quality of corrected labels by using the model’s training dynamics. ACM ensures that low-quality corrected samples are appropriately masked out during training.

Considering that the model is inevitable to fit noisy samples in the later stage of training, we resort to the *Exponential Moving Average* (EMA) to achieve more reliable label correction. The corrected labels for noisy samples are formulated as:

$$\hat{y}^t = \alpha \hat{y}^{t-1} + (1-\alpha)p(A_w(x), \theta), (x, y) \in D_n. \quad (12)$$

\hat{y}^t is the soft corrected label in the t -th epoch. α is the EMA coefficient. $A_w(x)$ represents the weakly augmented view of the sample x . By introducing the prediction history to alleviate the misguidance from erroneously predicted outputs, the correction results are encouraged to be more robust.

Nevertheless, corrected labels with low confidence are not beneficial for model training, as they are essentially akin to noisy labels. Accordingly, inspired by [65]–[67], we introduce the metric of *Confidence Margin* (CM) to measure the quality of corrected labels as follows:

$$CM_j^t(x) = \begin{cases} \hat{y}_j^t - \max_{c \neq j}(\hat{y}_c^t), j = \arg \max(\hat{y}^t), \\ \hat{y}_j^t - \max(\hat{y}^t), j \neq \arg \max(\hat{y}^t). \end{cases} \quad (13)$$

\hat{y}_j^t is the confidence corresponding to the j -th class of the soft corrected label \hat{y}^t . $CM_{\arg \max(\hat{y}^t)}^t$ quantifies the confidence margin between classes with the largest and the second-largest confidence scores in the corrected label distribution. Consequently, a lower $CM_{\arg \max(\hat{y}^t)}^t$ value indicates greater ambiguity in the model prediction, making the corresponding label correction less reliable. We also compute CM_j^t for the remaining classes $j \neq \arg \max(\hat{y}^t)$, aiming to reflect how these classes confuse the model prediction.

We find that CM only considers the model predictions at the current epoch, making it potentially unstable. Thus, we further propose *average confidence margin* (ACM) to average all the margins with respect to the corrected label from the beginning of training until the current epoch t as follows:

$$ACM^t(x) = \frac{1}{t} \sum_{k=1}^t CM_{\arg \max \hat{y}^k}^k(x). \quad (14)$$

ACM implements an iterative estimation method for assessing the contribution of corrected labels to model learning and generalization during training, providing a more stable measure of confidence for corrected labels.

In practice, we maintain a vector of confidence margins for all classes accumulated during training. We dynamically retrieve the accumulated confidence margin of the predicted class (*i.e.*, $\arg \max \hat{y}^t$) at epoch t to obtain ACM^t . Eq.13 illustrates that CM_j^t is positive for $j = \arg \max(\hat{y}^t)$ and negative for $j \neq \arg \max(\hat{y}^t)$. Hence, when predictions of the model frequently disagree across different epochs, the confidence margins for $\arg \max \hat{y}^t$ in previous epochs may not consistently be positive, leading to a low ACM^t . In cases where the model predictions show uniformity and stability across epochs in the corrected label, the confidence margins for $\arg \max \hat{y}^t$ in previous epochs are more likely to be positive, resulting in a higher ACM^t . Accordingly, ACM is evidenced to dynamically capture the characteristics of erroneously corrected labels that adversely affect the training process. We take a linear interpolation of all corrected labels' ACM at t -th epoch as a threshold \mathcal{T} to mask out corrected labels with low confidence.

$$\mathcal{T}^t = \min(ACM) + (\max(ACM) - \min(ACM)) * \tau, \quad (15)$$

where τ is set to control the value of \mathcal{T} (In our implementation, τ is empirically set to 0.2).

After integrating our proposed ACM, we further enhance the model performance through a consistency regularization loss \mathcal{L}_{reg} between the weakly and strongly augmented sample views. \mathcal{L}_{reg} ensures that reliable corrected noisy data is effectively utilized as follows:

$$\mathcal{L}_{reg} = -\frac{1}{|D'_n|} \sum_{(x,y) \in |D'_n|} \hat{y} \log p(A_s(x), \theta). \quad (16)$$

$D'_n = \{(x, y) | ACM^t(x) > \mathcal{T}^t, (x, y) \in D_n\}$. $A_s(x)$ denotes the strongly augmented view of the sample x . $ACM^t > \mathcal{T}^t$ is used to mask out corrected labels with low confidence, hindering their induced harm to the model training. It is worth noting that the samples masked out are only a portion of the samples in the noisy subset, whose corrected labels are deemed unreliable. By employing this consistency regularization design, we achieve prediction consistency between different sample views, implicitly enhancing the feature extraction of the network. Furthermore, we also attain epoch-wise prediction consistency for noisy samples, strengthening the stability of the model optimization. The epoch-wise label consistency is implicitly realized by Eq. (12), which integrates historical and current model prediction results. As noted in Eq. (12), \hat{y} contains predictions from previous epochs to alleviate the misguidance from error prediction, thereby enhancing model stability and reliability.

E. Overall Framework

The learning procedure of our proposed method is illustrated in Algorithm 1 and Fig. 2. The final objective loss function in our method is:

$$\mathcal{L} = \mathcal{L}_{D_c} + \alpha \mathcal{L}_{reg}. \quad (17)$$

Algorithm 1: Our proposed algorithm

Input: The training set D_{train} , the test set D_{test} , the neural network $\mathcal{F}(\cdot, \theta)$, warm-up epochs T_w , total epochs T_{total} , the sample selection ratio ρ , and the batch size bs .

```

1: for epoch = 1, 2, ..., Ttotal do
2:   if epoch ≤ Tw then
3:     for iteration = 1, 2, ..., do
4:       Fetch B = {(xi, yi)}1bs from Dtrain
5:       Calculate Lce = -∑i=1bs yi log p(xi, θ)
6:       Calculate Lcp = -∑i=1bs p(xi, θ) log p(xi, θ)
7:       Calculate L = Lce + Lcp
8:       Update θ by optimizing L
9:       Obtain CM by Eq. (13)
10:      Obtain ACM by Eq. (14)
11:    end for
12:   end if
13:   if Tw < epoch ≤ Ttotal then
14:     Obtain Dc and Dn based on Eqs. (4) and (5).
15:     for iteration = 1, 2, ..., do
16:       Fetch B = {(xi, yi)}1bs from Dtrain
17:       Obtain B̃ ⊆ B by Eqs. (8) and (9)
18:       Obtain ŷ by Eq. (12)
19:       Obtain CM by Eq. (13)
20:       Obtain ACM by Eq. (14)
21:       Obtain Tt by Eq. (15)
22:       Calculate LDc and Lreg using Eqs. (11) and (16)
23:       Calculate L = LDc + Lreg
24:       Update θ by optimizing L
25:     end for
26:   end if
27: end for

```

\mathcal{L}_{D_c} and \mathcal{L}_{reg} denote the classification loss term and the consistency regularization loss term, respectively. α is the loss weighting factor

As presented in Algorithm 1, similar to existing methods [17]–[19], our method starts from a warm-up stage. Besides the cross-entropy loss L_{ce} , we additionally leverage an entropy loss L_{cp} in the warm-up. By minimizing L_{ce} and L_{cp} during warm-up, we enhance model prediction confidence. After warm-up, our proposed method first divides the noisy training set into clean and noisy subsets in a class-balanced manner based on the proposed class-balance-based sample selection (CBS) method. Then, for samples in the clean subset, we employ confidence-based sample augmentation (CSA) to increase the reliability of selected clean samples. Subsequently, we correct the labels of noisy samples based on EMA. We introduce the average confidence margin (ACM) to filter noisy samples whose corrected labels are of low quality by leveraging the model's evolving training dynamics. Finally, we impose consistency regularization from two perspectives: (1) we encourage sample-view-wise prediction consistency to improve the feature extraction ability; (2) we enforce epoch-wise prediction consistency on noisy samples to stabilize the

TABLE I

THE AVERAGE TEST ACCURACY (%) ON SYNTHETIC CIFAR10 WITH VARIOUS NOISE RATES AND IMBALANCE FACTORS OVER THE LAST TEN EPOCHS. THE BEST AND SECOND-BEST RESULTS ARE BOLDED AND UNDERLINED, RESPECTIVELY.

Imbalance Factor	Publication	1			10			50		
		0%	20%	60%	0%	20%	60%	0%	20%	60%
Standard	-	91.35	81.18	45.34	83.58	67.33	30.90	69.50	52.20	24.00
Decoupling [68]	NeurIPS 2017	91.00	85.54	69.12	81.73	75.11	35.25	70.49	60.87	31.03
Co-teaching [17]	NeurIPS 2018	91.68	88.82	75.43	82.94	76.91	32.46	68.91	55.47	21.34
Co-teaching+ [20]	ICML 2019	91.20	89.04	74.07	81.65	72.94	24.33	66.91	48.79	18.87
JoCoR [19]	CVPR 2020	91.95	89.09	77.19	83.14	77.17	30.22	68.24	59.38	19.48
DivideMix [21]	ICLR 2020	92.96	<u>91.63</u>	79.27	87.43	79.49	50.61	68.51	62.79	30.88
CDR [69]	ICLR 2021	94.11	89.02	81.27	85.55	75.11	47.28	73.44	59.69	31.81
Jo-SRC [18]	CVPR 2021	93.88	90.57	82.47	<u>87.79</u>	76.02	40.20	<u>78.10</u>	60.75	28.67
Co-LDL [70]	TMM 2022	92.40	90.49	79.14	<u>82.86</u>	75.83	41.44	73.77	53.71	25.72
AGCE [47]	TPAMI 2023	92.79	90.09	82.68	86.08	78.95	52.57	74.12	57.77	29.53
TCL [57]	CVPR 2023	93.06	89.47	<u>85.66</u>	83.10	82.20	52.79	72.35	<u>66.41</u>	<u>35.61</u>
Robust LR [71]	AAAI 2023	<u>94.88</u>	91.06	84.25	87.74	<u>82.44</u>	<u>56.94</u>	73.82	64.25	31.71
Ours	-	95.45	94.30	91.79	89.13	86.42	72.49	82.31	75.36	54.13

TABLE II

THE AVERAGE TEST ACCURACY (%) ON SYNTHETIC CIFAR100 WITH VARIOUS NOISE RATES AND IMBALANCE FACTORS OVER THE LAST TEN EPOCHS. THE BEST AND SECOND-BEST RESULTS ARE BOLDED AND UNDERLINED, RESPECTIVELY.

Imbalance Factor	Publication	1			10			50		
		0%	20%	60%	0%	20%	60%	0%	20%	60%
Standard	-	68.40	52.91	19.59	50.82	34.74	11.31	37.50	23.64	9.42
Decoupling [68]	NeurIPS 2017	67.88	54.09	22.82	52.55	39.16	13.88	40.44	29.60	10.86
Co-teaching [17]	NeurIPS 2018	67.94	61.33	47.10	50.75	43.41	18.13	38.30	28.44	11.21
Co-teaching+ [20]	ICML 2019	67.57	56.97	35.74	51.28	38.59	14.24	39.83	26.64	9.68
JoCoR [19]	CVPR 2020	68.98	61.63	44.34	50.87	42.37	20.10	37.73	28.68	13.79
DivideMix [21]	ICLR 2020	75.86	<u>69.46</u>	39.38	<u>58.33</u>	48.66	16.12	44.51	30.51	10.26
CDR [69]	ICLR 2021	74.86	<u>63.68</u>	42.66	57.11	41.42	20.10	42.25	28.38	12.94
Jo-SRC [18]	CVPR 2021	75.05	67.95	48.71	57.12	<u>50.91</u>	23.21	<u>46.96</u>	<u>37.86</u>	12.61
Co-LDL [70]	TMM 2022	71.02	65.01	40.07	46.06	<u>37.24</u>	17.72	31.54	<u>25.78</u>	12.17
AGCE [47]	TPAMI 2023	72.49	67.07	47.37	57.37	44.04	22.38	43.09	31.86	11.91
TCL [57]	CVPR 2023	74.12	63.52	<u>50.20</u>	56.53	47.36	24.47	45.84	29.71	<u>18.58</u>
Robust LR [71]	AAAI 2023	<u>76.87</u>	68.91	48.07	54.44	46.06	<u>33.34</u>	39.05	29.52	14.39
Ours	-	78.37	75.27	66.57	63.42	56.43	38.67	48.07	42.52	27.30

training process. The final objective loss integrates the classification loss on clean samples and the consistency regularization loss on noisy samples.

IV. EXPERIMENTS

This section focuses on experimental evaluations. We first introduce our experimental setup, including datasets, implementation details, evaluation metrics, and baselines. Afterward, we present experimental results on synthetic datasets (*i.e.*, CIFAR10 and CIFAR100 [72]) and real-world datasets (*i.e.*, Web-Aircraft, Web-Bird, and Web-Car [73]). These results firmly verify the effectiveness of our method in alleviating noisy labels in class-imbalanced datasets. Moreover, we conduct extensive ablation studies to investigate the effectiveness of each component and hyper-parameters in our method.

A. Experimental Setup

Synthetic Datasets: Synthetic datasets are mainly derived from CIFAR10 and CIFAR100 [72]. These two datasets consist of 60,000 RGB images (50,000 for training and 10,000 for testing). Images are equally distributed among 10 categories and 100 categories. We randomly corrupt the sample labels from their ground-truth categories to other categories using

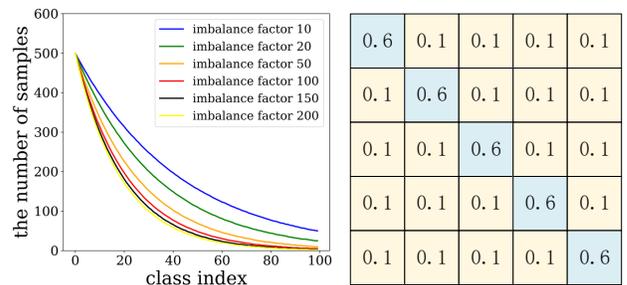


Fig. 3. The number of samples belonging to each class in CIFAR100 under various imbalance factor settings (left) and an example of the uniform noise transition matrix (right).

a pre-defined noise rate (NR) η . In our experiments, we adopt the uniform noise, which randomly corrupts labels from their ground-truth classes to other ones with the pre-defined noise rate η on CIFAR10 and CIFAR100. To construct class-imbalanced datasets, we take an exponential function $n_i = n_0 \mu^i$ to reduce the number of samples per category, where n_i is the sample number of class i and $\mu \in (0, 1]$. We use the class imbalance factor (IF), defined as $\frac{\max(n_i)}{\min(n_i)}$, to measure how imbalanced a dataset is. Fig. 3 (left) presents the sample distribution of synthetic CIFAR100 under different

TABLE III

COMPARISON WITH SOTA APPROACHES IN TEST ACCURACY (%) ON REAL-WORLD NOISY DATASETS: WEB-AIRCRAFT, WEB-BIRD, AND WEB-CAR. THE BEST AND SECOND-BEST RESULTS ARE BOLDED AND UNDERLINED, RESPECTIVELY.

Methods	Publication	Backbone	Performances(%)			
			Web-Aircraft	Web-Bird	Web-Car	Average
Standard	-	ResNet50	60.80	64.40	60.60	61.93
Decoupling [68]	NeurIPS 2017	ResNet50	75.91	71.61	79.41	75.64
Co-teaching [17]	NeurIPS 2018	ResNet50	79.54	76.68	84.95	80.39
Co-teaching+ [20]	ICML 2019	ResNet50	74.80	70.12	76.77	73.90
PENCIL [33]	CVPR 2019	ResNet50	78.82	75.09	81.68	78.53
Hendrycks <i>et al.</i> [74]	NeurIPS 2019	ResNet50	73.24	70.03	73.81	72.36
mCT-S2R [75]	WACV 2020	ResNet50	79.33	77.67	82.92	79.97
JoCoR [19]	CVPR 2020	ResNet50	80.11	79.19	85.10	81.47
AFM [76]	ECCV 2020	ResNet50	81.04	76.35	83.48	80.29
DivideMix [21]	ICLR 2020	ResNet50	82.48	74.40	84.27	80.38
Self-adaptive [77]	NeurIPS 2020	ResNet50	77.92	78.49	78.19	78.20
Peer-learning [78]	ICCV 2021	ResNet50	78.64	75.37	82.48	78.83
Co-LDL [70]	TMM 2022	ResNet50	81.97	80.11	<u>86.95</u>	83.01
NCE [54]	ECCV 2022	ResNet50	84.94	<u>80.22</u>	86.38	<u>83.85</u>
SOP [79]	ICML 2022	ResNet50	84.06	79.40	85.71	83.06
SPRL [80]	PR 2023	ResNet50	84.40	76.36	86.84	82.53
AGCE [47]	TPAMI 2023	ResNet50	84.22	75.60	85.16	81.66
TCL [57]	CVPR 2023	ResNet50	84.51	79.22	85.13	82.95
Robust LR [71]	CVPR 2023	ResNet50	<u>85.78</u>	78.65	86.13	83.52
Ours	-	ResNet50	87.46	80.50	87.96	85.31

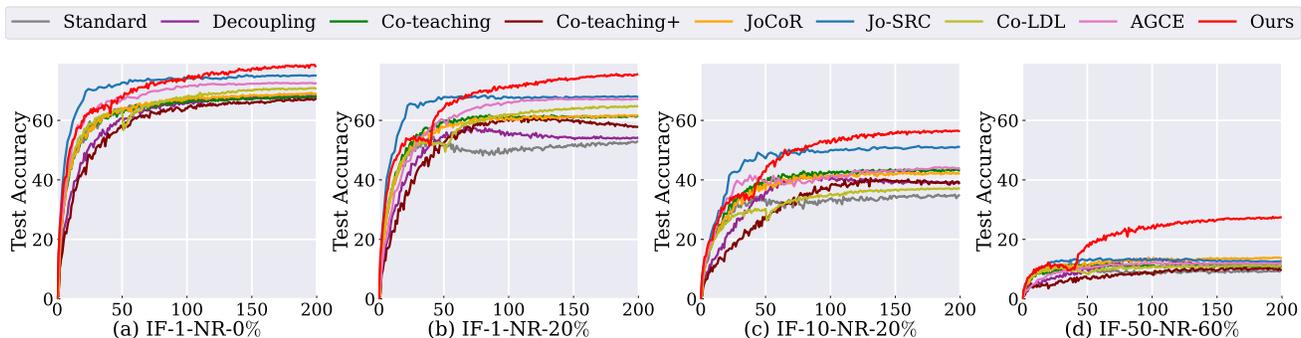


Fig. 4. The test accuracy (%) vs. epochs on CIFAR100 with IF-1-NR-0% (a), IF-1-NR-20% (b), IF-10-NR-20% (c) and IF-50-NR-60% (d) during the training process. (IF-X-NR-Y% means that the imbalance factor and the noise rate are X and Y%, respectively.)

imbalance factors (*i.e.*, 10, 20, 50, *etc.*). Fig. 3 (right) shows an example of the uniform noise transition matrix.

Real-world Datasets: To further verify the effectiveness of our method in practical scenarios, we conduct experiments on real-world noisy datasets: Web-Aircraft, Web-Bird, Web-Car [73] and Clothing1M [41]. These three Web-datasets are subsets of the web-image-based fine-grained image dataset WebFG-496 [73]. Their training images are crawled from web image search engines, making label noise inevitable. Web-Aircraft is a fine-grained aircraft dataset containing 13,503 training images and 3,333 test images belonging to 100 different aircraft models. Web-Bird is a fine-grained bird dataset containing 200 different classes. There are 18,388 noisy training instances, whereas the test set consists of 5794 accurately-labeled samples. Web-Car is a fine-grained car dataset composed of 21,448 samples, and the test set consists of 8,041 samples belonging to 196 car classes. Clothing1M comprises 1M clothing images of 14 categories, which are collected from several online shopping websites and include many mislabelled samples. These real-world datasets contain both corrupted labels and class imbalance.

Implementation Details: On synthetic datasets, we conduct experiments on CIFAR10 and CIFAR100 with various imbalance factors and noise rates. We use ResNet18 as our backbone. The network is trained using SGD with a momentum of 0.9 for 200 epochs. Our warm-up stage lasts for 40 epochs. The initial learning rate and batch size are 0.01 and 128, respectively. During the robust learning stage, we decay the learning rate in a cosine annealing manner. We set ρ and τ as $1 - \eta$ and 0.2, respectively. On real-world datasets (*i.e.*, Web-Aircraft, Web-Bird, Web-Car, and Clothing1M), we follow [70] and select ResNet50 pre-trained on ImageNet as our backbone and SGD with a momentum of 0.9 as the optimizer. We set the initial learning rate as 0.001, and adopt the cosine schedule to adjust the learning rate during training. The training lasts for 100 epochs (including 10 warm-up epochs). We use random cropping and horizontal flipping as weak augmentation, and adopt AutoAugment [81] as strong augmentation. AutoAugment designs a search space where a policy consists of many sub-policies (*i.e.*, translation, rotation, or shearing), one of which is randomly chosen for each image in each mini-batch.

TABLE IV
PERFORMANCE COMPARISON WITH SOTA METHODS IN TEST ACCURACY (%) ON CLOTHING1M.

Method	Publications	Performance
Standard	-	68.94
Decoupling [68]	NeurIPS 2017	69.84
Co-teaching [17]	NeurIPS 2018	69.21
JoCoR [19]	CVPR 2020	70.30
DivideMix [21]	ICLR 2020	74.76
JNPL [82]	CVPR 2021	74.15
UNICON [64]	CVPR 2022	74.98
TCL [57]	CVPR 2023	74.80
Ours	-	74.99

Evaluation Metrics: We adopt test accuracy as the evaluation metric. On synthetic CIFAR10 and CIFAR100 datasets, we additionally employ the average test accuracy over the last epochs to evaluate the performance more comprehensively.

Baselines: We compare our method with state-of-the-art (SOTA) methods for evaluation. On synthetic CIFAR10 and CIFAR100, we compare our method with Decoupling [68], Co-teaching [17], Co-teaching+ [20], JoCoR [19], DivideMix [21], Jo-SRC, CDR [69], [18], Co-LDL [70], AGCE [47], TCL [57], and Robust LR [71]. On Web-Aircraft, Web-Bird, and Web-Car, the following SOTA methods are adopted for comparison: Decoupling [68], Co-teaching [17], PENCIL [33], Hendrycks *et al.* [74], mCT-S2R [75], JoCoR [19], AFM [76], DivideMix [21], Self-adaptive [77], Peer-learning [78], Co-LDL [70], NCE [54], SOP [79], SPRL [80], AGCE [47], TCL [57], and Robust LR [71]). On Clothing1M, We compare the following methods: Decoupling [68], Co-teaching [17], JoCoR [19], DivideMix [21], JNPL [82], UNICON [64] and TCL [57]). Moreover, we perform conventional training using the entire noisy dataset. The result is provided as a baseline (denoted as ‘‘Standard’’).

B. Experimental Results on Synthetic Datasets

For evaluating the performance of our proposed method in learning with imbalanced noisy data, we conduct extensive experiments on the synthetic CIFAR10 and CIFAR100 using different imbalance factors (*i.e.*, 1, 10, 50) and noise rates (*i.e.*, 0%, 20%, 60%). We compare our proposed method with existing SOTA methods, which are re-implemented using their open-sourced code and default hyper-parameters. Results are shown in Table I, Table II and Fig.4.

Results presented in Table I and Table II illustrate: (1) When the dataset contains only noisy labels (*i.e.*, the imbalance factor is 1), both existing SOTA methods and our method achieve robust performance. Notably, our method achieves the best performance. (2) When the dataset is noisy and imbalanced, SOTA methods for learning with noisy labels dreadfully degrade their performance when we increase the imbalance factor and the noise rate. In particular, methods employing the small-loss criterion (*e.g.*, Co-teaching, Co-teaching+, and JoCoR) generally exhibit inferior performance. These methods erroneously discard samples from tail classes due to their large loss values, thereby introducing a learning bias against tail classes.

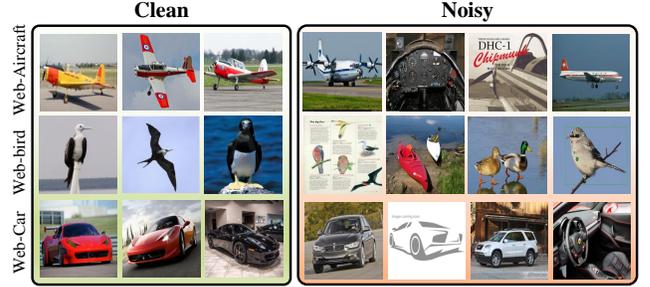


Fig. 5. Some visualization results of clean and noisy samples selected by our sample selection methods on Web-Aircraft, Web-Bird, and Web-Car. The corresponding fine-grained class names are *DHC-1*, *frigatebird*, and *Ferrari 458 Italia Coupe 2012*.

Our method consistently sustains robust performance, surpassing all competing SOTA methods across all experimental settings. Results from Table I and Table II clearly illustrate the effectiveness and superiority of our proposed method. This is mainly attributed to our proposed class-balanced sample selection approach. CBS ensures that the tail class samples sufficiently participate in the model training. Fig. 4 presents the test accuracy vs. epochs in four different scenarios on CIFAR100. It further demonstrates that our proposed method consistently performs better than competing approaches during the training process. Especially in the most challenging situation (*i.e.*, IF-50-NR-60%), our method still performs superiorly throughout the training process. These advances in performance validate the superiority of our method.

C. Experimental Results on Real-World Datasets

In addition to evaluating our method on synthetic noisy and imbalanced datasets (*i.e.*, CIFAR10 and CIFAR100), we also utilize three real-world web-image-based datasets (*i.e.*, Web-Aircraft, Web-Bird, and Web-Car) to corroborate the effectiveness and superiority of our method. Table III presents the performance comparison on real-world datasets. These datasets contain at least 25% of unknown noisy labels and do not provide any label verification information, making them practical and challenging label noise scenarios. It should be emphasized that these three datasets are fine-grained ones, which undoubtedly makes them more challenging. Results of existing methods shown in Table III are obtained under the same experimental settings. As shown in Table III, our method consistently outperforms SOTA methods on the three datasets. Our method achieves 87.46%, 80.50%, and 87.96% accuracy on test sets of Web-Aircraft, Web-Bird, and Web-Car, respectively. It achieves a significant performance boost of +1.68% / +0.28% / +1.01% over the best SOTA method. In the average test accuracy of these three datasets, our proposed approach outperforms Co-LDL [70] and NCE [54] by 2.30% and 1.46%, respectively. Table IV shows the additional experimental results on another real-world noisy and imbalanced dataset Clothing1M. By comparing with SOTA methods in Table IV, we can find that our method can achieve competitive performance against SOTA approaches. It should be noted that DivideMix, UNICON and TCL involve two simultaneously trained networks, while our method trains only one network.

TABLE V
EFFECTS OF DIFFERENT INGREDIENTS IN TEST ACCURACY (%) ON CIFAR10 AND CIFAR100 (20%-10 MEANS THAT NOISE RATE AND IMBALANCE FACTOR ARE 20% AND 10, RESPECTIVELY). RESULTS AT THE BEST EPOCHS ARE PRESENTED.

Model	CIFAR10			CIFAR100		
	20%-1	20%-10	20%-50	0%-10	20%-10	60%-10
Standard	81.18	67.33	52.20	50.82	34.74	11.31
Standard+CBS	90.95	75.31	60.71	54.84	42.51	22.47
Standard+CBS+CSA	91.11	82.58	68.94	56.96	50.60	34.72
Standard+CBS+CSA+CR	94.30	85.46	71.87	62.15	54.46	38.67
Standard+CBS+CSA+CR+ACM	94.82	86.42	75.36	63.42	56.43	39.70

In order to further visualize the performance of our method, we provide qualitative analysis on three real-world datasets. As shown in Fig. 5, we provide several visualization results of clean and noisy samples selected by our sample selection methods on three fine-grained categories (*i.e.*, *DHC-1*, *frigate-bird*, and *Ferrari 458 Italia Coupe 2012*). It is evident that our proposed selection method can effectively distinguish clean and noisy samples.

D. Ablation Studies

In this section, we study the influence of each proposed component (CBS, CSA, CR, and CM) and each hyper-parameter (ρ , τ and α) in our method. We conduct ablation experiments on CIFAR10 and CIFAR100 with various imbalance factors (*i.e.*, 1, 10, 50) and noise rates (*i.e.*, 0%, 20%, 60%). The results are provided in Table V and Fig. 6. Standard represents the conventional forward training using the cross-entropy loss. CBS denotes the class-balance-based sample selection. CSA indicates the confidence-based sample augmentation. CR means consistency regularization. ACM denotes the proposed average confidence margin.

Effects of Class-Balance-based Sample Selection: Based on the aforementioned analysis, loss-based sample selection tends to cause learning bias in imbalanced datasets. In particular, tail class samples are prone to be identified as noise because under-learning leans to result in large loss values. Our class-balance-based sample selection method ensures that the tail classes are not neglected during training. Accordingly, the learning bias issue is effectively alleviated. As shown in Table V, employing CBS achieves notable performance gains compared to Standard in all experimental settings.

Effects of Confidence-based Sample Augmentation: The sample selection process may inevitably cause some noisy samples to be mistakenly selected into the clean subset. Thus, we propose confidence-based sample augmentation to enhance the reliability of selected clean samples. Table V illustrates that adopting CSA boosts model performance, proving the effectiveness of CSA. In particular, as the noise rate and imbalance factor increase, the performance gain is more significant.

Effects of Consistency Regularization: Although we leverage EMA to re-assign labels for selected noisy samples, the corrected labels may still be inaccurate. The imperfect label correction is prone to causing performance degradation. Accordingly, we propose consistency regularization to achieve enhancement in both feature extraction and model prediction. From Table V, we can find CR successfully boosts the model

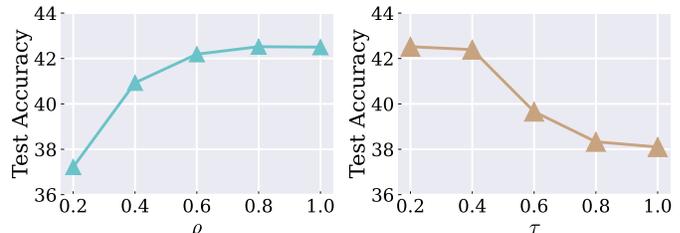


Fig. 6. Hyper-parameter sensitivities of ρ (left) and τ (right). Experiments are conducted on CIFAR100 (imbalance factor is 50 and noise rate is 20%).

performance. For example, employing CR yields a 3.95% performance gain on CIFAR100 when noise rate and imbalance are 60% and 10.

Effects of Average Confidence Margin: In addition to the employment of consistency regularization, we introduce the average confidence margin to measure the confidence of corrected labels, aiming to promote model robustness when learning from label-corrected noisy samples. Our training process uses a dynamic mechanism to continuously assess the impact of corrected labels on learning and generalization, rather than solely relying on model predictions at the current epoch. Accordingly, we discard label-corrected noisy samples with low confidence from training, alleviating their potential damage to the model. As shown in Table V, employing ACM achieves considerable performance gains.

Effects of Hyper-parameters: We investigate the effects of hyper-parameters ρ , τ and α in our proposed method. We take the ρ to control the proportion of selected clean samples per class. τ and α are set to control the number of reliable corrected labels and the loss weight in Eq.17, respectively. We provide the model performance under different ρ and τ settings in Fig. 6. We can observe that a properly selected ρ and τ can boost the model performance further. When ρ is 0.8 (*i.e.*, $1-\eta$), and τ is 0.2, our method achieves the highest performance on the test set on synthetic noisy CIFAR100, whose noise rate is 20% and imbalance factor is 10. Additionally, we further demonstrate the performance of our CBS under different loss weights α in Eq. 17. When the values of α are 0.5, 1.0, 1.5, and 2.0, the test accuracy is 41.37%, 42.52%, 40.41%, and 39.81%, respectively. This further validates the effectiveness of our loss function in Eq. (17).

V. CONCLUSION

In this paper, we focused on the challenge of learning with noisy and imbalanced datasets. To address label noise and class imbalance simultaneously, we proposed a simple yet effective

method based on balanced sample selection. Our proposed method followed the semi-supervised learning paradigm and trained only one network in the training process. Specifically, we proposed a class-balance-based sample selection strategy to divide samples into clean and noisy subsets in a class-balanced manner. We then performed confidence-based sample augmentation to enhance the reliability of selected clean samples. Afterward, we employed EMA to relabel selected noisy samples and filtered those with low confidence based on the average confidence margin metric. Finally, consistency regularization was adopted on label-corrected noisy samples with high confidence to improve the robustness and stability of the model training. Extensive experiments and ablation studies were conducted to substantiate the effectiveness and superiority of our proposed method.

REFERENCES

- [1] X. Zhong, C. Gu, M. Ye, W. Huang, and C. Lin, "Graph complemented latent representation for few-shot image classification," *IEEE Trans. Multimedia*, vol. 25, pp. 1979–1990, 2023. **1**
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Adv. Neural Inform. Process. Syst.*, 2012, pp. 1106–1114. **1**
- [3] Z. Shao, Y. Pu, J. Zhou, B. Wen, and Y. Zhang, "Hyper RPCA: joint maximum coreentropy criterion and laplacian scale mixture modeling on-the-fly for moving object detection," *IEEE Trans. Multimedia*, vol. 25, pp. 112–125, 2023. **1**
- [4] J. Redmon and A. Farhadi, "YOLO9000: better, faster, stronger," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017, pp. 6517–6525. **1**
- [5] F. Boutros, N. Damer, F. Kirchbuchner, and A. Kuijper, "Elasticface: Elastic margin loss for deep face recognition," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022, pp. 1577–1586. **1**
- [6] J. Qian, S. Zhu, C. Zhao, J. Yang, and W. K. Wong, "Otface: Hard samples guided optimal transport loss for deep face representation," *IEEE Trans. Multimedia*, vol. 25, pp. 1427–1438, 2023. **1**
- [7] H. Ying, Z. Huang, S. Liu, T. Shao, and K. Zhou, "Embedmask: Embedding coupling for instance segmentation," in *IJCAI*, 2021, pp. 1266–1273. **1**
- [8] K. Zhang, C. Yuan, Y. Zhu, Y. Jiang, and L. Luo, "Weakly supervised instance segmentation by exploring entire object regions," *IEEE Trans. Multimedia*, vol. 25, pp. 352–363, 2023. **1**
- [9] T. Chen, Y. Yao, and J. Tang, "Multi-granularity denoising and bidirectional alignment for weakly supervised semantic segmentation," *IEEE Trans. Image Process.*, vol. 32, pp. 2960–2971, 2023. **1**
- [10] T. Chen, Y. Yao, L. Zhang, Q. Wang, G. Xie, and F. Shen, "Saliency guided inter- and intra-class relation constraints for weakly supervised semantic segmentation," *IEEE Trans. Multimedia*, vol. 25, pp. 1727–1737, 2023. **1**
- [11] E. L. Malfa, R. Michelmoro, A. M. Zbrzezny, N. Paoletti, and M. Kwiatkowska, "On guaranteed optimal robust explanations for NLP models," in *IJCAI*, 2021, pp. 2658–2665. **1**
- [12] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2009, pp. 248–255. **1, 2**
- [13] T. Wu, B. Dai, S. Chen, Y. Qu, and Y. Xie, "Meta segmentation network for ultra-resolution medical images," in *IJCAI*, 2020, pp. 544–550. **1**
- [14] P. Welinder, S. Branson, S. J. Belongie, and P. Perona, "The multidimensional wisdom of crowds," in *Adv. Neural Inform. Process. Syst.*, 2010, pp. 2424–2432. **1**
- [15] R. Fergus, L. Fei-Fei, P. Perona, and A. Zisserman, "Learning object categories from internet image searches," *Proc. IEEE*, pp. 1453–1466, 2010. **1**
- [16] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," in *Int. Conf. Learn. Represent.*, 2017. **1**
- [17] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. W. Tsang, and M. Sugiyama, "Co-teaching: Robust training of deep neural networks with extremely noisy labels," in *Adv. Neural Inform. Process. Syst.*, 2018, pp. 8536–8546. **1, 2, 4, 6, 7, 8, 9**
- [18] Y. Yao, Z. Sun, C. Zhang, F. Shen, Q. Wu, J. Zhang, and Z. Tang, "Jo-src: A contrastive approach for combating noisy labels," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021, pp. 5192–5201. **1, 2, 4, 6, 7, 9**
- [19] H. Wei, L. Feng, X. Chen, and B. An, "Combating noisy labels by agreement: A joint training method with co-regularization," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020, pp. 13 723–13 732. **1, 2, 3, 4, 6, 7, 8, 9**
- [20] X. Yu, B. Han, J. Yao, G. Niu, I. W. Tsang, and M. Sugiyama, "How does disagreement help generalization against label corruption?" in *Int. Conf. Mach. Learn.*, 2019, pp. 7164–7173. **1, 2, 7, 8, 9**
- [21] J. Li, R. Socher, and S. C. Hoi, "Dividemix: Learning with noisy labels as semi-supervised learning," in *Int. Conf. Learn. Represent.*, 2020. **1, 2, 3, 4, 7, 8, 9**
- [22] Z. Sun, F. Shen, D. Huang, Q. Wang, X. Shu, Y. Yao, and J. Tang, "Pnp: Robust learning from noisy labels by probabilistic noise prediction," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022, pp. 5311–5320. **1, 3**
- [23] Y. Lu and W. He, "SELC: self-ensemble label correction improves learning with noisy labels," in *IJCAI*, 2022, pp. 3278–3284. **1**
- [24] Y. Liu, N. Xu, Y. Zhang, and X. Geng, "Label distribution for learning with noisy labels," in *IJCAI*, 2020, pp. 2568–2574. **1**
- [25] Y. Bai, E. Yang, B. Han, Y. Yang, J. Li, Y. Mao, G. Niu, and T. Liu, "Understanding and improving early stopping for learning with noisy labels," in *Adv. Neural Inform. Process. Syst.*, 2021, pp. 24 392–24 403. **1**
- [26] X. Xia, B. Han, N. Wang, J. Deng, J. Li, Y. Mao, and T. Liu, "Extended \$t\$: Learning with mixed closed-set and open-set noisy labels," *IEEE Trans. Pattern Anal. Mach. Intell.*, pp. 3047–3058, 2023. **1**
- [27] D. Cheng, T. Liu, Y. Ning, N. Wang, B. Han, G. Niu, X. Gao, and M. Sugiyama, "Instance-dependent label-noise learning with manifold-regularized transition matrix estimation," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022, pp. 16 609–16 618. **1**
- [28] S. Li, X. Xia, S. Ge, and T. Liu, "Selective-supervised contrastive learning with noisy labels," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022, pp. 316–325. **1**
- [29] E. Yang, D. Yao, T. Liu, and C. Deng, "Mutual quantization for cross-modal search with noisy labels," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022, pp. 7541–7550. **1**
- [30] C. Gong, Y. Ding, B. Han, G. Niu, J. Yang, J. You, D. Tao, and M. Sugiyama, "Class-wise denoising for robust learning under label noise," *IEEE Trans. Pattern Anal. Mach. Intell.*, pp. 2835–2848, 2023. **1**
- [31] D. Cheng, T. Liu, Y. Ning, N. Wang, B. Han, G. Niu, X. Gao, and M. Sugiyama, "Instance-dependent label-noise learning with manifold-regularized transition matrix estimation," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022, pp. 16 609–16 618. **1**
- [32] E. Arazo, D. Ortego, P. Albert, N. E. O'Connor, and K. McGuinness, "Unsupervised label noise modeling and loss correction," in *Int. Conf. Mach. Learn.*, 2019, pp. 312–321. **1**
- [33] K. Yi and J. Wu, "Probabilistic end-to-end noise correction for learning with noisy labels," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019, pp. 7017–7025. **1, 2, 8, 9**
- [34] Z. Sun, Y. Yao, X. Wei, F. Shen, H. Liu, and X.-S. Hua, "Boosting robust learning via leveraging reusable samples in noisy web data," *IEEE Trans. Multimedia*, 2022. **1**
- [35] G. Patrini, A. Rozza, A. Krishna Menon, R. Nock, and L. Qu, "Making deep neural networks robust to label noise: A loss correction approach," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017, pp. 1944–1952. **1, 2**
- [36] J. Goldberger and E. Ben-Reuven, "Training deep neural-networks using a noise adaptation layer," in *Int. Conf. Learn. Represent.*, 2017. **1, 2**
- [37] X. Xia, T. Liu, N. Wang, B. Han, C. Gong, G. Niu, and M. Sugiyama, "Are anchor points really indispensable in label-noise learning?" in *Adv. Neural Inform. Process. Syst.*, H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox, and R. Garnett, Eds., 2019, pp. 6835–6846. **1, 2**
- [38] Z. Sun, Y. Yao, X.-S. Wei, Y. Zhang, F. Shen, J. Wu, J. Zhang, and H.-T. Shen, "Webly supervised fine-grained recognition: Benchmark datasets and an approach," in *Int. Conf. Comput. Vis.*, 2021, pp. 10 602–10 611. **1, 2**
- [39] X. Gui, W. Wang, and Z. Tian, "Towards understanding deep learning from noisy labels with small-loss criterion," in *IJCAI*, 2021, pp. 2469–2475. **1**
- [40] Y. Cui, M. Jia, T. Lin, Y. Song, and S. J. Belongie, "Class-balanced loss based on effective number of samples," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019, pp. 9268–9277. **2**
- [41] T. Xiao, T. Xia, Y. Yang, C. Huang, and X. Wang, "Learning from massive noisy labeled data for image classification," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2015, pp. 2691–2699. **2, 8**

- [42] Z. Sun, X.-S. Hua, Y. Yao, X.-S. Wei, G. Hu, and J. Zhang, "Crssc: salvage reusable samples from noisy data for robust learning," in *ACM Int. Conf. Multimedia*, 2020, pp. 92–101. [2](#)
- [43] C. Zhang, Y. Yao, X. Xu, J. Shao, J. Song, Z. Li, and Z. Tang, "Extracting useful knowledge from noisy web images via data purification for fine-grained recognition," in *ACM Int. Conf. Multimedia*, 2021, pp. 4063–4072. [2](#)
- [44] H. Liu, C. Zhang, Y. Yao, X.-S. Wei, F. Shen, Z. Tang, and J. Zhang, "Exploiting web images for fine-grained visual recognition by eliminating open-set noise and utilizing hard examples," *IEEE Trans. Multimedia*, vol. 24, pp. 546–557, 2021. [2](#)
- [45] Y. Yao, X. Hua, G. Gao, Z. Sun, Z. Li, and J. Zhang, "Bridging the web data and fine-grained visual recognition via alleviating label noise and domain mismatch," in *ACM Int. Conf. Multimedia*, 2020, pp. 1735–1744. [2](#)
- [46] Z. Zhang and M. R. Sabuncu, "Generalized cross entropy loss for training deep neural networks with noisy labels," in *Adv. Neural Inform. Process. Syst.*, 2018, pp. 8792–8802. [2, 3](#)
- [47] X. Zhou, X. Liu, D. Zhai, J. Jiang, and X. Ji, "Asymmetric loss functions for noise-tolerant learning: Theory and applications," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 7, pp. 8094–8109, 2023. [2, 3, 7, 8, 9](#)
- [48] H. Zhang, M. Cissé, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," in *Int. Conf. Learn. Represent.*, 2018. [2, 3, 5](#)
- [49] S. Liu, J. Niles-Weed, N. Razavian, and C. Fernandez-Granda, "Early-learning regularization prevents memorization of noisy labels," in *Adv. Neural Inform. Process. Syst.*, 2020. [2, 3](#)
- [50] E. Zheltonozhskii, C. Baskin, A. Mendelson, A. M. Bronstein, and O. Litany, "Contrast to divide: Self-supervised pre-training for learning with noisy labels," in *IEEE Winter Conference on Applications of Computer Vision*, 2022, pp. 387–397. [2, 3](#)
- [51] S. Li, X. Xia, S. Ge, and T. Liu, "Selective-supervised contrastive learning with noisy labels," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022, pp. 316–325. [2, 3](#)
- [52] C. Zhang, G. Lin, Q. Wang, F. Shen, Y. Yao, and Z. Tang, "Guided by meta-set: A data-driven method for fine-grained visual recognition," *IEEE Trans. Multimedia*, vol. 25, pp. 4691–4703, 2023. [2](#)
- [53] D. Tanaka, D. Ikami, T. Yamasaki, and K. Aizawa, "Joint optimization framework for learning with noisy labels," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018, pp. 5552–5560. [2](#)
- [54] J. Li, G. Li, F. Liu, and Y. Yu, "Neighborhood collective estimation for noisy label identification and correction," in *Eur. Conf. Comput. Vis.*, 2022, pp. 128–145. [3, 8, 9](#)
- [55] D. Patel and P. S. Sastry, "Adaptive sample selection for robust learning under label noise," in *IEEE Winter Conference on Applications of Computer Vision*, 2023, pp. 3921–3931. [3](#)
- [56] X. Zhou, X. Liu, C. Wang, D. Zhai, J. Jiang, and X. Ji, "Learning with noisy labels via sparse regularization," in *Int. Conf. Comput. Vis.*, 2021, pp. 72–81. [3](#)
- [57] Z. Huang, J. Zhang, and H. Shan, "Twin contrastive learning with noisy labels," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2023, pp. 11 661–11 670. [3, 7, 8, 9](#)
- [58] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in *Int. Conf. Mach. Learn.*, 2018, pp. 4331–4340. [3](#)
- [59] 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 *Adv. Neural Inform. Process. Syst.*, 2019, pp. 1917–1928. [3](#)
- [60] 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, pp. 7024–7032. [3](#)
- [61] Y. Huang, B. Bai, S. Zhao, K. Bai, and F. Wang, "Uncertainty-aware learning against label noise on imbalanced datasets," in *AAAI*, 2022, pp. 6960–6969. [3](#)
- [62] X. Xia, T. Liu, B. Han, M. Gong, J. Yu, G. Niu, and M. Sugiyama, "Sample selection with uncertainty of losses for learning with noisy labels," in *Int. Conf. Learn. Represent.*, 2022. [3](#)
- [63] X. Xia, B. Han, Y. Zhan, J. Yu, M. Gong, C. Gong, and T. Liu, "Combating noisy labels with sample selection by mining high-discrepancy examples," in *Int. Conf. Comput. Vis.*, 2023, pp. 1833–1843. [3](#)
- [64] N. Karim, M. N. Rizve, N. Rahnavard, A. Mian, and M. Shah, "UNI-CON: combating label noise through uniform selection and contrastive learning," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022, pp. 9666–9676. [3, 9](#)
- [65] T. Sosea and C. Caragea, "Marginmatch: Improving semi-supervised learning with pseudo-margins," in *IEEE Conf. Comput. Vis. Pattern Recog.* IEEE, 2023, pp. 15 773–15 782. [5](#)
- [66] H. Chen, R. Tao, Y. Fan, Y. Wang, J. Wang, B. Schiele, X. Xie, B. Raj, and M. Savvides, "Softmatch: Addressing the quantity-quality trade-off in semi-supervised learning," in *Int. Conf. Learn. Represent.*, 2023. [5](#)
- [67] D. Berthelot, N. Carlini, I. J. Goodfellow, N. Papernot, A. Oliver, and C. Raffel, "Mixmatch: A holistic approach to semi-supervised learning," in *Adv. Neural Inform. Process. Syst.*, 2019, pp. 5050–5060. [5](#)
- [68] E. Malach and S. Shalev-Shwartz, "Decoupling "when to update" from "how to update"," in *Adv. Neural Inform. Process. Syst.*, 2017, pp. 960–970. [7, 8, 9](#)
- [69] X. Xia, T. Liu, B. Han, C. Gong, N. Wang, Z. Ge, and Y. Chang, "Robust early-learning: Hindering the memorization of noisy labels," in *Int. Conf. Learn. Represent.*, 2020. [7, 9](#)
- [70] Z. Sun, H. Liu, Q. Wang, T. Zhou, Q. Wu, and Z. Tang, "Co-ldl: A co-training-based label distribution learning method for tackling label noise," *IEEE Trans. Multimedia*, pp. 1093–1104, 2022. [7, 8, 9](#)
- [71] 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," in *AAAI*, B. Williams, Y. Chen, and J. Neville, Eds., 2023, pp. 14 765–14 773. [7, 8, 9](#)
- [72] A. Krizhevsky, "Learning multiple layers of features from tiny images," *Technical report, University of Toronto*, 2009. [7](#)
- [73] Z. Sun, X. Hua, Y. Yao, X. Wei, G. Hu, and J. Zhang, "CRSSC: salvage reusable samples from noisy data for robust learning," in *ACM Int. Conf. Multimedia*, 2020, pp. 92–101. [7, 8](#)
- [74] D. Hendrycks, M. Mazeika, S. Kadavath, and D. Song, "Using self-supervised learning can improve model robustness and uncertainty," in *Adv. Neural Inform. Process. Syst.*, 2019, pp. 15 637–15 648. [8, 9](#)
- [75] D. Mandal, S. Bharadwaj, and S. Biswas, "A novel self-supervised re-labeling approach for training with noisy labels," in *IEEE Winter Conference on Applications of Computer Vision*, 2020, pp. 1370–1379. [8, 9](#)
- [76] X. Peng, K. Wang, Z. Zeng, Q. Li, J. Yang, and Y. Qiao, "Suppressing mislabeled data via grouping and self-attention," in *Eur. Conf. Comput. Vis.*, 2020, pp. 786–802. [8, 9](#)
- [77] L. Huang, C. Zhang, and H. Zhang, "Self-adaptive training: beyond empirical risk minimization," in *Adv. Neural Inform. Process. Syst.*, 2020. [8, 9](#)
- [78] Z. Sun, Y. Yao, X. Wei, Y. Zhang, F. Shen, J. Wu, J. Zhang, and H. T. Shen, "Webly supervised fine-grained recognition: Benchmark datasets and an approach," in *Int. Conf. Comput. Vis.*, 2021, pp. 10 582–10 591. [8, 9](#)
- [79] S. Liu, Z. Zhu, Q. Qu, and C. You, "Robust training under label noise by over-parameterization," in *Int. Conf. Mach. Learn.*, 2022, pp. 14 153–14 172. [8, 9](#)
- [80] X. Shi, Z. Guo, K. Li, Y. Liang, and X. Zhu, "Self-paced resistance learning against overfitting on noisy labels," *Pattern Recognition*, vol. 134, p. 109080, 2023. [8, 9](#)
- [81] E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le, "Autoaugument: Learning augmentation strategies from data," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019, pp. 113–123. [8](#)
- [82] Y. Kim, J. Yun, H. Shon, and J. Kim, "Joint negative and positive learning for noisy labels," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021, pp. 9442–9451. [9](#)