# Tripartite: Tackle Noisy Labels by a More Precise Partition

Xuefeng Liang, Longshan Yao, Xingyu Liu, Ying Zhou.
School of Artificial Intelligence, Xidian University, China.

xliang@xidian.edu.cn

## Abstract

*Samples in large-scale datasets may be mislabeled due to various reasons, and Deep Neural Networks can easily over-fit to the noisy labeled data. The key solution is to alleviate the harm of these noisy labels. Many existing methods try to divide training data into clean and noisy subsets in terms of loss values, and then process the noisy labeled data variedly. We observe that a reason hindering a better performance is the hard samples. As hard samples usually have relatively large losses whether their labels are clean or noisy, these methods could not divide them accurately. Instead, we propose a Tripartite solution to partition training data into three subsets: hard, noisy, and clean. The partition criteria are based on the inconsistent predictions of two networks, and the inconsistency between the prediction of a network and the given label. To minimize the harm of noisy labels but maximize the value of noisy labeled data, we apply a low-weight learning on hard data and a self-supervised learning on noisy labeled data without using the given labels. Extensive experiments demonstrate that Tripartite can filter out noisy labeled data more precisely, and outperforms most state-of-the-art methods on five benchmark datasets, especially on real-world datasets.*

## 1. Introduction

Thanks to the large-scale datasets with human precisely annotated labels, DNNs achieve a great success. However, collecting high-quality and extensive data is considerably costly and time-consuming. To alleviate this issue, some cheaper alternatives are often employed, such as web-crawling [23], online queries [4, 20, 34], crowdsourcing [41,45] and so on. Unfortunately, they inevitably introduce some noisy labels. Many studies [3,19,33,46] have reported that DNNs could easily over-fit to the noises, which significantly degrades their generalization performance.

There have been many efforts to tackle noisy labels. Many of them reach a consensus that the key is to alleviate the impact of noisy labels for network training. Loss-correction based methods [9, 27, 28, 32] aim to rectify the losses of noisy labeled data in the training stage, but may mistakenly rectify some clean data. Sample-selection based methods [10, 18, 37, 44] tackle noisy labels by partitioning the training data into clean and noisy subsets, then using them for network training in different ways. The mainstream partition criteria are two types: 1) *Small-loss criterion* [10] assumes that the losses of noisy labeled data are significantly higher than those of clean data during training. Therefore, they try to find a threshold, $T$, and select the samples, whose losses $< T$, as clean data. The others are treated as noisy labeled data. 2) *Gaussian Mixture Model* (GMM) *criterion* [18] believes that the statistical distribution of noisy data losses is different from that of clean data losses. It aims to find the probability of a sample being noisy or clean by fitting a mixture model.

However, we observe that many real-world noisy labels are introduced between similar categories, especially happen among hard samples. For instances, a dolphin is mislabeled as a whale, vice versa, as shown in Fig. 1. In this paper, we define "hard samples" as data that distribute close to the decision boundary and are difficult to be distinguished. Our investigation shows, regardless of clean or noisy, the training losses of hard samples are neither small nor significantly different. So, existing sample-selection based methods have two flaws: 1) Low quality of training data partition. They are likely to mistake the hard noisy labeled samples as clean ones solely based on losses, and vice versa. This downgrades the performance of a network because the network will learn certain noisy labels and discard some clean data. 2) Ineffective usage of noisy labeled data. Small-loss criterion often drops the noisy samples which is a waste of valuable information. GMM criterion applies the semi-supervised learning to reuse noisy labeled samples by assigning pseudo-labels to them. But it heavily relies on the discriminative ability of networks. The incorrect relabeling will cause a severe harm to networks as well.

To tackle above problems, we propose a novel method: *Tripartite* that mainly addresses hard samples in real-world datasets. Since hard samples distribute around the decision boundary, predictions of varied networks are often inconsistent at the early training stage. Meanwhile, the prediction
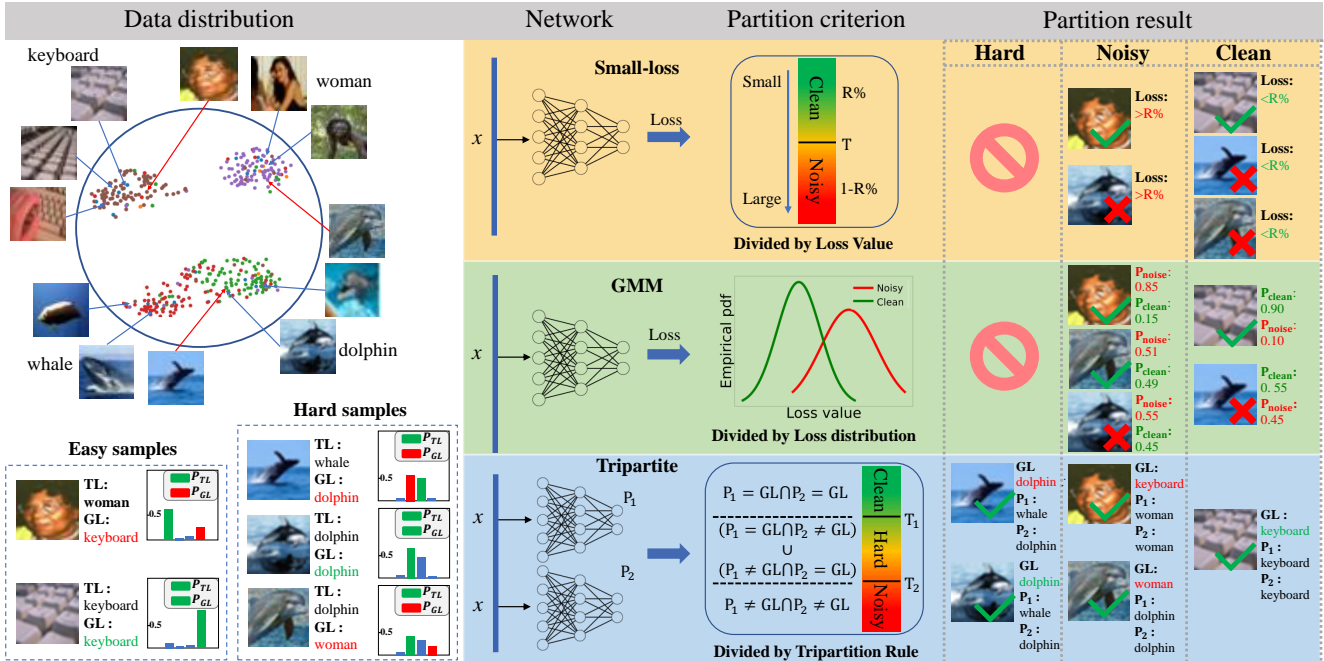
Figure 1. Comparison with existing partition criteria. From top to bottom are Small-loss criterion, GMM criterion and the proposed Tripartition criterion. TL denotes the true label, GL denotes the given label. $P_1$ and $P_2$ are the predicted labels of inputs. $P_{TL}$ and $P_{GL}$ denote the predicted probabilities of true and given labels, respectively. $P_{clean}$ and $P_{noise}$ denote the predicted probabilities of a sample to be clean and noisy, respectively. In the partition result, the check mark and cross mark denote that the sample is partitioned correctly and incorrectly, respectively. Small-loss criterion sorts losses and selects R% samples with small losses as clean data for training. GMM criterion finds the probability of a sample being clean or noisy by fitting a mixture model. Neither of them can filter out hard samples. By contrast, Tripartite divides the training set into hard, noisy, and clean subsets according to the relations between $P_1$, $P_2$ and GL.

of a network of an easy noisy sample is usually inconsistent with the given label. Based on this observation, Tripartite can divide training data into hard, noisy and clean subsets. It has the advantage of improving the quality of clean and noisy subsets. To effectively use data, we apply a low-weight training strategy for samples in the hard subset, and employ a self-supervised training strategy for samples in the noisy subset without using the given labels. Hence, Tripartite is designed toward minimizing the harm of noisy labels and maximizing the value of noisy labeled data. The extensive experiments on five benchmark datasets demonstrate the superior performance of Tripartite compared with the state-of-the-art (SOTA) methods. Especially, it shows robustness at a wide range of real-world datasets. The key contributions of our work are threefold.

- We propose a novel partition criterion, which divides training data into three subsets: hard, noisy, and clean. It alleviates the hard sample selection problem of other criteria, and largely improves the quality of clean and noisy subsets.

- We design a low-weight training strategy for hard data and a self-supervised training strategy for noisy labeled data, which aim at minimizing the harm of noisy labels and maximizing the value of noisy labeled data.

- To mimic the noisy label of hard sample in real-world datasets, we create a synthetic class-dependent label noise on CIFAR datasets, called realistic noise. It flips labels of samples, which are from two different classes, at controlled ratios according to their similarity. The details are shown in Supplementary.

## 2. Related work

There have been many studies to address noisy labels. They all try to lower the impact of noisy labels to improve the recognition performance of methods during the training stage. To this end, two different ideas were proposed.

**Loss correction**. The specific methods include *noise transition matrix* [25–27, 38], *robust loss functions* [8, 35, 36, 40, 48], *label correction* [21, 32, 43], etc. *Noise transition matrix* methods construct the label transition matrix to estimate the possibilities of noisy labels transiting among multiple classes. F-correction [27] performs forward correction by multiplying the transition matrix with the softmax outputs in the forward propagation. Later, Hendrycks *et al.* [13] improved the corruption matrix using a small

clean dataset. *Robust-loss* methods aim to design loss functions that are robust to noisy labels, such as Mean Absolute Error (MAE) [8], Improved MAE [35] that is a reweighted MAE, Generalized Cross-Entropy loss (GCE) [48] that is a generalization of MAE, Symmetric Cross-Entropy [36] that adds a reverse cross-entropy term to the usual cross-entropy loss. *Label correction* methods try to rectify the noisy labels according to the network predictions during the training. Joint optimization [32] learns network parameters and infers the true labels simultaneously. PENCIL [43] adopts label probability distributions to supervise network learning and update these distributions through back-propagation in each epoch. Inspired by "early learning", ELR+ [21] uses the network predictions in the early training stage to correct the noisy labels.

Above methods do not distinguish training data into clean and noisy subsets. They may mistakenly rectify the losses of clean data and introduce new noisy labels into training data. Therefore, sample selection methods were proposed.

**Sample selection**. These methods try to divide training data into "clean" and "noisy" subsets according to a specific partition criterion. Afterwards, they apply different strategies to train the model on the two subsets separately. The existing partition criteria are mainly based on the training loss, e.g. *small-loss criterion* [10] and *Gaussian Mixture Model* (GMM) *criterion* [18]. *Small-loss criterion* selects training samples with small loss as clean ones by setting a threshold $T$. In particular, MentorNet [14] pretrains a teacher network for selecting clean samples to guide the training of a student network. Co-teaching [10] trains two networks where each network selects small-loss samples to feed its peer network for updating parameters. Furthermore, Co-teaching+ [44] emphasizes the help of inconsistency to the network, and JoCoR [37] emphasizes the consistency of the two networks. However, all these methods only use the data with smaller losses for training without considering the data with larger losses. Meanwhile, finding a feasible $T$ is very challenging. GMM assumes that the distributions of losses of noisy labeled data and clean data follow two normal distributions, respectively. DivideMix [18] employs two networks, each selects data for another in the training, and then applies a semi-supervised method to learn the noisy labeled data. It greatly improves the performance of networks. However, in real-world datasets, the noisy labels often appear between similar classes, such as *whale* and *dolphin*, *oak_tree* and *maple_tree* in CIFAR-100. GMM may mistake the hard noisy samples and the hard clean samples because of the little difference of their losses. It leads to a low quality of training data partition. By contrast, we propose the Tripartition criterion that divides training data into clean, hard, and noisy subsets. A more precise partition can effectively lower the harm of noisy labels for network

training.

## 3. The proposed method

To explain the mechanism of Tripartite, we firstly analyze the distributions of hard samples and noisy labeled data, and then derive the logic of data partition in Sec. 3.2. Sec. 3.3 details the proposed Tripartition criterion. Sec. 3.4 presents the training strategies for data in three subsets. The pseudo-code of Tripartite is shown in Algorithm 1.

### 3.1. Preliminaries

In our method, we train two networks, denoted by $f(\theta_k), k = \{1, 2\}$, which are initialized with different weights. We assume that they could learn different views that include reliable information on noisy training data but produce inconsistent predictions on hard samples. $P_k$ is the predicted label of $f(x, \theta_k)$ on an input $x$. The proposed methods mentioned hereinafter are all based on the predictions of these two networks. Studies [3, 46] report that DNNs tend to learn simple patterns first before memorizing noisy labeled data. And, DNN's optimizations are content-aware, taking advantage of patterns shared by multiple training samples [21, 46]. Therefore, we assume that networks do not over-fit to the noisy labeled samples in the early learning stage because noisy labels are less related to the data features themselves.

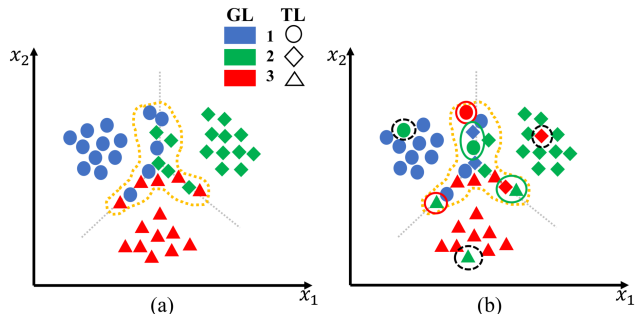### 3.2. Distribution of training data



Figure 2. The distribution of training data. Colors denote the given labels; Shapes denote the true labels. (a) Data with clean labels. Hard samples are grouped by a yellow dotted line. (b) Data with noisy labels. The hard samples with noisy labels from similar/dissimilar classes are circled in green/red. The easy samples with noisy labels are circled in black.

Training data include hard and easy samples. As hard samples share certain features of two or more similar classes, they distribute very close to the decision boundaries. Fig. 2 (a) illustrates the data distribution, where hard samples are grouped by a yellow dotted line, the others are easy samples.

3

For easy samples, there are two cases. 1) *Easy samples with clean labels* whose color and shape are corresponded in Fig. 2 (b). As the label information is rather correlated with the sample features, networks can fit these samples in the early training stage. Their losses should be the smallest in the training set. 2) *Easy samples with noisy labels*, circled in black in Fig. 2 (b), for instance, a lady is mislabeled as keyboard in Fig. 1. As the information provided by the given labels is irrelevant with the features of these samples, they are anomalies in training data. Their losses should be the largest in the training set.

For hard samples, there are three cases. 1) *Hard samples with clean labels*. The label information and sample features are correlated. As distributing close to the decision boundary, they also have some shared features with the similar classes. Therefore, the features of these data are related to both two classes. 2) *Hard samples with noisy labels from similar classes*, circled in green in Fig. 2 (b). They distribute close to the decision boundaries and are mislabeled to similar classes. Thus, their features are related to these two classes to some extent. Our observation shows that the losses of both hard samples are greater than the losses of easy clean samples but less than the losses of easy noisy labeled samples. 3) *Hard samples with noisy labels from dissimilar classes*, circled in red in Fig. 2 (b), for example, a dolphin is mislabeled as woman. These samples also distribute close to the decision boundaries but are mislabeled to dissimilar classes. The information provided by the given label is irrelevant with the features of these samples, which is analogous to the case of easy noisy samples.

Above analysis illustrates five possible cases of training data, two of easy samples and three of hard samples. In order to improve the quality of data partition and minimize the harm of noisy labels, we then group the five cases into three subsets according to the losses. They are: *hard subset* including hard samples with clean labels, and hard samples with noisy labels from similar classes; *noisy subset* including easy samples with noisy labels, and hard samples with noisy labels from dissimilar classes; *clean subset* including easy samples with clean labels. One can see there exists an order: losses of noisy subset > losses of hard subset > losses of clean subset. In practice, there are two thresholds $T_1$ and $T_2$ in each epoch, where $T_1 < T_2$. Then, training data can be partitioned shown as in Fig. 1.

### 3.3. Tripartition method and criteria

Theoretically, training data can be partitioned into three subsets according to $T_1$ and $T_2$. Nevertheless, selecting the feasible $T_1$ and $T_2$ is non-trivial. Our investigation shows the predictions of different networks on hard samples usually are inconsistent. There exists a correlation between the losses and the network predictions of training data, which explains why bipartition methods are difficult to handle
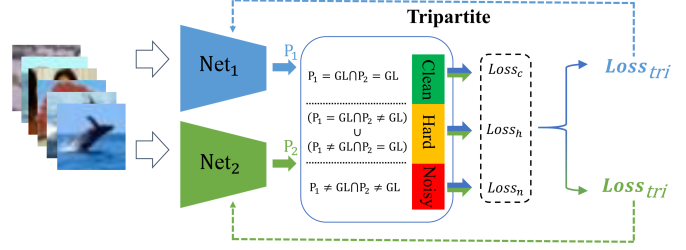


Figure 3. The logic of our Tripartition method. GL denotes the given label.

hard samples. Please refer to Supplementary. Therefore, determining $T_1$ and $T_2$ becomes measuring the consistency of predictions of different networks. Thus, we propose a Tripartition method and the selection criteria. The logic of Tripartite is shown in Fig. 3. We train two networks with different initializations. They give their own predictions in each epoch. Training data are partitioned into three subsets accordingly.

**Selection criterion for hard subset**

$$(P_1 = GL \cap P_2 \neq GL) \cup (P_1 \neq GL \cap P_2 = GL) \quad (1)$$

Hard subset mainly includes the hard samples with clean labels and the hard samples with noisy labels from similar classes. The features of these samples and the information provided by labels are somewhat correlated but not very consistent. So, different networks may not give a consistent prediction on such samples, especially at the early training stage. We then apply the inconsistent predictions of two networks, Eq. (1), to select them. The losses of selected data should follow $T_1 < loss < T_2$. Our experiments indicate that the population of these data will gradually decrease with the increase of the discriminative performance of networks. Please refer to Supplementary for details.

**Selection criterion for noisy subset**

$$P_1 \neq GL \cap P_2 \neq GL \quad (2)$$

Noisy subset mainly includes the easy samples with noisy labels and the hard samples with noisy labels from dissimilar classes. As the features of these samples and the information provided by labels are irrelevant and inconsistent, the networks will not fit these data in the early training stage. So, the prediction of a network would not be consistent with the given label. We then apply this inconsistency, Eq. (2), to select them. The losses of selected data should follow $loss > T_2$.

**Selection criterion for clean subset**

$$P_1 = GL \cap P_2 = GL \quad (3)$$

Clean data mainly include the easy samples with clean labels. As the features of these samples and the label information are rather consistent, networks can learn the mapping between data features and labels in the early training stage. So, the predictions of two networks and the given label should be consistent. We then apply this consistency, Eq. (3), to select clean data. The losses of selected data should follow $loss < T_1$.

### 3.4. Learning strategies for three subsets

We design different learning strategies for three subsets, respectively. Let's consider a classification problem with $C$ classes. The training set consists of $n$ examples, $\{x^i, y^i\}$ is the $i$th input data, $y^i \in \{0,1\}^C$ is a one-hot label of $C$ dimensions vector corresponding to the class of $x^i$. The network maps the input $x^i$ into a $C$-dimensional encoding and feeds it into a softmax function to estimate the probability $p^i$ of $x^i$ to each class.

**Learning strategy for hard subset**

The hard data subset includes both samples with noisy labels and clean labels. It is difficult to distinguish them. We wish to utilize the valuable information of hard data to improve the discriminative ability of networks, meanwhile, lower the harm of noisy labels in this subset. Therefore, we apply a low-weight cross-entropy loss function for hard data.

$$Loss_h = \lambda_h(-\frac{1}{n}\sum_{i=1}^{n}\sum_{c=1}^{C}y_c^i log\left(p_c^i\right)), \qquad (4)$$

where, $\lambda_h$ is a weight in the range (0, 1). $p_c^i$ is the estimated probability of $x^i$ to be the class $c$.
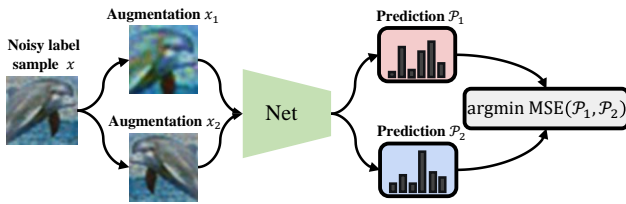


Figure 4. The logic of our self-supervised strategy. It learns samples in noisy subset without using the given labels.

**Learning strategy for noisy subset**

To lower the impact of noisy labels, many works apply semi-supervised learning to replace the given labels by pseudo-labels. However, the quality of pseudo-labels heavily relies on the discriminative ability of networks. New noisy labels may be introduced into training data when the discriminative ability is weak. Instead, we design a self-supervised method to learn these data without using the given labels.

Given a noisy labeled sample $x$ shown in Fig. 4, we randomly augment it by rotation, flip, cropping, desaturation, contrast, blurring, and MixUp [47], etc. These augmentations can create many positive sample-pairs, $\langle \bar{x}, \hat{x} \rangle$. They are fed into the network, which predicts a probability distribution, $\mathcal{P} = (p_1, p_2, ..., p_C)$, for each augmentation. As aiming at learning the common discriminative features from similar samples, we constrain the consistency of two probability distributions $\langle \bar{\mathcal{P}}, \hat{\mathcal{P}} \rangle$ using the MSE loss function.

$$Loss_n = \frac{1}{n}\sum_{i=1}^{n}(\bar{\mathcal{P}}^i - \hat{\mathcal{P}}^i)^2. \qquad (5)$$

Please note that two networks still learn a bit different features because of the random augmentation in each batch. It helps the two networks to give inconsistent predictions on most of the hard samples.

**Learning strategy for clean subset**

As Tripartite ensures the high quality of samples in clean subset, we apply the cross-entropy loss for clean data.

$$Loss_c = -\frac{1}{n}\sum_{i=1}^{n}\sum_{c=1}^{C}y_c^i log\left(p_c^i\right). \qquad (6)$$

Finally, the total loss is

$$Loss_{tri} = Loss_c + \lambda_n Loss_n + Loss_h, \qquad (7)$$

where, $\lambda_n$ is a parameter to control the contribution of noisy subset.

---

**Algorithm 1:** Tripartite

**Input:** Network $f(\theta_k), k = \{1,2\}$, learning rate $\eta$, loss weight $\lambda_h$ of hard subset, loss weight $\lambda_n$ of noisy subset, end epoch of warm up $T_k$, and epoch $T_{max}$, iteration $I_{max}$

1 **Warm up** $f(\theta_1), f(\theta_2)$
2 **for** $t = T_k, ..., T_{max}$ **do**
3 $\quad$ $P_1 = f(x, \theta_1), \forall x \in D$
4 $\quad$ $P_2 = f(x, \theta_2), \forall x \in D$
5 $\quad$ **Obtain** $Tripart = \textbf{Tripartition}(P_1, P_2, D)$
6 $\quad$ **for** $k = 1, 2$ **do**
7 $\quad\quad$ **for** $i = 1, 2, ..., I_{max}$ **do**
8 $\quad\quad\quad$ **Fetch** mini-batch $D_n$ from $D$
9 $\quad\quad\quad$ $D_{clean}, D_{hard}, D_{noise} = Tripart(D_n)$
10 $\quad\quad\quad$ **Calculate** $Loss_c$ by Eq. (6),$\forall x \in D_{clean}$
11 $\quad\quad\quad$ **Calculate** $Loss_h$ by Eq. (4),$\forall x \in D_{hard}$
12 $\quad\quad\quad$ **Obtain** $x_1, x_2 = $ Augment$(x),\forall x \in D_{noise}$
13 $\quad\quad\quad$ **Calculate** $Loss_n$ by Eq. (5)
14 $\quad\quad\quad$ **Calculate** $Loss_{tri}$ by Eq. (7)
15 $\quad\quad\quad$ **Update** $\theta_k = \theta_k - \eta\nabla Loss_{tri}$
16 $\quad\quad$ **end**
17 $\quad$ **end**
18 **end**

**Output:** $\theta_1, \theta_2$

---

## 4. Experiments

### 4.1. Implementation details and datasets

**Datasets and noise types:** We verify the effectiveness of our Tripartite extensively on two benchmark datasets (CIFAR-10 [16], CIFAR-100 [16]) and three real-world datasets (Food-101N [17], Clothing1M [39], WebVision1.0 [20]).

Both CIFAR-10 and CIFAR-100 contain 60,000 images of size 32×32. CIFAR-10 includes 10 super-classes. CIFAR-100 has 100 subclasses. We validate two noise settings on them. In symmetric noise setting, we randomly flipping the labels of a given proportion (20%, 50%, 80%) to other labels uniformly in the training set. In realistic noise setting, we designed a new method that ensures the labels are replaced based on the similarities between classes. The details are shown in Supplementary.

Food-101N is a benchmark for visual food classification using the Food-101 [5] taxonomy. It contains about 310,000 images of food recipes classified into 101 classes. The estimated noise ratio is about 20%. Clothing1M consists of 1 million training images collected from online shopping websites with labels generated with the surrounding text. Its noise level is about 38.5%. WebVision1.0 is a large-scale dataset with real-world noisy labels. It contains more than 2.4 million images crawled from the web using the 1,000 concepts in ImageNet ILSVRC2012 [7]. Its noise level is about 20%. For ease of comparison with competing methods, we follow the previous work [6] and then do the test on a subset which is the first 50 classes crawled from Google image.

**Models and parameters:** Following previous works, we use an 18-layer PreAct Resnet [12] for CIFAR-10, CIFAR-100, and train it by the SGD optimizer with a learning rate of 0.02, a momentum of 0.9 and a weight-decay of 0.0005. The batch size is set to 128. We train the network for 300 epochs. The architectures of the two networks used in our method are the same. The learning rate is reduced by a factor of 10 after 150 and 200 epochs. We set the warm-up epochs to 10 for CIFAR-10 and to 30 for CIFAR-100. The most critical parameters in Tripartite are $\lambda_n$ and $\lambda_h$. The $\lambda_n$ is selected from $\{1, 10, 20, 40, 60, 80, 100\}$, $\lambda_h$ is 0.6 in both CIFAR datasets.

For real-world datasets, the experiment settings and parameters are the same with CIFAR datasets except training epochs, learning rate and network architectures. We train the network for 100 epochs. The batch size is set to 64. The learning rate is reduced by a factor of 10 after 30 and 60 epochs. We set the warm-up epochs to 5 for Food-101N and Clothing1M, and 30 for WebVsion1.0. Resnet-50 [11] pretrained on ImageNet (following previous work [19]) is for Food-101N and Clothing1M, and inception-resnet v2 [30] is for WebVsion1.0.

### 4.2. Comparison with SOTA methods

We compare our method with the baseline Standard CE (the backbone + Cross Entropy), and the recent SOTA methods including Bootstrap [29], F-correction [27], P-correction (PENCIL) [43], Meta-Learning [19], M-correction [2], Co-teaching [10], JoCoR [37], DivideMix [18], ELR+ [21], Co-learning [31], DSOS [1] and JNPL [15]. We repeat following experiments 5 times, and report the average test accuracy of 5 trials for each experiment.

| Dataset *symmetric* Method | CIFAR-10 | | | CIFAR-100 | | |
|---|---|---|---|---|---|---|
| | 20% | 50% | 80% | 20% | 50% | 80% |
| Standard CE | 86.8 | 79.4 | 62.9 | 62.0 | 46.7 | 19.9 |
| Bootstrap(2015) [29] | 86.8 | 79.8 | 63.3 | 62.1 | 46.6 | 19.9 |
| F-correction(2017) [27] | 86.8 | 79.8 | 63.3 | 61.5 | 46.6 | 19.9 |
| Co-teaching+(2019) [44] | 89.5 | 85.7 | 67.4 | 65.6 | 51.8 | 27.9 |
| P-correction(2019) [43] | 92.4 | 89.1 | 77.5 | 69.4 | 57.5 | 31.1 |
| Meta-Learning(2019) [19] | 92.9 | 89.3 | 77.4 | 68.5 | 59.2 | 42.4 |
| M-correction(2019) [2] | 94.0 | 92.0 | 86.8 | 73.9 | 66.1 | 48.2 |
| DivideMix(2020) [18] | <u>96.1</u> | 94.6 | <u>93.2</u> | 77.3 | <u>74.6</u> | <u>60.2</u> |
| ELR+(2020) [21] | 95.8 | <u>94.8</u> | **93.3** | <u>77.6</u> | 73.6 | **60.8** |
| Co-learning(2021) [31] | 92.5 | 84.83 | 63.5 | 66.7 | 55.0 | 36.2 |
| DSOS(2022) [1] | 92.7 | 87.4 | 62.7 | 75.1 | 66.2 | 38.0 |
| Tripartite | **96.3** | **94.9** | 92.6 | **78.7** | **74.7** | 59.8 |

Table 1. Accuracy comparisons on CIFAR-10 and CIFAR-100 with symmetric noise. We ran the open code of DSOS, and report the results. The results of ELR+ and Co-learning are from their papers, the rest of results are from [18]. The best results are in bold and the second-best results are underlined.

| Dataset *realistic* Method | CIFAR-10 | | | CIFAR-100 | | |
|---|---|---|---|---|---|---|
| | 20% | 40% | 50% | 20% | 40% | 50% |
| Co-teaching(2018) [10] | 82.7 | 74.4 | 55.8 | 50.3 | 40.9 | 32.5 |
| JoCoR(2020) [37] | 82.2 | 68.7 | 55.8 | 49.7 | 35.1 | 29.1 |
| DivideMix(2020) [18] | 91.1 | 92.3 | <u>91.8</u> | <u>76.2</u> | 66.1 | 59.5 |
| ELR+(2020) [21] | <u>95.1</u> | <u>92.9</u> | 91.2 | 75.8 | <u>72.5</u> | <u>60.6</u> |
| Co-learning(2021) [31] | 91.6 | 75.5 | 61.7 | 68.5 | 58.5 | 50.2 |
| DSOS(2022) [1] | 92.5 | 89.0 | 81.1 | 75.3 | 64.8 | 52.6 |
| Tripartite | **96.2** | **95.9** | **94.4** | **78.8** | **74.9** | **68.6** |

Table 2. Accuracy comparisons on CIFAR-10 and CIFAR-100 with realistic noise. We strictly implemented the competing methods according to their open codes and papers. The backbone is the Pre-ResNet18.

**Results on CIFAR-10 and CIFAR-100**

Tab. 1 shows the test accuracy on CIFAR-10 and CIFAR-100 with different levels of symmetric noise ranging from 20% to 80%. We can see that Tripartite outperforms competing methods on 20% and 50% settings. However, it is slightly worse than the top two methods on 80% setting. The reason is that most networks have difficulties of learning valuable information during the early training stage when the noise ratio is too high (80%). This hinders Tripartite to divide training data well. Since Tripartite aims at real-world noisy data, in which the noise ratio over 50% is
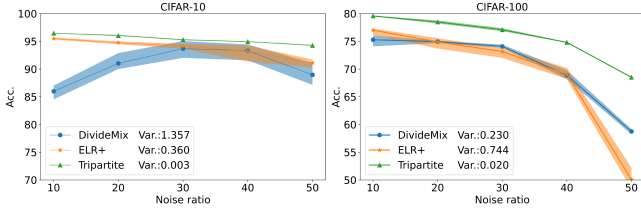
Figure 5. Accuracy comparisons among the top 3 methods on CIFAR with realistic noise. The noise ratio ranges from 10% to 50%. The solid lines are the average accuracies. The shaded regions denote the variances.

unlikely to happen [31], we focus on the data with the noise ratio less than 50%.

Tab. 2 shows the test accuracy on CIFAR-10 and CIFAR-100 with varied levels of realistic noise. We can see that Tripartite achieves the best across all noise levels. To have a more intuitive comparison, we plot the top 3 methods in Fig. 5. Compared with DivideMix and ELR+, Tripartite steadily improves the accuracies on every noise ratio, and considerably outperforms them about 3.08% on CIFAR-10 and 4.54% on CIFAR-100. Moreover, the variance of Tripartite performances is much lower than those of DivideMix and ELR+. This experiment confirms the superiority of Tripartite on hard noisy labeled data.

| Food-101N | | Clothing1M | |
|---|---|---|---|
| Methods | Acc. | Methods | Acc. |
| Standard CE | 84.03 | Standard CE | 69.21 |
| CleanNet [17] | 83.95 | F-correction [27] | 69.84 |
| Decoupling [24] | 85.53 | Joint [32] | 72.16 |
| Co-teaching [10] | 61.91 | Meta-Learning [19] | 73.47 |
| Co-teaching+ [44] | 81.61 | P-correction [43] | 73.49 |
| JoCoR [37] | 77.94 | DivideMix [18] | 74.76 |
| Jo-SRC [42] | 86.66 | ELR+ [21] | 74.81 |
| Co-learning [31] | 87.57 | JNPL [15] | 74.15 |
| DSOS [1] | 87.70 | DSOS [1] | 73.63 |
| Tripartite | **88.34** | Tripartite | **75.23** |

Table 3. Accuracy comparison on the Food-101N. The results of Jo-SRC and DSOS are from [42] and [1]. The rest of results are from [31].

Table 4. Accuracy comparison on the Clothing1M. The results of JNPL and DSOS are from [15] and [1]. The rest of results are from [21].

**Results on real-world datasets**

Tabs. 3 to 5 list the results on real-world noisy datasets, Food-101N, Clothing1M and WebVision1.0, respectively. We also evaluate the generalization ability of these methods on ImageNet ILSVRC12. Tripartite consistently outperforms competing methods across all large-scale datasets. One can see the performance gains are not very significant on Food-101N and Clothing1M. The possible reason is that images in them often contain more than one subject. Especially, images in Clothing1M usually are full-body shots that include top and bottom clothing but are labeled by either of them. In this case, networks are difficult to learn the

| Dataset | WebVision | | ILSVRC12 | |
|---|---|---|---|---|
| Method | Top-1 | Top-5 | Top-1 | Top-5 |
| F-correction(2017) [27] | 61.12 | 82.68 | 57.36 | 82.36 |
| Decoupling(2017) [24] | 62.54 | 84.74 | 58.26 | 82.26 |
| D2L(2018) [22] | 62.68 | 84.00 | 57.80 | 81.36 |
| MentorNet(2018) [14] | 63.00 | 81.40 | 57.80 | 79.92 |
| Co-teaching(2018) [10] | 63.58 | 85.20 | 61.48 | 84.70 |
| Iterative-CV(2019) [6] | 65.24 | 85.34 | 61.60 | 84.98 |
| DivideMix(2020) [18] | 77.32 | 91.64 | 75.20 | 90.84 |
| ELR+(2020) [21] | 77.78 | 91.68 | 70.29 | 89.76 |
| DSOS(2022) [1] | 77.76 | 92.04 | 74.36 | 90.80 |
| Tripartite | **78.96** | **93.20** | **75.92** | **92.88** |

Table 5. Top-1 (Top-5) accuracy comparison on the WebVision1.0 and ILSVRC12. The result of DSOS is from [1], the rest of results are from [21].

discriminative feature. Hence, Tripartite may partition data with some uncertainty.

As WebVision has more and diverse categories, it allows Tripartite to achieve a considerable improvement than other methods. Please note that ELR+ performs well on WebVision test set, but much worse on ILSVRC12 Top-1 validation set. On the contrary, Tripartite demonstrates its robustness and steady superiority on both WebVision and ILSVRC12.

## 4.3. Ablation study

To verify the effectiveness of our training strategies, the ablation study is conducted on CIFAR-100 in symmetric and realistic settings with varied noise ratios.
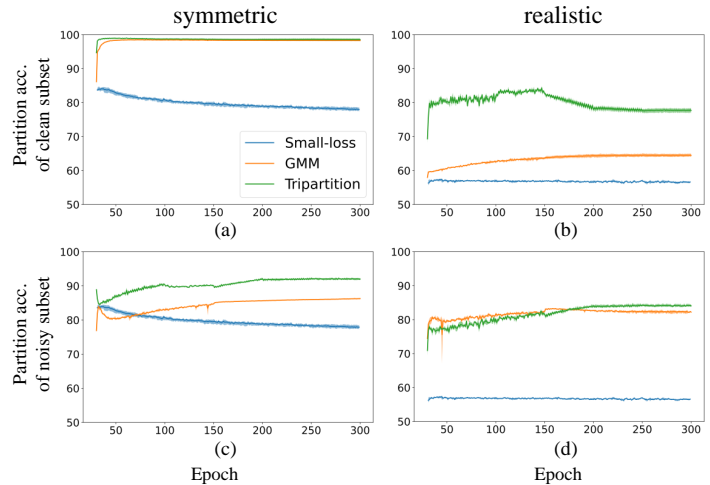
**The quality of training data partition**



Figure 6. Comparison of three partition criteria on the CIFAR-100 dataset with symmetric and realistic noise types and 50% noise ratio. (a) and (b) plot the partition accuracy of clean samples in clean subset against training epochs. (c) and (d) plot the partition accuracy of noisy labeled samples in noisy subset against training epochs.

To test the partition accuracies of clean and noisy sub-

| $\lambda_h =$ Noise type | | 2 | 1 | 0.8 | 0.6 | 0.4 | 0.2 |
|---|---|---|---|---|---|---|---|
| Symmetric | 20% | 73.78 | 76.62 | 77.54 | **78.72** | 78.4 | 76.91 |
| | 50% | 60.12 | 66.19 | 70.29 | **74.73** | 72.92 | 71.06 |
| Realistic | 20% | 75.65 | 78.07 | 78.35 | **78.76** | 78.22 | 77.08 |
| | 50% | 59.67 | 64.8 | 65.54 | **68.62** | 67.17 | 65.31 |

sets, we compare Tripartite with *small-loss criterion* and GMM *criterion*. The results are shown in Fig. 6. We can see Tripartite achieves the best partition accuracies on both clean subsets. On the symmetric noise, our Tripartite outperforms the GMM *criterion* slightly. The reason is that the symmetric noises are generated by a symmetric label transfer. The noisy label and clean label can fit the two Gaussian distributions in GMM better in terms of losses. The *small-loss criterion* is the worst because it heavily relies on the selection of R%. It is worth noting that Tripartite performs better than others, although the population of hard samples with noisy labels from similar classes is smaller on the symmetric type. On the realistic noise, our Tripartite largely outperforms other two criteria on clean subset. As the population of hard noisy labeled data becomes much larger, it narrows the difference between losses of clean and noisy labeled data in the training set, and then further weakens the discriminations of two loss-based criteria.

The comparison of results in Tab. 7 and Tab. 2 also shows that the higher quality of data partition does improve the performance of a model. On realistic noise with a ratio of 50%, we apply the strategy of co-teaching and JoCoR which trains network without using the noisy subset partitioned by Tripartite. The result (65.04% in Tab. 7) is much higher than the results (32.5% and 29.1% in Tab. 2) of Co-teaching and JoCoR. A similar trend can be found on the strategy of DivideMix, which applies the semi-supervised learning on the noisy subset partitioned by Tripartite. The result (65.45% in Tab. 7) is higher than the results (59.5% in Tab. 2) of DivideMix.

**Training strategy for hard subset**

| Noise type Strategy | Symmetric | | Realistic | |
|---|---|---|---|---|
| | 20% | 50% | 20% | 50% |
| Drop | 78.01 | 68.69 | 77.92 | 65.04 |
| Semi-supervised | 78.23 | 73.79 | 78.06 | 65.45 |
| Self-supervised | **78.72** | **74.73** | **78.76** | **68.62** |

Table 7. Evaluation of training strategy for noisy subset.

The weight $\lambda_h$ is a key parameter for training hard subset. Our motivation is to lower the harm of noisy labels, therefore, we test the $\lambda_h$ in the range of (0,2], and list results in Tab. 6. Obviously, $\lambda_h = 2$ reaches the worst result because it amplifies the harm. When $\lambda_h = 0.6$, we have the best performance. However, $\lambda_h = 0.2$ results in a worse result. The possible reason is that there exist some

hard clean samples, which are critical to the performance of the network. If $\lambda_h$ is too low, the network learns little valuable information from them.

**Training strategy for noisy subset**

To evaluate the effectiveness of training strategy for noisy subset, we test three strategies and report the results in Tab. 7. One can see that our self-supervised strategy achieves the best because it is label free and learns the common features from similar samples. It can minimize the harm of noisy labels and fully use these data. The $\lambda_n$ is not highly sensitive to Tripartite, see details in Supplementary. The semi-supervised strategy is worse because it generates pseudo-labels for noisy labeled data to reuse them. The quality of a pseudo-label heavily relies on the performance of the networks. Some incorrect pseudo-labels will degrade the final performance. The strategy of dropping noisy labeled data performs the worst, because it does not use these data.

## 5. Conclusion

Existing loss-based methods divide training data into clean and noisy subsets, but may mistake the hard noisy labeled samples and the hard clean samples. To further improve the quality of two subsets, we propose a novel method, Tripartite, to mainly addresses hard samples. Unlike bipartition, Tripartite partitions training data into three subsets: *hard*, *noisy*, and *clean*. The partition criteria are based on the inconsistent predictions of two networks, and the inconsistency between the prediction of a network and the given label. To minimize the harm of noisy labels but maximize the value of noisy labeled data, Tripartite applies a low-weight strategy on hard data and a self-supervised strategy on noisy labeled data. The extensive experiments demonstrate that Tripartite outperforms SOTA methods and provides a new perspective of data partition and data utilization.

## References

[1] Paul Albert, Diego Ortego, Eric Arazo, Noel E O'Connor, and Kevin McGuinness. Addressing out-of-distribution label noise in webly-labelled data. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 392–401, 2022.

[2] Eric Arazo, Diego Ortego, Paul Albert, Noel O'Connor, and Kevin McGuinness. Unsupervised label noise modeling and loss correction. In *ICML*, pages 312–321. PMLR, 2019.

[3] Devansh Arpit, Stanisław Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. In *ICML*, pages 233–242. PMLR, 2017.

[4] Avrim Blum, Adam Kalai, and Hal Wasserman. Noise-tolerant learning, the parity problem, and the statistical query model. *Journal of the ACM*, 50(4):506–519, 2003.

[5] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101–mining discriminative components with random forests. In *ECCV*, pages 446–461. Springer, 2014.

[6] Pengfei Chen, Ben Ben Liao, Guangyong Chen, and Shengyu Zhang. Understanding and utilizing deep neural networks trained with noisy labels. In *ICML*, pages 1062–1070. PMLR, 2019.

[7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. IEEE, 2009.

[8] Aritra Ghosh, Himanshu Kumar, and PS Sastry. Robust loss functions under label noise for deep neural networks. In *AAAI*, volume 31, pages 1919–1925, 2017.

[9] Jacob Goldberger and Ehud Ben-Reuven. Training deep neural-networks using a noise adaptation layer. In *ICLR*, 2017.

[10] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor W Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *NeurIPS*, 2018.

[11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778. IEEE, 2016.

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *ECCV*, pages 630–645. Springer, 2016.

[13] Dan Hendrycks, Mantas Mazeika, Duncan Wilson, and Kevin Gimpel. Using trusted data to train deep networks on labels corrupted by severe noise. In *NeurIPS*, 2018.

[14] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *ICML*, pages 2304–2313. PMLR, 2018.

[15] Youngdong Kim, Juseung Yun, Hyounguk Shon, and Junmo Kim. Joint negative and positive learning for noisy labels. In *CVPR*, pages 9442–9451. IEEE, 2021.

[16] Alex Krizhevsky and Geoffrey E. Hinton. Learning multiple layers of features from tiny images. *Technical Report, University of Toronto*, 2009.

[17] Kuang-Huei Lee, Xiaodong He, Lei Zhang, and Linjun Yang. Cleannet: Transfer learning for scalable image classifier training with label noise. In *CVPR*, pages 5447–5456. IEEE, 2018.

[18] Junnan Li, Richard Socher, and Steven CH Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. In *ICLR*, 2020.

[19] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S Kankanhalli. Learning to learn from noisy labeled data. In *CVPR*, pages 5051–5059. IEEE, 2019.

[20] Wen Li, Limin Wang, Wei Li, Eirikur Agustsson, and Luc Van Gool. Webvision database: Visual learning and understanding from web data. *arXiv preprint arXiv:1708.02862*, 2017.

[21] Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-learning regularization prevents memorization of noisy labels. In *NeurIPS*, 2020.

[22] Xingjun Ma, Yisen Wang, Michael E Houle, Shuo Zhou, Sarah Erfani, Shutao Xia, Sudanthi Wijewickrema, and James Bailey. Dimensionality-driven learning with noisy labels. In *ICML*, pages 3355–3364. PMLR, 2018.

[23] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, pages 181–196. Springer, 2018.

[24] Eran Malach and Shai Shalev-Shwartz. Decoupling" when to update" from" how to update". In *NeurIPS*, 2017.

[25] Aditya Menon, Brendan Van Rooyen, Cheng Soon Ong, and Bob Williamson. Learning from corrupted binary labels via class-probability estimation. In *ICML*, pages 125–134. PMLR, 2015.

[26] Nagarajan Natarajan, Inderjit S. Dhillon, Pradeep Ravikumar, and Ambuj Tewari. Learning with noisy labels. In *NeurIPS*, 2013.

[27] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep neural networks robust to label noise: A loss correction approach. In *CVPR*, pages 1944–1952. IEEE, 2017.

[28] Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping. In *ICLR*, 2014.

[29] Scott E Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping. In *ICLR (Workshop)*, 2015.

[30] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, pages 4278–4284, 2017.

[31] Cheng Tan, Jun Xia, Lirong Wu, and Stan Z Li. Co-learning: Learning from noisy labels with self-supervision. In *ACM MM*, pages 1405–1413. ACM, 2021.

[32] Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa. Joint optimization framework for learning with noisy labels. In *CVPR*, pages 5552–5560. IEEE, 2018.

[33] Ryutaro Tanno, Ardavan Saeedi, Swami Sankaranarayanan, Daniel C Alexander, and Nathan Silberman. Learning from noisy labels by regularized estimation of annotator confusion. In *CVPR*, pages 11244–11253. IEEE, 2019.

[34] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, 59(2):64–73, 2016.

[35] Xinshao Wang, Yang Hua, Elyor Kodirov, and Neil M Robertson. Imae for noise-robust learning: Mean absolute error does not treat examples equally and gradient magnitude's variance matters. *arXiv preprint arXiv:1903.12141*, 2019.

[36] Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross entropy for robust learning with noisy labels. In *ICCV*, pages 322–330, 2019.

[37] Hongxin Wei, Lei Feng, Xiangyu Chen, and Bo An. Combating noisy labels by agreement: A joint training method with co-regularization. In *CVPR*, pages 13726–13735. IEEE, 2020.

[38] Xiaobo Xia, Tongliang Liu, N. Wang, Bo Han, Chen Gong, Gang Niu, and Masashi Sugiyama. Are anchor points really indispensable in label-noise learning? In *NeurIPS*, 2019.

[39] Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. Learning from massive noisy labeled data for image classification. In *CVPR*, pages 2691–2699. IEEE, 2015.

[40] Yilun Xu, Peng Cao, Yuqing Kong, and Yizhou Wang. Ldmi: A novel information-theoretic loss function for training deep nets robust to label noise. In *NeurIPS*, pages 6222–6233, 2019.

[41] Yan Yan, Rómer Rosales, Glenn Fung, Ramanathan Subramanian, and Jennifer Dy. Learning from multiple annotators with varying expertise. *Machine learning*, 95(3):291–327, 2014.

[42] Yazhou Yao, Zeren Sun, Chuanyi Zhang, Fumin Shen, Qi Wu, Jian Zhang, and Zhenmin Tang. Jo-src: A contrastive approach for combating noisy labels. In *CVPR*, pages 5192–5201. IEEE, 2021.

[43] Kun Yi and Jianxin Wu. Probabilistic end-to-end noise correction for learning with noisy labels. In *CVPR*, pages 7017–7025. IEEE, 2019.

[44] Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor Tsang, and Masashi Sugiyama. How does disagreement help generalization against label corruption? In *ICML*, pages 7164–7173. PMLR, 2019.

[45] Xiyu Yu, Tongliang Liu, Mingming Gong, and Dacheng Tao. Learning with biased complementary labels. In *ECCV*, pages 68–83. Springer, 2018.

[46] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.

[47] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018.

[48] Zhilu Zhang and Mert R Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. In *NeurIPS*, 2018.

# Supplementary Materials

## 1. The design of realistic noisy

In the reality, the noisy labels are often introduced on ambiguous data. Thus, the artificial noise types are unlikely to happen in the real world. For instance, the symmetric noise randomly replaces the labels of a certain percentage of the training data by other labels. Although many existing methods can achieve a remarkable performance on such artificial noise types, they may not perform well on the real-world noisy datasets. Instead, we propose a new noise type on CIFAR datasets to mimic the real-world noisy datasets called realistic noise. It flips labels in class-pairs at a varying ratio according to the similarity between two classes. The detailed procedure to generate the noise transition matrix of realistic noise is shown below.

**Step 1.** We firstly train a ResNet50 on CIFAR-10/CIFAR-100 until it converges. Tab. 1 shows the test accuracy of the trained ResNet50 on both CIFAR-100 and CIFAR-10 datasets. We then consider the weight vector between the last fully connected layer and a node (a class) in the output layer as the prototype of the class [1].

**Step 2.** We pair all prototypes into $N$ class-pairs $(c_i, c_j)$, $i \neq j$, then compute their cosine similarities and sort them in descending order. We select the top $K$ pairs, and flip some labels in each pair, where $K$ is set to 10 for CIFAR-10 (Tab. 2) and 60 for CIFAR-100 (Tab. 3).

**Step 3.** The $K$ pairs is partitioned into three similarity levels. Each level has a weight $w$, which is proportional to the similarity of the class-pair. Higher $w$ is, more labels will be flipped. $w$ is selected from $\{0.9, 0.6, 0.3\}$ for CIFAR-100, and $\{0.9, 0.8, 0.7\}$ for CIFAR-10. As more similar pairs are selected from CIFAR-100, the weight step between different levels become larger.

**Step 4.** To construct the noise transition matrix, we temporally fill in the positions of the top $K$ class-pairs with the corresponding weights $w^{i,j}$, where $i$ and $j$ denote the $ith$ row and $jth$ column in the matrix. Other positions are filled in with zero. In order to let the sum of all weights in each row to be 1, the $w^{i,j}$ is replaced by

| Dataset | Acc.(%) |
|---------|---------|
| CIFAR-10 | 94.67 |
| CIFAR-100 | 78.85 |

Table 1. Accuracy of ResNet50 on CIFAR-10/CIFAR-100.

$$\hat{w}^{i,j} = \frac{w^{i,j}}{\sum_{j=1}^{n} w^{i,j}}, \tag{1}$$

where $n$ is the length of columns.

**Step 5.** Given a noise ratio $r$, we fill in the diagonal of the matrix with $(1 - r)$ and multiply the other elements by $r$.

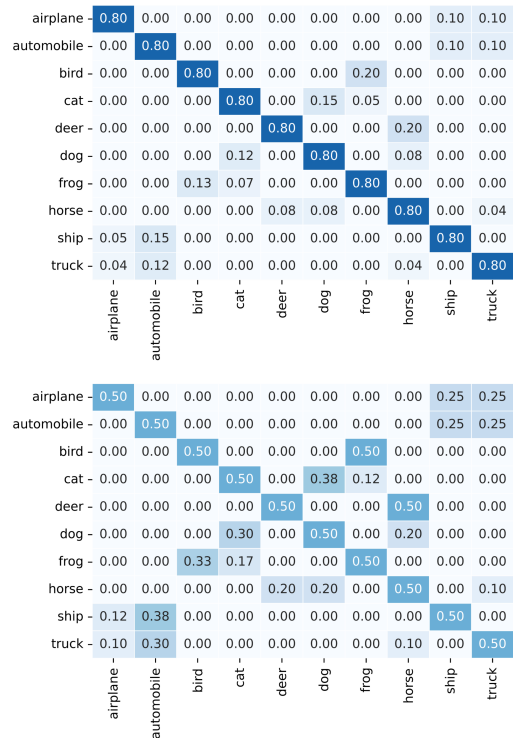Fig. 1 shows the generated noise transition matrices of CIFAR-10 with noise ratios 20% and 50%.



Figure 1. The generated noise transition matrices of CIFAR-10 with noise ratios 20% and 50%.

Fig. 2 visualizes the feature distributions of three similar class-pairs in CIFAR-100 using t-SNE. We train a ResNet50 for 200 epochs on CIFAR-100 without noisy labels, and use the output of the last convolutional layer as the feature of each sample. We can see that many hard data distribute close to the decision boundary between two similar classes.
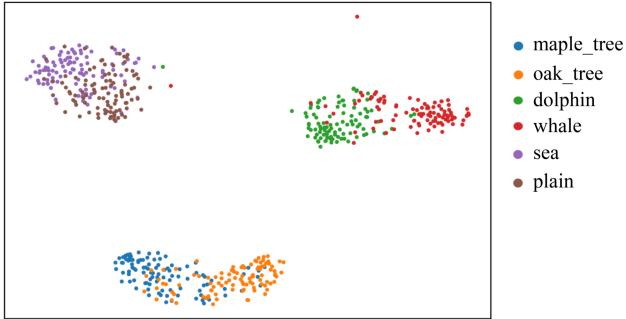
Figure 2. The visualization of the feature distributions of three similar class-pairs in CIFAR-100 using t-SNE [2].



Figure 3. Some samples selected from the top 3 similar class-pairs in CIFAR-100.

| Similar classes | | Similarity |
|---|---|---|
| automobile | truck | 0.075 |
| automobile | ship | 0.062 |
| cat | dog | 0.062 |
| dog | horse | 0.054 |
| deer | horse | 0.052 |
| bird | frog | 0.046 |
| airplane | ship | 0.044 |
| horse | truck | 0.025 |
| cat | frog | 0.021 |
| airplane | truck | 0.016 |

Table 2. The top 10 similar class-pairs in descending order from CIFAR-10. The number is the cosine similarity.

## 2. The motivation of Tripartition criterion

Figs. 5 to 9 show the distribution of normalized losses of all samples from CIFAR-100 and WebVision, Food-101N, Clothing1M datasets partitioned by our Tripartition criterion. We can see that the losses of samples in clean, hard and noisy subsets gradually increase. This observation is consistent with our hypothesis in the paper, i.e. losses of clean subset < losses of hard subset < losses of noisy subset. To filter out the hard subset, we need to determine a feasible range $[T_1, T_2]$.

One can see the loss distributions of three subsets gradually change with the progressing of training. It is nontrivial to determine a feasible range $[T_1, T_2]$ dynamically. Instead, our proposed Tripartition criterion wisely addresses this problem. It can dynamically partition the data and result in more accurate clean and noisy subsets.

In addition, Figs. 6 and 7 show that the loss distribution of realistic noise is more similar to the one of WebVision dataset compared to symmetric noise. Especially, the mean and variance of loss distributions of noisy subset and hard subset partitioned by Tripartition criterion are close to the ones of WebVision. This demonstrates that the proposed realistic noise is more consistent with the real-world noise type from the perspective of the loss distribution.

## 3. The changes of hard subset

Fig. 4 shows that the population of hard subset is becoming smaller with the progressing of training. The test is carried out 3 times on CIFAR-100 with symmetric and realistic noise types and 50% noise ratio.
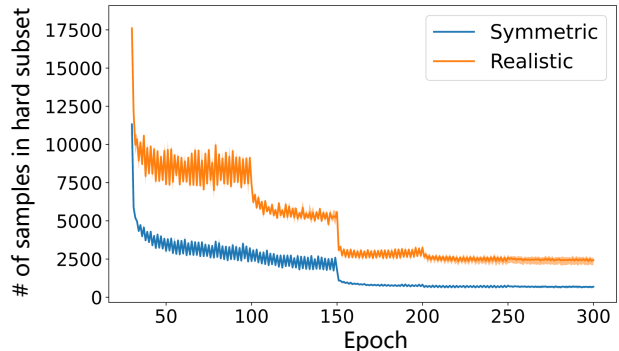


Figure 4. The changes of population of hard subset in training stage.

## 4. The effectiveness of Tripartition criterion.

To demonstrate the effectiveness of our proposed Tripartition criterion, we conduct a control experiment for Small-loss based and GMM based methods. Specifically, we replace the training data by the one partitioned by our Tripartition criterion for each method, and apply their own learning strategies to verify if their performances would be improved. We choose three versions of a network (PreActResNet-18) which was trained in the early, middle, and late training stages on CIFAR-100 with 50% realistic noise. These three versions achieve accuracies of 31.38%, 41.93%, and 63.50%, respectively. Taking the version trained in the middle stage (41.93%) as an example,

| Similar classes | | Similarity | Similar classes | | Similarity | Similar classes | | Similarity |
|---|---|---|---|---|---|---|---|---|
| maple_tree | oak_tree | 0.337 | poppy | tulip | 0.249 | beaver | otter | 0.222 |
| dolphin | whale | 0.312 | bed | couch | 0.246 | lion | tiger | 0.220 |
| plain | sea | 0.306 | poppy | rose | 0.245 | clock | plate | 0.220 |
| bus | pickup_truck | 0.301 | snake | worm | 0.245 | pickup_truck | tank | 0.219 |
| girl | woman | 0.298 | bottle | can | 0.242 | flatfish | ray | 0.217 |
| bus | streetcar | 0.296 | apple | pear | 0.242 | plain | road | 0.216 |
| bicycle | motorcycle | 0.275 | dolphin | shark | 0.237 | chimpanzee | woman | 0.216 |
| baby | girl | 0.273 | baby | boy | 0.232 | shark | trout | 0.215 |
| maple_tree | willow_tree | 0.270 | poppy | sunflower | 0.232 | lawn_mower | tractor | 0.214 |
| castle | house | 0.260 | bowl | plate | 0.231 | house | streetcar | 0.211 |
| oak_tree | pine_tree | 0.258 | cloud | plain | 0.231 | leopard | lion | 0.210 |
| cloud | sea | 0.257 | hamster | mouse | 0.230 | otter | seal | 0.210 |
| oak_tree | willow_tree | 0.255 | chair | couch | 0.230 | beetle | spider | 0.209 |
| orchid | tulip | 0.255 | bowl | cup | 0.229 | bed | chair | 0.209 |
| beetle | cockroach | 0.254 | rose | tulip | 0.229 | orange | pear | 0.209 |
| streetcar | train | 0.253 | apple | orange | 0.229 | maple_tree | rose | 0.207 |
| leopard | tiger | 0.253 | television | wardrobe | 0.227 | pickup_truck | tractor | 0.205 |
| man | woman | 0.251 | keyboard | telephone | 0.227 | forest | willow_tree | 0.205 |
| orange | sweet_pepper | 0.250 | boy | girl | 0.224 | crab | lobster | 0.205 |
| apple | sweet_pepper | 0.250 | mouse | shrew | 0.223 | mountain | sea | 0.204 |

Table 3. The top 60 similar class-pairs in descending order from CIFAR-100. The number is the cosine similarity.
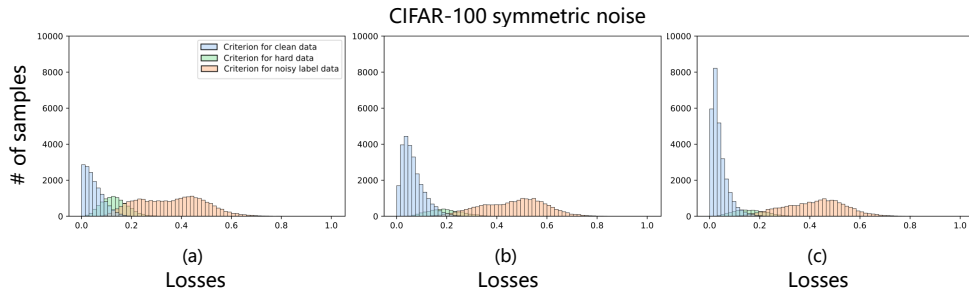


Figure 5. The distribution of normalized losses of all samples from CIFAR-100 partitioned by our Tripartition criterion. A PreActResNet-18 is trained with 30% symmetric noise. (a) The distribution at the end of the tenth epoch. (b) The distribution at the end of the fifieth epoch. (c) The distribution of epochs at the optimal accuracy.
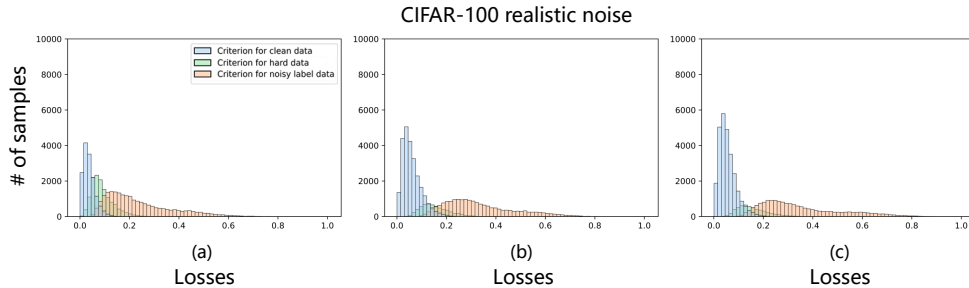


Figure 6. The distribution of normalized losses of all samples from CIFAR-100 partitioned by our Tripartition criterion. A PreActResNet-18 is trained with 30% realistic noise. (a) The distribution at the end of the tenth epoch. (b) The distribution at the end of the fifieth epoch. (c) The distribution of epochs at the optimal accuracy.
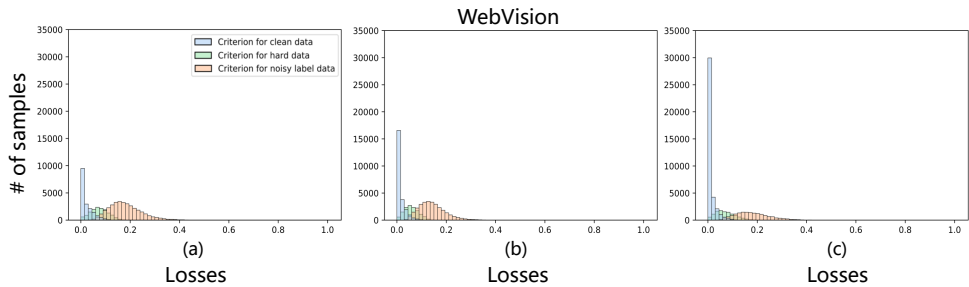


Figure 7. The distribution of normalized losses of all samples from WebVision partitioned by our Tripartition criterion. (a) The distribution at the end of the sceond epoch. (b) The distribution at the end of the fourth epoch. (c) The distribution of epochs at the optimal accuracy.
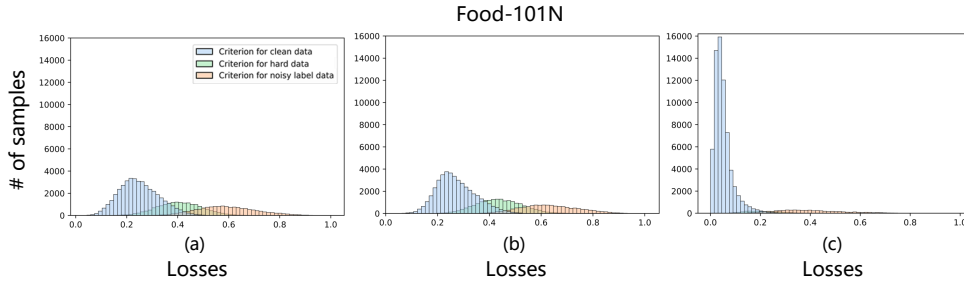
Figure 8. The distribution of normalized losses of all samples from Food-101N partitioned by our Tripartition criterion. The model is pre-trained on ImageNet. (a) The distribution at the end of the second epoch. (b) The distribution at the end of the fifth epoch. (c) The distribution of epochs at the optimal accuracy.
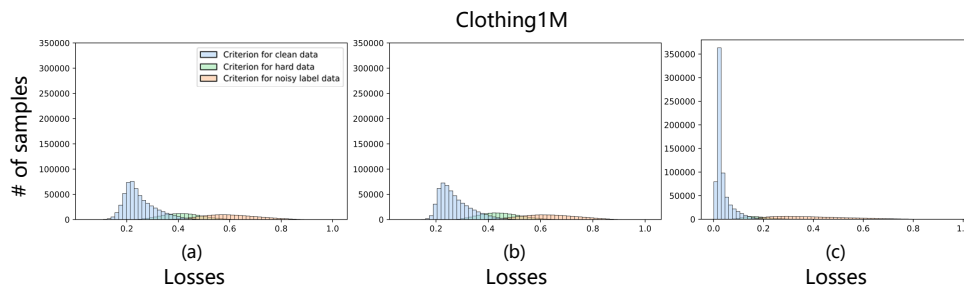


Figure 9. The distribution of normalized losses of all samples from Clothing1M partitioned by our Tripartition criterion. The model is pre-trained on ImageNet. (a) The distribution at the end of the second epoch. (b) The distribution at the end of the fifth epoch. (c) The distribution of epochs at the optimal accuracy.

the training data are divided by this version using three partition criteria, i.e. Small-loss (the partition criterion used by Co-teaching), GMM (the partition criterion used by DivideMix), Tripartition, respectively. We then have three different partitioned training datasets for the control experiment. For Co-teaching, it is further trained by the dataset divided by the Small-loss criterion. During the training, the dataset is not divided anymore. The control group is trained by the dataset divided by our Tripartition criterion. Two networks will converge after about 200 epochs. Their classification accuracies are listed in Tab. 4. One can see the network trained by dataset divided by our Tripartition criterion has 1.7% performance gain. For DivideMix, we follow the same protocol. The network is trained by the datasets divided by the GMM criterion and Tripartition criterion, respectively. The results show that the network trained by dataset divided by our Tripartition criterion achieves 1.04% performance gain.

We can observe similar results of the version trained in the early stage (31.38%) and the late stage (63.50%). Our Tripartition criterion improves Co-teaching and DivideMix about 0.72% and 0.23%, 0.6% and 0.7%, respectively. Above experiments demonstrate our Tripartition does reach

| Base_Model (Acc.%) | Method | Partition criterion | Acc.(%) |
|---|---|---|---|
| Early(31.38) | Co-teaching | Small-loss | 26.52 |
| | | Tripartition | **27.20** |
| | DivideMix | GMM | 33.85 |
| | | Tripartition | **34.08** |
| Middle(41.93) | Co-teaching | Small-loss | 32.70 |
| | | Tripartition | **34.40** |
| | DivideMix | GMM | 41.26 |
| | | Tripartition | **42.30** |
| Late(63.50) | Co-teaching | Small-loss | 37.92 |
| | | Tripartition | **38.54** |
| | DivideMix | GMM | 49.45 |
| | | Tripartition | **50.15** |

Table 4. The effectiveness of Tripartition criterion on other models when they are in different training stages.

a better training data partition and help a model to achieve a better performance.

## 5. The selection of $\lambda_n$ for noisy subset

The $\lambda_n$ is a parameter to control the contribution of noisy subset. Tab. 5 lists the selected values of $\lambda_n$ in varied noise settings on CIFAR datasets. The dataset with higher noise ratio and more complex data requires a larger $\lambda_n$ to learn

| Hyperparameter | Dataset | CIFAR-10 | | | | | CIFAR-100 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $\lambda_n$ | Realistic | 10% | 20% | 30% | 40% | 50% | 10% | 20% | 30% | 40% | 50% |
| | | 1 | 1 | 1 | 1 | 1 | 10 | 30 | 30 | 60 | 60 |
| | Symmetric | 20% | | 50% | | 80% | 20% | | 50% | | 80% |
| | | 1 | | 1 | | 1 | 10 | | 60 | | 120 |

Table 5. The selection of $\lambda_n$ in varied noise setting.

| Noise type | Symmetric | | Realistic | |
|---|---|---|---|---|
| Strategy | 20% | 50% | 20% | 50% |
| $\lambda_n = 1$ | 78.21 | 71.95 | 78.11 | 68.08 |
| $\lambda_n = 10$ | **78.65** | 72.95 | 78.21 | 68.30 |
| $\lambda_n = 40$ | 78.45 | 74.01 | **78.74** | 68.45 |
| $\lambda_n = 60$ | 78.31 | **74.73** | 78.24 | **68.62** |
| $\lambda_n = 100$ | 78.23 | 72.79 | 77.42 | 67.80 |

Table 6. Evaluation of $\lambda_n$ for training strategy of noisy subset in CIFAR-100.

| Datasets | # of training | # of test | # of class | size |
|---|---|---|---|---|
| CIFAR-10 | 50,000 | 10,000 | 10 | 32x32 |
| CIFAR-100 | 50,000 | 10,000 | 100 | 32x32 |
| Food-101N | 310,000 | 5,000 | 101 | 64x64 |
| Clothing1M | 1,000,000 | 10,526 | 14 | * |
| WebVision1.0 | 2,400,000 | 50,000 | 1000 | * |

Table 8. The detailed configurations of five benchmark datasets in the paper.

sufficient information. Tab. 6 shows more detailed results with varied $\lambda_n$. One can see that Tripartite is not sensitive to $\lambda_n$, particularly on low noise ratio and realistic noise. On symmetric noise with a ratio of 50%, the difference between the best and worst performances is less than 3%.

## 6. Additional experiment results

### 6.1. Results on CIFAR-10/CIFAR-100 with asymmetric noise

Tab. 7 shows the test accuracy on CIFAR-10 and CIFAR-100 with 40% asymmetric noise. (The noise ratio over 50% is not used because some classes become theoretically indistinguishable in such setting.)

## References

[1] Guy Bukchin, Eli Schwartz, Kate Saenko, Ori Shahar, Rogério Schmidt Feris, Raja Giryes, and Leonid Karlinsky. Fine-grained angular contrastive learning with coarse labels. In *CVPR*, pages 8726–8736, 2021. 1

[2] Laurens van der Maaten and Geoffrey E. Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9:2579–2605, 2008. 2

| Method(Backbone) \ Dataset | CIFAR-10 | CIFAR-100 |
|---|---|---|
| Cross entropy(ResNet34) | 80.11 | 42.74 |
| Bootstrap(ResNet34) | 81.21 | 45.12 |
| Forward(ResNet34) | 83.55 | 34.44 |
| GSE(ResNet34) | 76.74 | 47.22 |
| SL(ResNet34) | 82.51 | 69.32 |
| DivideMix(PreActResNet-18) | 93.40 | 72.10 |
| ELR+(PreActResNet-18) | 92.70 | 76.50 |
| Tripartite(PreActResNet-18) | **94.01** | **76.89** |

Table 7. Accuracy comparisons on CIFAR-10 and CIFAR-100 with 40% asymmetric noise.

### 6.2. Additional results of Tripartition criterion

Tab. 9 and Tab. 10 are the numerical versions of Fig. 6 in the paper.

## 7. Dataset details

The details of datasets used in the paper are shown in Tab. 8.

| Epoch | Acc.(%) | Clean_subset | | Noisy_subset | | Hard_subset |
|---|---|---|---|---|---|---|
| | | Partition acc. | Total nums | Partition acc. | Total nums | Total nums |
| 30 | 44.06 | 90.19 | 14873 | 83.65 | 23413 | 11714 |
| 50 | 46.49 | 97.79 | 17421 | 82.32 | 29543 | 3036 |
| 100 | 51.67 | 97.96 | 18901 | 85.32 | 28501 | 2598 |
| 150 | 53.71 | 97.75 | 19833 | 87.15 | 27765 | 2402 |
| 200 | 68.57 | 97.71 | 22509 | 90.83 | 26212 | 1279 |
| 250 | 73.23 | 97.84 | 22663 | 91.26 | 26030 | 1307 |
| 300 | 73.68 | 97.97 | 23132 | 92.87 | 25329 | 1539 |

Table 9. The results of Tripartition on CIFAR-100 with symmetric noise type and 50% noise ratio.

| Epoch | Acc.(%) | Clean subset | | Noisy subset | | Hard subset |
|---|---|---|---|---|---|---|
| | | Partition acc. | Total nums | Partition acc. | Total nums | Total nums |
| 30 | 39.20 | 68.48 | 17456 | 72.69 | 15485 | 17059 |
| 50 | 45.02 | 83.14 | 18801 | 76.11 | 25730 | 5469 |
| 100 | 46.89 | 83.78 | 20026 | 78.17 | 25390 | 4584 |
| 150 | 48.17 | 85.91 | 20361 | 80.20 | 25358 | 4281 |
| 200 | 67.10 | 81.76 | 24738 | 83.53 | 22896 | 2366 |
| 250 | 68.01 | 80.71 | 25646 | 84.07 | 22416 | 1938 |
| 300 | 67.76 | 80.21 | 25952 | 84.16 | 22288 | 1760 |

Table 10. The results of Tripartition on CIFAR-100 with realistic noise type and 50% noise ratio.