

Regularly Truncated M-estimators for Learning with Noisy Labels

Xiaobo Xia*, Pengqian Lu*, Chen Gong, *Senior Member, IEEE*, Bo Han, Jun Yu, *Senior Member, IEEE*, Jun Yu†, Tongliang Liu, *Senior Member, IEEE*

Abstract—The *sample selection* approach is very popular in learning with noisy labels. As deep networks “*learn pattern first*”, prior methods built on sample selection share a similar training procedure: the small-loss examples can be regarded as clean examples and used for helping generalization, while the large-loss examples are treated as mislabeled ones and excluded from network parameter updates. However, such a procedure is *arguably debatable* from two folds: (a) it does not consider the bad influence of noisy labels in selected small-loss examples; (b) it does not make good use of the discarded large-loss examples, which may be clean or have meaningful information for generalization. In this paper, we propose regularly truncated M-estimators (RTME) to address the above two issues *simultaneously*. Specifically, RTME can *alternately switch modes between truncated M-estimators and original M-estimators*. The former can *adaptively* select small-losses examples without knowing the noise rate and reduce the side-effects of noisy labels in them. The latter makes the possibly clean examples but with large losses involved to help generalization. Theoretically, we demonstrate that our strategies are label-noise-tolerant. Empirically, comprehensive experimental results show that our method can outperform multiple baselines and is robust to broad noise types and levels. The implementation is available at https://github.com/xiaoboxia/RTM_LNL.

Index Terms—learning with noisy labels, sample selection, truncated M-estimators, regularly truncated M-estimators, generalization

1 INTRODUCTION

LEARNING with noisy labels is one of the hottest problems in weakly supervised learning [1], [2], [3], [4], [5], since noisy labels are ubiquitous in real-world datasets, which always arise in mistakes of manual or automatic annotators [6], [7], [8], [9], [10], [11], [12], [13], [14]. Noisy labels can impair the performance of models, especially deep learning models (e.g., convolutional and recurrent neural networks) which have large model capacities. General regularization techniques such as dropout and weight decay cannot address this issue well [15]. Different approaches therefore have been proposed for robust learning with noisy labels [16], [17], [18], [19], [20], [21], [22], [23], [24]. Among them, the *sample selection* approach attracted a lot of attention from researchers, since it always has a simple mechanism but promising performance, and is orthogonal to other approaches [25], [26], [27]. This approach is also *our focus* in this paper.

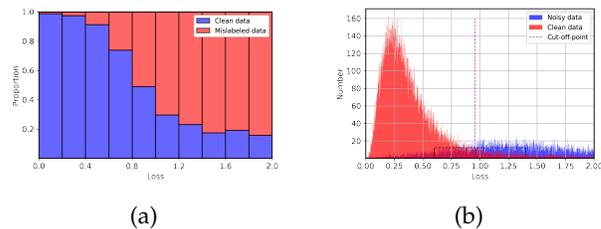


Fig. 1: Illustrations of training loss distributions. Experiments were conducted on synthetic *CIFAR-10* with instance-dependent label noise [17]. The noise rate is set to 30%. **(a):** Loss *vs* Proportion of clean/mislabeled data. Here, $proportion = (\# \text{ clean (resp. mislabeled data)}) / (\# \text{ all training data})$. The proportion of clean data is almost *negatively related* to the values of losses. **(b):** Loss *vs* Number. Noisy labels still exist in the selected small-loss examples, which hurt generalization.

The sample selection approach is based on selecting possibly clean examples from noisy examples for training. Intuitively, if we can exploit less noisy data for network parameter updates, the network will be more robust. At the present stage, the sample selection built on the small-loss criteria is the most common method, and has been verified to be effective in many circumstances [26], [28], [29], [30], [31]. Specifically, since deep networks *learn patterns first* [15], they would first memorize training data of clean labels and then those of noisy labels with the assumption that clean labels are of the majority in a noisy class. Small-loss examples can be regarded as clean examples *with high probability*. Therefore, in each iteration for a mini-batch data, the small-loss examples are selected for robust training *with equal importance*. By contrast, the large-loss examples are *treated to be mislabeled and excluded from training*.

* Equal contributions.

† Corresponding author.

- X. Xia and T. Liu are with the Sydney AI Center, School of Computer Science, Faculty of Engineering, The University of Sydney, Darlington, NSW2008, Australia (e-mail: xxia5420@uni.sydney.edu.au; tongliang.liu@sydney.edu.au).
- P. Lu is with the Australian AI Institute, Faculty of Engineering and IT, The University of Technology Sydney, Broadway, NSW, 2007, Australia (e-mail: pengqian.lu@student.uts.edu.au).
- C. Gong is with the PCA Lab, the Key Laboratory of Intelligent Perception and Systems for High-Dimensional Information of Ministry of Education, School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing, 210094, P.R. China; and is also with the Department of Computing, Hong Kong Polytechnic University, Hong Kong SAR, China (e-mail: chen.gong@njust.edu.cn).
- B. Han is with the Department of Computer Science, Hong Kong Baptist University, Hong Kong, China (email: bhanml@comp.hkbu.edu.hk).
- J. Yu is with the School of Computer Science and Technology, Hangzhou Dianzi University, Hangzhou, 310018, China (e-mail: yujun@hdu.edu.cn).
- J. Yu is with the Department of Automation, University of Science and Technology of China, Hefei, 230026, China (e-mail: harryjun@ustc.edu.cn).

However, such a selection procedure is *debatable* from *two folds*. First, the equal importance should not be assigned to different small-loss examples. Specifically, although we *rank* the losses of all examples and regard a proportion of examples as clean examples, such a way does not guarantee that selected examples are *completely clean* [29], especially the examples have *relatively large* losses but still are seen to be clean [32]. Recall the selection procedure, the purity of an example is *negatively correlated* with its loss, i.e., the example with a smaller loss is regarded to be clean with a higher degree of confidence (Fig. 1). Therefore, we should assign larger weights to the examples with smaller losses to make use of these “confirmable” clean examples to help generalization. Second, the large-loss examples should not be discarded directly. Specifically, although the large-loss examples may be mislabeled, the instances (e.g., images) may be helpful for generalization [33]. This opinion is motivated by the prior work [30], which shows that the images of mislabeled data may have meaningful information (e.g., scene information), even though such images come from a *different instance space*. For our task, mislabeled examples and clean examples share the same instance space. Such mislabeled examples is thereby more reasonable to be considered useful, and can be exploited for training.

In this paper, to relieve the above two issues simultaneously while keeping *end-to-end*, we propose *regularly truncated M-estimators*. More specifically, we borrow the statistical robust M-estimators in statistical learning [34], which can adaptively assign larger weights to examples with smaller losses. Based on the multiple robust M-estimators, to perform sample selection, we develop novel truncated M-estimators. By performing truncation on magnitudes of losses meanwhile *without knowing or estimating the noise rate*, our truncated M-estimators can concern the purity of small-loss examples and assign *zero weights* to possibly mislabeled examples to enhance networks. Since truncated M-estimators only consider the better use of small-loss examples, but do not make use of meaningful large-loss examples, we *regularly* switch robust M-estimators between truncated ones and original ones to achieve the proposed regularly truncated M-estimators. In this way, we can assign different weights to “clean” examples after sample selection (with truncated M-estimators). Additionally, the large-loss examples can be introduced regularly into network parameter updates for helping generalization (with original robust M-estimators). As large-loss examples are not introduced into training all the time, but are introduced regularly, and have *smaller* weights compared with small-loss examples, the side effect of possibly mislabeled examples can be reduced effectively, following better generalization.

Before delving into details, we highlight the main contributions of this paper in three folds:

- We show that the most frequently used sample selection procedure still has some potential weaknesses and discuss them carefully. Based on this, novel regularly truncated M-estimators are proposed to address the mentioned issues.
- Theoretical analysis is presented to demonstrate that the proposed methods are label-noise-tolerant. We also discuss that this work actually provides a new

and interesting perspective to make one loss function robust to label noise using the truncation mechanism.

- Extensive experiments on datasets with synthetic label noise and real-world label noise are conducted to verify the effectiveness of the proposed methods. Experimental results justify our claims well. Codes are open-source for future research.

1.1 Previous work

In this subsection, we briefly review prior approaches to learning with noisy labels, including robust loss functions, loss correction, and label correction. Our focus, i.e., the sample selection approach, will be introduced in detail later.

Robust loss functions. Some efforts have been made to design robust loss functions to handle noisy labels, e.g., the generalized cross-entropy loss [35], the normalized loss [36], the curriculum loss [37], the symmetric (cross-entropy) loss [38], [39], the negative loss [40], the peer loss [41], and the mutual information loss [42], etc.

Loss correction. This approach improves the robustness of networks by modifying the training loss. The modification can be achieved by reweighting losses [43], [44], estimating the noise transition matrix [17], [45], [46], [47], [48], [49], [50], [51], and adding an adaption layer [52], etc.

Label correction. The label correction approach [53], [54] aims to correct wrong labels into correct ones. The correction can be obtained by using directed graphical models [55], conditional random fields [56], knowledge graphs [57], and joint optimization methods [58], etc.

Integrated approach. Nowadays, state-of-the-art methods of handling noisy labels [59], [60], [61] are often designed by integrating various techniques at the same time. For example, they can simultaneously involve Mixup [62], soft labels [63], and semi-supervised learning [64], or involve sample selection and self-supervised learning [65]. We suggest that the readers refer to [66], [67] for more details about learning with noisy labels.

Compared with these prior effects, this paper offers an inspiring perspective to handle noisy labels, i.e., regularly truncated M-estimators, which successfully connects the classical statistical M-estimators and learning with noisy labels. Conceptually, this connection is new and valuable, and contributes to the research field.

1.2 Organization

The rest of this paper is organized as follows. In Section 2, we introduce the problem setting and some background of the proposed methods. In Section 3, we present the proposed methods formally. Experimental results are discussed in Section 4. The conclusion is given in Section 5.

2 PRELIMINARIES

In this section, we first introduce the notations (Section 2.1) and problem setting (Section 2.2). Then the sample selection approach for learning with noisy labels is discussed in detail (Section 2.3). Finally, we provide a brief introduction for the M-estimator (Section 2.4) and employed M-estimators in this work (Section 2.5).

2.1 Notations

Vectors and matrices are denoted by bold-faced letters. We use $\|\cdot\|_p$ as the ℓ_p norm of vectors or matrices. Let $[z] = \{1, 2, \dots, z\}$. For a function g , we use ∇g to denote its gradient. Let $\mathbb{1}_{\{\cdot\}}$ be the indicator function and “mod” be the math operation of taking the remainder.

2.2 Problem setup

Let \mathcal{X} and \mathcal{Y} be the instance and label space respectively. We consider a k -class classification problem, i.e., $\mathcal{Y} = [k]$. Let (\mathbf{x}, y) be the random variable pair of interest, and $p(\mathbf{x}, y)$ be the underlying joint density from which test data will be sampled. In *learning with noisy labels*, the labels of training data are corrupted. The training data are sampled from a *corrupted joint density* $p(\mathbf{x}, \tilde{y})$ rather than $p(\mathbf{x}, y)$, where \tilde{y} denotes the random variable of the noisy labels. Here, $p(\mathbf{x})$ remains the same, but $p(y|\mathbf{x})$ is corrupted into $p(\tilde{y}|\mathbf{x})$ [33], [68]. Therefore, we have an observed noisy training sample as follows:

$$S = \{(\mathbf{x}_i, \tilde{y}_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, \tilde{y}) = p(\tilde{y}|\mathbf{x})p(\mathbf{x}), \quad (1)$$

where n denotes the sample size of training data.

Let $f : \mathcal{X} \rightarrow \mathbb{R}^k$ be a classifier with parameters \mathbf{w} . Let $\ell : \mathbb{R}^k \rightarrow \mathbb{R}$ be a surrogate loss function for k -class classification. In this paper, we use the *softmax cross entropy loss* (abbreviated as the CE loss) [69]. Given an arbitrary training example $(\mathbf{x}_i, \tilde{y}_i)$, with parameters \mathbf{w} , we can obtain its CE loss:

$$L_i = \ell(f(\mathbf{w}; \mathbf{x}_i), \tilde{y}_i). \quad (2)$$

2.3 Sample selection for handling noisy labels

Prior effects exploited the sample selection approach to handle noisy labels [26], [28], [29], [33], [70], [71], which only used the “clean” examples (with relatively small losses) from each mini-batch for training. These clean examples have the same weights to contribute to optimization. Such methods employ the memorization effects of deep networks [15], which show that they would first memorize training data with clean labels and then those with noisy labels. We use a self-teach version of MentorNet [28] to give a better understanding for readers. The main procedure is shown in Algorithm 1.

Let us look at this procedure more closely. When a mini-batch data are formed (Step 5), we start to select possibly clean examples. In **Step 6**, we select a proportion of small-loss examples (controlled by the function $R(T)$) based on the network predictions. The large-loss examples are *abandoned directly* from optimization. In **Step 7**, the selected small-losses examples in the previous step are exploited for parameter updates. Their importance is seen to be the *same* for generalization. In Step 9, we update $R(T)$. Note that the function $R(T)$ needs to be designed carefully to better use the memorization effects of deep networks, and always is task-dependent [25]. For instance, in [26], [29], [31], $R(T) = 1 - \min\{T/T_k * \tau, \tau\}$, where τ is the noise rate. In practice, we cannot know the noise rate and have to estimate it [32]. Unfortunately, in some cases, e.g., the label noise is *instance-dependent*, the noise rate is hard to be estimated accurately [17], [72]. Accordingly, the effect of the

Algorithm 1 The main procedure of self-teach MentorNet for combating noisy labels.

-
- 1: **Input:** initialized classifier f , epoch T_k and T_{\max} , iteration t_{\max} .
 - 2: **for** $T = 0, \dots, T_{\max} - 1$ **do**
 - 3: **Shuffle** training dataset S ;
 - 4: **for** $t = 0, \dots, t_{\max} - 1$ **do**
 - 5: **Draw** a mini-batch \bar{S} from S ;
 - 6: **Select** $R(T)$ small-loss examples \bar{S}_f from \bar{S} based on classifier’s predictions;
 - 7: **Update** classifier parameters only using \bar{S}_f ;
 - 8: **end for**
 - 9: **Update** $R(T)$ with T_k ;
 - 10: **end for**
 - 11: **Output:** trained classifier f .
-

sample selection process will be influenced, which is never our desideratum.

As mentioned above, it is argued that the sample selection procedure (Algorithm 1) does not take care of the mislabeled data in the selected one and does not make use of large-loss data. Our methods tackle the two issues directly and are more advanced in that (1) the side-effect of mislabeled data belonging to selected data is reduced; (2) the meaningful formation of large-loss examples can be employed to help generalization. The technical implementation of our methods will be carefully discussed later.

2.4 The M-estimator

In statistics, M-estimators are a broad class of extremum estimators for which the objective function is a sample average [73]. We use a classical example (i.e., the estimation of the geometric median) to give an explanation for the M-estimator (cf. [34], [74]). For a dataset $\mathcal{A} = \{\mathbf{a}_i\}_{i=1}^N \subset \mathbb{R}^d$, the geometric median is the minimizer of the following function of $\mathbf{b} \in \mathbb{R}^d$:

$$\sum_{i=1}^N \|\mathbf{b} - \mathbf{a}_i\|_2. \quad (3)$$

This is a typical example of an M-estimator, that is a minimizer of a function of the form $\sum_{i=1}^N \rho(r_i)$, where r_i is a residual of the i -th data point, from the parametrized object (3). We have $r_i = \|\mathbf{b} - \mathbf{a}_i\|_2$ and $\rho(r_i) = r_i$. If there are some outliers in \mathcal{A} , the residuals of some data points may be unusually large and cause the minimizer to be unable to be learned accurately. Therefore, we need to give smaller weights to such data points to make results more robust, i.e., using robust M-estimators.

Below, we give a formal definition of the M-estimator in the context of learning with noisy labels.

Definition 1 (M-estimator). *In learning with noisy labels, an estimator is called the M-estimator, if it is an extremum estimator and can improve the robustness of the model by mitigating the side effect of mislabeled data during empirical risk minimization.*

For our task, we use the CE loss to measure the difference between predictions and given labels, the minimizer is the classifier f . During empirical risk minimization, the loss of the i -th data point is L_i accordingly. The data point with

M-estimators	$\Phi(L)$	$\nabla\Phi(L)$
CE [69]	L	∇L
Catoni's [75]	$\log(1 + L + L^2/2)$	$\frac{1+L}{1+L+L^2/2} \nabla L$
Log-sum Penalty [76]	$\log(1 + L/\epsilon)$	$\frac{\epsilon}{\epsilon+L} \nabla L$
Welsch+ [77]	$1 - \exp\{-L/\alpha^2\}$	$\frac{1}{\alpha^2} \exp\{-L/\alpha^2\} \nabla L$

TABLE 1: The definitions of used robust M-estimators.

an extremum of the loss is likely to be mislabeled. Its bad impacts on model robustness should be handled with the M-estimator which is discussed later.

It is worth noting that, for technical implementation, the M-estimators share a similar idea with prior robust loss functions in tackling noisy labels, i.e., making the contributions of mislabeled data into optimization smaller (but not zero) for robustness enhancement. The difference between the M-estimators and robust loss functions is that the original M-estimators perform a subsequent reweighting process based on the magnitude of the loss, while robust loss functions output the loss in robust training directly.

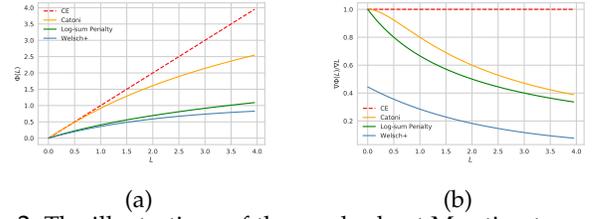
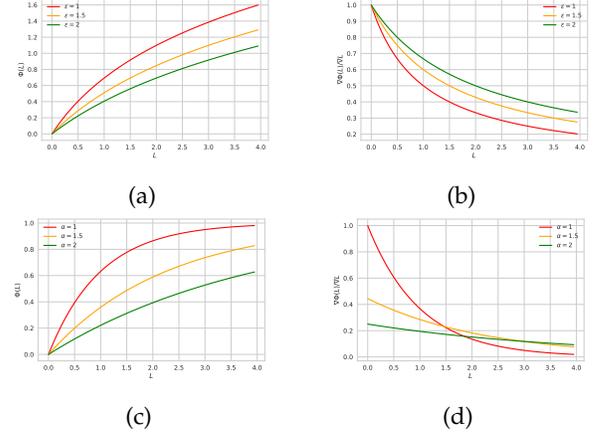
2.5 Representative M-estimators

The robustness of statistical M-estimators has been carefully studied for several decades [34]. One mainstream is to assign smaller weights to the data points with larger residuals [75] to make estimation results more robust. The reason for this is straightforward: the data points with large residuals are more likely to be *outliers*. If we reduce their contributions to the optimization of the objective function, the results will be less influenced by outliers, and naturally will be more robust. We borrow some representative examples of robust M-estimators in this paper, which will be introduced as follows. The robust M-estimators used are denoted by $\Phi(\cdot)$. To make the description clearer, we will directly use the notations in learning with noisy labels, i.e., $\Phi(L)$.

We compare assigned weights by robust M-estimators from an optimization viewpoint. That is to say, we compare the contributions to gradients brought by different examples, i.e., $\nabla\Phi(L)$. We exploit three robust M-estimators, i.e., Catoni's [75], Log-sum Penalty [76], and Welsch [77]. For Welsch, we change L^2 to L for weights assignments. The modified version is named Welsch+. The details of robust M-estimators used in this paper are provided in Table 1.

Note that $\epsilon \in [1, +\infty)$ and $\alpha \in (0, +\infty)$ are parameters of Log-sum Penalty and Welsch+ respectively. For a better understanding of used robust M-estimators, we provide illustrations for $\Phi(L)$ and $\nabla\Phi(L)$, which are shown in Fig. 2. From the illustrations, we can see that robust M-estimators can change the behaviors of losses integrally. When the loss of an example is large, the example may be mislabeled. Robust M-estimators can reduce its loss value and its contribution to optimization during training.

Besides, the curves of Log-sum Penalty and Welsch+ with different parameters are plotted in Fig. 3. As can be seen, different parameters can control different penalties for large-loss examples. The choices of parameters of the estimators Log-sum Penalty and Welsch+, i.e., ϵ and α , will be discussed in more detail later.


 Fig. 2: The illustrations of the used robust M-estimators, with $\epsilon = 2$ and $\alpha = 1.5$ for Log-sum Penalty and Welsch+. (a): L vs $\Phi(L)$. (b): L vs $\Phi(L)/\nabla L$.

 Fig. 3: The illustrations of Log-sum Penalty and Welsch+ with different parameters, i.e., ϵ and α . (a): Log-sum Penalty L vs $\Phi(L)$. (b): Log-sum Penalty L vs $\Phi(L)/\nabla L$. (c): Welsch+ L vs $\Phi(L)$. (d): Welsch+ L vs $\Phi(L)/\nabla L$.

3 METHODOLOGY

In this section, we formally present the proposed methods. We first propose how to perform truncation on the loss distribution automatically and achieve regularly truncated M-estimators (Section 3.1). Afterward, the analyses of parameters of regularly truncated M-estimators are presented (Section 3.2).

3.1 The proposed algorithms

We have discussed the mechanism of robust M-estimators. Nevertheless, we have two aspects that need to be considered carefully:

- How to reduce the side effects of noisy labels in selected small-loss examples?
- How to make good use of large-losses examples to help generalization?

The first question can be answered immediately by using robust M-estimators on selected small-loss examples. For the second question, we need to think prudently. Specifically, large-losses examples may be clean as discussed. Moreover, even they may be mislabeled, their instances (e.g., images) still may be helpful [30]. However, due to the harmful influence of incorrect labels, large-loss examples should be used *conservatively*. In this section, we formally present the proposed regularly truncated M-estimators to handle the mentioned problems at the same time.

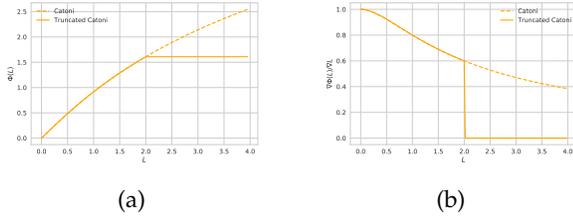


Fig. 4: Truncated Catoni vs Catoni. The truncation is performed at $\sigma = 2$.

3.1.1 Truncated M-estimators

To handle the first problem, i.e., using robust M-estimators on selected small-loss examples, we propose *truncated M-estimators*. Namely, we perform truncation on the loss distribution. The truncation divides all examples into small-loss ones and large-loss ones. We then can employ robust M-estimators to reduce the side effects of noisy labels in selected small-loss examples. By using the M-estimators Catoni’s, the truncated M-estimators are defined as follows:

$$\Phi^T(L) = \begin{cases} \log(1 + L + L^2/2) & L \leq \sigma \\ \log(1 + \sigma + \sigma^2/2) & L > \sigma \end{cases} \quad (4)$$

where $\sigma > 0$ is a hyperparameter related to the loss distribution to control the truncated point (or threshold). Other truncated robust M-estimators are provided in Table 2. The comparison between truncated Catoni’s and Catoni’s is provided in Fig. 4. As can be seen, truncated Catoni’s reserves the nice properties so that it can assign different weights to selected small-loss examples to relieve the influence of noisy labels (Fig. 4a). Meanwhile, in fact, it directly removes large-loss examples from training since such examples have no contribution to optimization (Fig. 4b).

Theoretical properties. We discuss the theoretical properties of the proposed truncated M-estimators. We demonstrate that they are noise-tolerant. That is, the minimizers of the risk under the truncated M-estimators with noisy labels would be the same as those with noise-free labels.

Lemma 1. *In a multi-class classification problem, the truncated M-estimators are noise-tolerant under symmetric (or uniform) label noise, if $c_2 - c_1 - k\Delta(\psi, f) > 0$ and the noise rate $\eta < \frac{(1-k)\Delta(\psi, f)}{c_2 - c_1 - k\Delta(\psi, f)}$. Here c_1 and c_2 denote the lower and upper bounds of the sum of the losses obtained by predictions on all classes, and $\Delta(\psi, f)$ is the gap between the clean risk w.r.t. (ψ, f) and the minimum clean risk brought by the global minimizer about f .*

Note that inspired by [78], the theoretical analysis can be extended to simple non-uniform noise under some conditions. Due to the limited page of the main paper, more background knowledge and detailed proofs of Lemma 1 are provided in Appendix A.

Remark 1. *The philosophy of noise tolerance of the truncated M-estimators is similar to the noise tolerance of some other robust loss functions that make the value of the loss sum bounded. For example, [78] considers the value of loss sum to be a constant C . Besides, [36] considers it to be 1. Differently, our truncated M-estimators employ a truncation mechanism to restrain the value of the loss sum, since the loss that is larger than σ is limited to be fixed, e.g., $\log(1 + \sigma + \sigma^2/2)$ for truncated Catoni’s. Based on*

this, the paper provides a new perspective to make one loss function robust to noisy labels.

3.1.2 Regularly truncated M-estimators.

As truncated M-estimators cannot handle the second problem, to handle the first and second problem at the same time, we further propose *regularly* truncated M-estimators. Here, the term “regularly” means that we alternately exploit truncated robust M-estimators and original robust M-estimators. Formally, we define

$$\Phi^R(L) = \mathbb{1}_{\{T \bmod R \neq 0\}} \Phi^T(L) + \mathbb{1}_{\{T \bmod R = 0\}} \Phi(L), \quad (5)$$

where $R \in \mathbb{N}_+$ is the hyperparameter about the frequency of using different kinds of robust M-estimators. Apparently, if the value of R is large, the large-loss examples will be involved in optimization *infrequently*, i.e., truncated M-estimators are often employed to perform reweighting on selected small-loss examples. Oppositely, if the value of R is small, large-loss examples will be involved in optimization *more frequently*. Distinctly, we need to choose a suitable R to achieve a great balance between truncated ones and original ones, which can be determined with a clean or noisy validation set. We will discuss this in Section 4.

3.2 Parameters analyses

For the truncated M-estimators, we have two types of parameters that need to be determined. The first type of parameter is the truncation point σ . The second type of parameter is the intrinsic parameter of M-estimators, i.e., ϵ and α . In this subsection, we discuss how to determine them.

We discuss how to determine σ . We borrow the “three-sigma” rule from the probability theory and statistics [79] rather than estimating the noise rate, since the noise rate is hard to be estimated in some cases [17]. The “three-sigma” rule has been verified to be effective to remove underlying outliers [80]. Specifically, let $\mathcal{L} \subset \mathbb{R}^n$ denote the losses of all training examples for each epoch. We first find the subset $\Gamma = \{0 \leq L_i \leq M | L_i \in \mathcal{L}\}$, where M represents the *median* of the set \mathcal{L} . Then we calculate the mean μ and standard deviation δ of the losses in Γ . Finally, we set the threshold $\sigma = \mu + 3\delta$. The threshold σ can be updated at every epoch according to the loss distribution. In Section 4, we will provide the experimental results for the justification of determining σ .

We then discuss how to determine ϵ and α during training. The parameter determination problem has been studied for several decades [81]. There are two main ways to determine this. On the one hand, we can empirically set the parameter in a reasonable range. On the other hand, we can assume the distribution of data to help determine. We follow both ways for the determination of ϵ and α . For the first way, we simply set $\epsilon = \alpha = 1$. For the second way, we assume that the outputs of M-estimators for selected small loss examples (denoted by Γ') obey a *Gaussian* distribution. More specifically, we calculate the mean μ' and standard deviation δ' . Then we tune ϵ (resp. α) to make that the distribution of Γ' is closer to $\mathcal{N}(\mu', \delta'^2)$.

The overall procedure of the proposed method is provided in Algorithm 2. As can be seen, in each epoch, we first determine the needed parameters (Step 4). Then when the mini-batch data is formed (Step 6), we use the proposed

Truncated M-estimators	$\Phi^T(L)$	$\nabla\Phi^T(L)$
Truncated Catoni's	$\begin{cases} \log(1 + L + L^2/2) & L \leq \sigma \\ \log(1 + \sigma + \sigma^2/2) & L > \sigma \end{cases}$	$\begin{cases} \frac{1+L}{1+L+L^2/2} \nabla L & L \leq \sigma \\ 0 & L > \sigma \end{cases}$
Truncated Log-sum Penalty	$\begin{cases} \log(1 + L/\epsilon) & L \leq \sigma \\ \log(1 + \sigma/\epsilon) & L > \sigma \end{cases}$	$\begin{cases} \frac{\epsilon}{\epsilon+L} \nabla L & L \leq \sigma \\ 0 & L > \sigma \end{cases}$
Truncated Welsch+	$\begin{cases} 1 - \exp\{-L/\alpha^2\} & L \leq \sigma \\ 1 - \exp\{-\sigma/\alpha^2\} & L > \sigma \end{cases}$	$\begin{cases} \frac{1}{\alpha^2} \exp\{-L/\alpha^2\} \nabla L & L \leq \sigma \\ 0 & L > \sigma \end{cases}$

TABLE 2: The definitions of proposed truncated M-estimators.

Algorithm 2 The procedure of regularly truncated M-estimators for learning with noisy labels.

- 1: **Input:** initialized classifier f , epoch T_{\max} , iteration t_{\max} , and the frequency R .
- 2: **for** $T = 0, \dots, T_{\max} - 1$ **do**
- 3: **Shuffle** the training dataset S ;
- 4: **Determine** the truncation parameter σ and intrinsic parameters ϵ/α as discussed in Section 3.2;
- 5: **for** $t = 0, \dots, t_{\max} - 1$ **do**
- 6: **Draw** a mini-batch \bar{S} from S ;
- 7: **Perform** regularly truncated M-estimators on \bar{S} with Eq. (5);
- 8: **Update** classifier parameters;
- 9: **end for**
- 10: **end for**
- 11: **Output:** trained classifier f .

regularly truncated M-estimators on it (Step 7). The proposed method is easy to follow and can keep an *end-to-end* manner.

4 EXPERIMENTS

In this section, we experimentally explore both the robustness and effectiveness of the proposed method. We first introduce the methods for comparison in the experiments (Section 4.1). We then introduce the details of the experiments on synthetic datasets (Section 4.2). The experiments on real-world datasets are finally presented (Section 4.3).

4.1 Comparison methods

We compare our method with multiple baselines, which include broad types of advanced approaches for combating noisy labels. The overview and publication locations of the baselines are summarized as follows.

- **Robust loss functions.** (1). APL (ICML 2020) [36], which combines two mutually reinforcing robust loss functions. (2). PCE (ICLR 2020) [82], which boosts the standard cross entropy loss with a partial trick. The tuning parameter of PCE is set to 2 in experiments. (3) AUL (T-PAMI 2023) [83], which are tailored to satisfy the Bayes-optimal condition and thus are robust to noisy labels under some conditions. (4) CELC (ICML 2023) [84], which induces a loss bound at the logit level, thus enhancing the noise robustness of the softmax cross entropy loss.
- **Loss correction.** (1). Revision (NeurIPS 2019) [85], which introduces a slack variable to revise the noise

transition matrix, leading to a better classifier. (2). Identifiability (ICML 2023) [86], which improves the estimation of the transition matrix using properly disentangled features.

- **Label correction.** (1). Joint (CVPR 2018) [58], which jointly optimizes the network parameters and the sample labels. The hyperparameters α and β for Joint are set to 1.2 and 0.8 respectively.
- **Sample selection.** (1). Co-teaching (NeurIPS 2018) [29], which trains two networks simultaneously and cross-updates parameters of peer networks. (2). SIGUA (ICML 2020) [33], which exploits stochastic integrated gradient underweighted ascent to handle noisy labels. We use self-teach SIGUA in this paper. (3). Co-Dis (ICCV 2023) [87], which selects possibly clean data that simultaneously have high-variance prediction probabilities between two networks. For these methods, we reserve their hyperparameter and optimization settings for selecting clean examples during training. Besides, we use an estimated noise rate [43] for them to ensure a fair comparison.

As a simple baseline, we compare our method with the standard deep network that directly trains on noisy datasets by using the softmax cross entropy loss function (abbreviated as CE). Note that we do not directly compare the proposed method with some state-of-the-art methods, e.g., SELF [88] and DivideMix [59]. It is because their proposed methods are aggregations of multiple techniques, e.g., Mixup [60], [62], soft labels [63], and semi-supervised learning [64]. We mainly focus on sample selection in learning with noisy labels. Therefore, the comparison is not fair. To make a fair comparison, we combine our method with semi-supervised learning and self-supervised learning to combat noisy labels. More details will be shown in Sections 4.2.8 and 4.2.9.

4.2 Experiments on simulated noisy datasets

4.2.1 Experimental setup

Datasets. We verify the effectiveness of our methods on the manually corrupted version of the following datasets: MNIST [89], SVHN [90], CIFAR-10 [91], CIFAR-100 [91], and NEWS [92], because these datasets are popularly used for the evaluation of learning with noisy labels in the literature [26], [29], [45], [93]. For NEWS, we borrowed the pre-trained word embeddings from GloVe [94]. The important statistics of the used synthetic datasets are summarized in Table 3.

Generating noisy labels. We consider two kinds of *class-dependent* label noise and one kind of *instance-dependent* label noise here. (1) Symmetric noise (abbreviated as Sym.) [68]:

Datasets	Type	# Of training	# Of testing	# Of class	Size
<i>MNIST</i>	image	60,000	10,000	10	$28 \times 28 \times 1$
<i>SVHN</i>	image	73,257	26,032	10	$32 \times 32 \times 3$
<i>CIFAR-10</i>	image	50,000	10,000	10	$32 \times 32 \times 3$
<i>CIFAR-100</i>	image	50,000	10,000	100	$32 \times 32 \times 3$
<i>NEWS</i>	text	11,314	7,532	20	300-D

TABLE 3: Summary of simulated noisy datasets used in the experiments.

this kind of label noise is generated by flipping labels in each class uniformly to incorrect labels of other classes. (2) Pairflip noise (abbreviated as Pair.) [26], [29]: the noise flips each class to its adjacent class. (3) Instance noise (abbreviated as Ins.) [72]: the noise is quite realistic, where the probability that an instance is mislabeled depends on its instances/features. We generate this type of label noise as did in [17]. For symmetric noise and instance noise, we set the noise rate τ to 30% and 50%. While, for pairflip noise, we set the noise rate τ to 30% and 45%, which aims to ensure that clean labels are diagonally dominant in noisy classes [29], [45], [70]. We leave out 10% of the noisy training data a validation set, which is used for model selection. Note that the correct labels are dominating in each noisy class and that label noise is random, the accuracy of the noisy validation set and the accuracy of the clean test data set are positively correlated. The noisy validation set therefore can be employed in experiments.

Network structure and optimizer. In terms of the five datasets with synthetic noise, for *MNIST*, we use a 3-layer MLP. Following [95], for *SVHN* and *CIFAR-10*, a ResNet-18 network is used. For *CIFAR-100*, a ResNet-50 network is used. Also, we employ a 3-layer MLP with the Softsign active function as did in [45]. We use SGD with momentum 0.9, weight decay 10^{-3} , batch size 128, and an initial learning rate 10^{-2} to train the networks. The learning rate is divided by 10 after the 40th epochs and 80th epochs. The maximum number of epochs is set to 200. For *SVHN*, *CIFAR-10*, and *CIFAR-100*, we perform data augmentation by horizontal random flips and 32×32 random crops after padding with 4 pixels on each side.

Measurement. As for performance measurement, we use test accuracy, i.e., $\text{test accuracy} = (\# \text{ of correct prediction}) / (\# \text{ of testing})$. All experiments are repeated five times. Intuitively, higher test accuracy means that the algorithm is more robust to noisy labels. We report the mean and standard deviation of the results. Besides, for fair comparison, we implement all methods with default parameters by PyTorch, and conduct all the experiments on NVIDIA Tesla V100 GPUs.

4.2.2 Analyses of experimental results

The results on *MNIST*, *SVHN*, *CIFAR-10*, *CIFAR-100*, and *NEWS* are presented in Table 4. For *MNIST*, as can be seen, the proposed methods achieve competitive classification performance. For *SVHN*, our methods outperform all baselines in all cases (one of the proposed methods works the best), which shows the effectiveness of our methods. For *CIFAR-10* and *CIFAR-100*, our methods also perform best. Lastly, for *NEWS*, our methods achieve the best results. Almost all the experimental results justify our claims well. Note that the performance on *NEWS* is a bit different from the results in [87]. This is because the optimization of the two works is

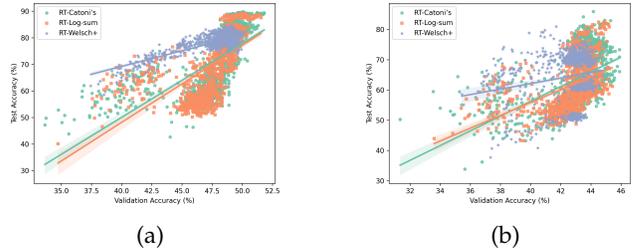


Fig. 5: Illustrations of positive correlation between validation accuracy and test accuracy. Straight lines are achieved by regression, which reflect the overall trend. The experiments are conducted on synthetic *CIFAR-10* with Pair.-45% noise (sub-figure (a)) and Ins.-50% noise (sub-figure (b)).

different. We use the SGD optimizer with momentum, while [87] uses the Adam optimizer.

4.2.3 Discussions of method selection

Note that different methods built on different M-estimators perform variably in different label noise cases. As discussed before, the differences between different M-estimators lie in the different ranges with respect to the training loss and different punishments on larger-loss examples. If we tend to choose the most suitable M-estimator, we need prior knowledge of data distributions, network architectures, and training dynamics (e.g., the loss distribution during optimization), which is rather hard or even impossible in practice. Fortunately, all methods built on M-estimators exhibit superior performance over baselines in most cases, which demonstrates the effectiveness of our M-estimator-based framework.

Moreover, here we propose to use the accuracy achieved on the noisy validation set for the selection of different methods. It is because the accuracy of the noisy validation set and the accuracy of the clean test data set are positively correlated. Therefore, for three methods, with the same training and validation data, we can choose the method that *overall* enjoys the higher validation accuracy. In Fig. 5, using the regression technology to mitigate the randomness and reflect the overall trend, we show the positive correlation between the accuracies of the noisy validation and test sets for method selection. Also, our RT-Catoni's enjoys both the highest validation accuracy and the highest test accuracy, which matches the results in Table 4.

4.2.4 The stability of our method

The stability about σ . As discussed in Section 3.2, our methods keep an adaptive manner to perform truncation and sample selection, i.e., using the “three-sigma” rule to determine a threshold. As an adaptive method, we do not need to estimate the noise rate. Prior works on sample selection show that if the noise rate cannot be estimated accurately, the classification performance will be affected largely [25]. Here, we show that our methods are stable even though the threshold is changed artificially during training.

The experiments are conducted on *MNIST* and *CIFAR-10* with 30% noise rates. Let Δ_σ be the disturbance added to the σ , where σ denotes the threshold determined by the algorithms. As shown in Fig. 6, the truncated M-estimators are

Datasets	Methods / Noise	Sym.-30%	Sym.-50%	Pair.-30%	Pair.-45%	Ins.-30%	Ins.-50%
MNIST	CE	96.29 ± 0.04	94.85 ± 0.12	95.15 ± 0.06	93.92 ± 0.39	95.87 ± 0.12	81.52 ± 5.52
	APL	96.22 ± 0.08	95.86 ± 0.26	96.28 ± 0.09	92.40 ± 0.69	90.07 ± 3.91	72.22 ± 15.36
	PCE	95.77 ± 0.62	95.07 ± 0.18	96.04 ± 0.17	93.92 ± 1.04	96.02 ± 0.47	78.93 ± 4.07
	AUL	94.07 ± 0.15	79.80 ± 3.64	60.42 ± 3.91	60.17 ± 3.69	92.16 ± 0.76	73.55 ± 7.17
	CELC	96.19 ± 0.11	95.35 ± 0.35	96.19 ± 0.98	95.84 ± 1.23	96.15 ± 1.38	89.15 ± 3.88
	Revision	96.47 ± 0.17	95.79 ± 0.24	96.08 ± 0.14	94.19 ± 0.93	96.49 ± 0.24	85.47 ± 3.04
	Identifiability	97.09 ± 0.35	95.86 ± 0.69	97.52 ± 0.47	97.48 ± 0.81	96.28 ± 0.56	88.44 ± 1.04
	Joint	96.26 ± 0.15	94.09 ± 0.47	94.02 ± 0.19	93.78 ± 0.92	96.03 ± 0.15	86.49 ± 4.15
	Co-teaching	96.04 ± 0.07	95.07 ± 0.24	96.09 ± 0.14	94.37 ± 0.58	94.53 ± 0.29	87.52 ± 2.44
	SIGUA	95.37 ± 0.93	95.07 ± 0.84	94.73 ± 0.39	90.04 ± 1.83	93.14 ± 1.29	80.47 ± 9.39
	Co-Dis	96.48 ± 0.15	95.37 ± 0.27	96.21 ± 0.14	94.20 ± 1.05	95.55 ± 1.03	90.33 ± 1.11
	RT-Catoni's	96.56 ± 0.04	96.33 ± 0.17	96.87 ± 0.12	95.72 ± 0.82	96.08 ± 0.18	93.25 ± 0.66
	RT-Log-sum	96.53 ± 0.12	96.21 ± 0.19	96.92 ± 0.09	95.70 ± 0.82	96.06 ± 0.14	93.77 ± 0.72
	RT-Welsch+	96.44 ± 0.06	95.90 ± 0.25	96.96 ± 0.04	96.56 ± 0.14	96.26 ± 0.05	93.02 ± 3.53
SVHN	CE	92.75 ± 0.31	90.63 ± 0.71	93.82 ± 0.13	70.95 ± 2.38	93.31 ± 0.37	63.16 ± 8.12
	APL	93.82 ± 0.19	91.34 ± 0.37	94.69 ± 0.19	86.77 ± 0.41	94.01 ± 0.36	67.61 ± 9.80
	PCE	93.81 ± 0.64	90.73 ± 0.19	94.24 ± 0.61	87.16 ± 1.14	93.31 ± 0.52	63.10 ± 7.26
	AUL	94.44 ± 0.52	92.75 ± 0.39	94.80 ± 1.20	82.77 ± 2.61	94.35 ± 0.16	64.33 ± 5.23
	CELC	95.06 ± 0.41	92.51 ± 0.89	94.36 ± 0.52	88.78 ± 1.37	94.16 ± 0.61	66.12 ± 3.44
	Revision	94.20 ± 0.22	94.06 ± 0.19	94.78 ± 0.30	81.36 ± 1.82	94.53 ± 0.57	67.21 ± 4.94
	Identifiability	93.18 ± 0.71	92.06 ± 1.33	92.66 ± 0.95	85.56 ± 1.40	92.01 ± 1.90	66.04 ± 5.71
	Joint	93.37 ± 0.27	92.11 ± 0.63	93.79 ± 0.29	75.86 ± 1.73	94.63 ± 0.82	62.19 ± 6.95
	Co-teaching	93.79 ± 0.67	92.63 ± 0.43	94.15 ± 0.62	88.36 ± 0.95	93.14 ± 0.12	61.55 ± 8.75
	SIGUA	94.04 ± 1.31	90.55 ± 2.44	92.19 ± 1.21	74.44 ± 5.72	92.66 ± 0.61	57.92 ± 11.68
	Co-Dis	94.77 ± 0.58	93.02 ± 0.82	94.78 ± 0.29	90.06 ± 1.03	93.77 ± 0.29	63.32 ± 8.80
	RT-Catoni's	95.54 ± 0.17	94.70 ± 0.20	95.29 ± 0.10	92.69 ± 0.83	94.69 ± 0.24	68.00 ± 13.15
	RT-Log-sum	95.51 ± 0.15	94.54 ± 0.21	95.25 ± 0.22	91.59 ± 2.17	94.92 ± 0.19	66.96 ± 12.57
	RT-Welsch+	95.44 ± 0.08	94.47 ± 0.16	95.16 ± 0.36	92.89 ± 0.71	94.99 ± 0.23	61.60 ± 15.63
CIFAR-10	CE	82.67 ± 0.48	76.01 ± 1.43	84.97 ± 1.04	61.76 ± 4.53	83.15 ± 0.55	54.29 ± 3.90
	APL	85.54 ± 0.51	78.36 ± 0.47	85.40 ± 0.14	80.84 ± 0.72	77.57 ± 0.15	39.45 ± 6.51
	PCE	86.12 ± 0.85	74.03 ± 4.96	85.03 ± 0.77	65.08 ± 3.41	85.64 ± 0.72	64.82 ± 4.13
	AUL	88.09 ± 0.78	82.81 ± 1.16	71.34 ± 1.91	56.80 ± 2.69	86.35 ± 0.90	60.75 ± 3.77
	CELC	89.46 ± 2.13	85.08 ± 3.95	89.77 ± 2.56	85.72 ± 4.52	86.67 ± 1.47	61.85 ± 4.98
	Revision	88.39 ± 0.38	83.40 ± 0.65	90.70 ± 0.47	83.61 ± 1.06	89.07 ± 0.35	66.93 ± 4.14
	Identifiability	87.12 ± 1.69	83.43 ± 2.11	86.45 ± 1.93	83.65 ± 2.46	80.47 ± 1.54	55.25 ± 3.78
	Joint	89.34 ± 0.52	85.06 ± 0.29	89.75 ± 0.63	80.52 ± 1.90	88.41 ± 1.02	64.12 ± 3.89
	Co-teaching	88.93 ± 0.56	85.73 ± 0.12	88.72 ± 0.61	84.19 ± 0.68	87.07 ± 0.35	60.09 ± 3.31
	SIGUA	83.19 ± 1.26	77.92 ± 3.11	83.93 ± 0.49	70.39 ± 1.94	82.90 ± 2.00	30.95 ± 9.70
	Co-Dis	89.20 ± 0.13	85.36 ± 0.94	89.20 ± 0.37	85.02 ± 1.33	87.13 ± 0.25	62.77 ± 3.90
	RT-Catoni's	89.39 ± 0.28	87.00 ± 0.08	90.83 ± 0.20	86.57 ± 0.92	89.34 ± 0.32	69.77 ± 2.14
	RT-Log-sum	89.60 ± 0.44	87.41 ± 0.30	90.49 ± 0.12	83.60 ± 1.38	89.65 ± 0.88	68.97 ± 3.82
	RT-Welsch+	90.65 ± 0.22	86.60 ± 0.51	90.15 ± 0.38	77.29 ± 6.52	89.56 ± 0.62	60.86 ± 10.60
CIFAR-100	CE	51.25 ± 0.50	40.28 ± 0.53	51.71 ± 0.63	38.54 ± 0.53	52.02 ± 0.44	36.35 ± 0.87
	APL	55.78 ± 0.91	46.96 ± 0.81	56.34 ± 0.68	49.55 ± 1.05	43.30 ± 1.57	29.01 ± 0.09
	PCE	58.84 ± 1.32	42.63 ± 2.02	54.23 ± 1.76	41.05 ± 2.83	55.72 ± 1.96	38.72 ± 3.01
	AUL	69.89 ± 0.21	60.00 ± 0.40	64.96 ± 0.55	39.37 ± 1.61	67.75 ± 1.84	40.27 ± 1.76
	CELC	67.96 ± 1.88	60.71 ± 2.39	67.96 ± 2.10	52.53 ± 3.17	66.25 ± 1.93	47.52 ± 3.93
	Revision	62.97 ± 0.46	43.60 ± 0.94	60.09 ± 1.21	49.33 ± 1.10	56.46 ± 1.45	40.78 ± 1.75
	Identifiability	50.53 ± 1.52	34.87 ± 2.36	52.88 ± 1.15	38.16 ± 2.68	52.48 ± 1.93	36.72 ± 3.10
	Joint	63.69 ± 0.84	55.62 ± 1.68	65.11 ± 1.79	49.77 ± 1.15	64.15 ± 1.11	45.47 ± 2.73
	Co-teaching	59.49 ± 0.36	52.19 ± 1.42	54.92 ± 2.84	47.53 ± 1.39	56.71 ± 1.26	42.09 ± 1.73
	SIGUA	54.22 ± 0.90	50.64 ± 3.92	47.92 ± 2.93	39.92 ± 2.33	53.19 ± 2.64	38.50 ± 1.69
	Co-Dis	64.02 ± 1.37	54.55 ± 2.06	58.72 ± 2.11	50.02 ± 2.80	59.15 ± 1.92	43.38 ± 1.25
	RT-Catoni's	70.04 ± 0.28	64.87 ± 0.52	71.75 ± 0.33	50.02 ± 0.95	71.66 ± 0.53	53.97 ± 0.45
	RT-Log-sum	70.30 ± 0.32	65.20 ± 0.44	71.68 ± 0.18	48.16 ± 1.26	71.22 ± 0.50	54.09 ± 0.37
	RT-Welsch+	69.17 ± 0.60	57.63 ± 0.92	69.34 ± 0.50	54.00 ± 1.50	69.22 ± 0.17	56.44 ± 1.78
NEWS	CE	43.16 ± 1.95	32.92 ± 0.86	42.86 ± 1.06	28.33 ± 3.58	44.08 ± 1.70	30.06 ± 7.92
	APL	54.04 ± 1.09	45.12 ± 2.17	51.98 ± 0.27	36.86 ± 2.31	52.18 ± 0.63	44.82 ± 3.61
	PCE	55.12 ± 0.94	49.77 ± 0.32	54.17 ± 0.98	37.92 ± 2.02	54.37 ± 0.95	46.14 ± 1.29
	AUL	53.77 ± 0.25	48.78 ± 1.62	53.72 ± 1.77	39.23 ± 1.06	55.19 ± 1.09	47.73 ± 2.11
	CELC	52.15 ± 0.86	47.25 ± 1.00	52.50 ± 0.84	38.10 ± 1.06	53.70 ± 1.81	47.00 ± 2.06
	Revision	55.19 ± 0.63	50.65 ± 0.97	53.77 ± 0.64	38.91 ± 1.38	53.29 ± 0.62	46.37 ± 2.94
	Identifiability	53.65 ± 1.65	50.84 ± 2.27	53.16 ± 1.95	39.16 ± 2.62	52.35 ± 1.92	44.87 ± 3.92
	Joint	53.15 ± 0.92	48.77 ± 1.47	51.90 ± 1.35	33.29 ± 3.45	52.92 ± 0.64	43.47 ± 2.94
	Co-teaching	53.81 ± 0.76	51.22 ± 0.61	53.90 ± 0.45	39.24 ± 1.19	53.99 ± 0.47	48.92 ± 2.04
	SIGUA	51.33 ± 1.41	47.47 ± 2.35	50.81 ± 2.19	32.12 ± 4.37	51.22 ± 2.61	30.82 ± 7.75
	Co-Dis	54.20 ± 0.39	51.97 ± 0.46	54.30 ± 0.15	41.04 ± 1.77	54.25 ± 0.25	49.03 ± 1.76
	RT-Catoni's	57.83 ± 0.45	53.16 ± 0.74	54.95 ± 0.85	44.25 ± 2.36	56.68 ± 0.58	48.85 ± 1.21
	RT-Log-sum	58.07 ± 0.32	53.30 ± 0.48	55.22 ± 0.31	44.21 ± 1.61	56.95 ± 0.75	49.01 ± 1.49
	RT-Welsch+	58.08 ± 0.67	54.22 ± 0.83	56.32 ± 0.27	42.75 ± 2.35	57.98 ± 0.57	50.13 ± 1.83

TABLE 4: Mean and standard deviations of test accuracy (%) on synthetic *MNIST*, *SVHN*, *CIFAR-10*, *CIFAR-100*, and *NEWS*. The best 3 experimental results are in bold while the best is underlined.

very sensitive to the values of thresholds. In particular, when $\Delta_\sigma = -20\%$, the classification performance of truncated M-estimators is greatly affected. On synthetic *CIFAR-10*, the

accuracies are reduced by almost 20%. As a comparison, the proposed regularly truncated M-estimators are very stable when the disturbance is added to the threshold. It is because

Methods	RT-Catoni's			RT-Log-sum			RT-Welsch+		
	Sym.-30%	Pair.-30%	Ins.-30%	Sym.-30%	Pair.-30%	Ins.-30%	Sym.-30%	Pair.-30%	Ins.-30%
Gaussian	89.36 ± 0.04	90.77 ± 0.27	89.25 ± 0.37	89.42 ± 0.10	90.65 ± 0.21	89.66 ± 0.17	90.38 ± 0.34	89.19 ± 0.09	88.77 ± 0.42
Fixed	89.05 ± 0.14	90.83 ± 0.20	89.07 ± 0.77	89.60 ± 0.44	90.49 ± 0.12	89.65 ± 0.88	90.65 ± 0.22	90.15 ± 0.38	89.56 ± 0.62

TABLE 5: Mean and standard deviations of test accuracy (%) with different parameter determination ways.

Methods	Catoni's Based			Log-sum Based			Welsch+ Based		
	Sym.-30%	Pair.-30%	Ins.-30%	Sym.-30%	Pair.-30%	Ins.-30%	Sym.-30%	Pair.-30%	Ins.-30%
Original	84.95 ± 0.59	85.37 ± 0.78	84.25 ± 0.66	85.31 ± 0.56	54.64 ± 0.58	84.19 ± 0.86	85.51 ± 0.40	71.91 ± 5.79	50.85 ± 18.60
T-CE	85.07 ± 0.31	86.32 ± 0.19	86.15 ± 1.48	85.77 ± 0.35	82.30 ± 1.95	86.33 ± 1.37	85.55 ± 0.89	84.12 ± 2.93	70.50 ± 4.05
T-M-estimators	86.43 ± 0.35	90.25 ± 0.22	88.37 ± 1.71	86.60 ± 0.41	88.44 ± 1.45	88.43 ± 1.72	86.74 ± 0.33	88.41 ± 1.48	88.50 ± 1.53
RT-M-estimators	89.39 ± 0.28	90.83 ± 0.20	89.34 ± 0.32	89.60 ± 0.44	90.49 ± 0.12	89.65 ± 0.88	90.65 ± 0.22	90.15 ± 0.38	89.56 ± 0.62

TABLE 6: Mean and standard deviations of test accuracy (%) with M-estimators (i.e., "Original"), truncated CE (abbreviated as "T-CE"), truncated M-estimators (abbreviated as "T-M-estimators"), and regularly truncated M-estimators (abbreviated as "RT-M-estimators"). The experiments are conducted on synthetic *CIFAR-10*.

we regularly introduce large-loss examples into training. The underlying clean examples can be exploited. Also, such a way can address the covariate shift issue effectively mentioned in [28], and therefore helps generalization.

The stability about ϵ and α . We present the sensitivity analyses on the intrinsic parameters of exploited M-estimators, i.e., ϵ in Log-sum Penalty and α in Welsch+. The experiments are conducted on *MNIST* and *CIFAR-10* with 30% noise rates. The range of ϵ and α is {1.5,2,2.5,3}. As can be seen in Tables 7 and 8, the M-estimators are robust to the choice of intrinsic parameters in a certain range, which implies that the proposed methods can be easily applied in practice.

4.2.5 Ablation study

We conduct detailed ablation studies to analyze and show the effects of different components to provide insights into what makes our methods successful.

The influence of R . We first analyze the effect of the frequency of using different kinds of robust M-estimators, i.e., R . The experimental results are shown in Fig. 7. As can be seen, with the increase of R , the test accuracies decrease clearly. In other words, the introduction of large-loss examples in a conservative way can improve the algorithm performance, which verifies the effectiveness of our methods.

Impact of each component. We then compare the results of M-estimators, truncated CE, truncated M-estimators, and the proposed regularly truncated M-estimators. The results are shown in Table 5. We can see that original M-estimators cannot work well when there are noisy labels, and truncated M-estimators can better handle noisy labels. Also, comparing the truncated CE with truncated M-estimators, we can see that assigning different weights on small-loss examples can improve performance. Additionally, the proposed regularly truncated M-estimators outperform truncated M-estimators, which shows the effectiveness of introducing large-loss examples into training. Note that compared with the results in Fig. 7 and Table 5, we can know that both the truncation and the introduction of large-loss examples are of importance against noisy labels. Besides, there is a trade-off switching frequency with a relatively small value.

Adaptive determination of ϵ and α . As discussed in Section 3.2, we can adaptively determine the ϵ and α by introducing a Gaussian distribution assumption. We use the

ℓ_2 distance to measure the distribution Γ' to $\mathcal{N}(\mu', \delta'^2)$ in this paper. The results are provided in Table 6. Accordingly, we can know that our methods are able to avoid tuning the hyperparameters ϵ and α artificially. Instead, they can be determined automatically by using the Gaussian distribution assumption. Also, our methods work well in such a way.

4.2.6 A closer look on the memorization effect

The memorization effect of the deep network [96] shows that it would first memorize clean data and then memorize mislabeled data. Therefore, in early training, the network is relatively robust with noisy labels, i.e., more memorization of clean data and less memorization of mislabeled data, following good test accuracy on clean test data. Here, we provide a closer look and show that exploited M-estimators (i.e., original M-estimators) can strengthen the memorization effect. The results in Fig. 8 show that used M-estimators make the deep network less memorize mislabeled data, leading to better test accuracy in early training. Interestingly, we find that the ways of enhancing model robustness of multiple M-estimators are slightly different, though all of them can tackle noisy labels successfully. In more detail, Catoni's and Log-sum can strengthen the memorization of clean data and reduce the memorization of mislabeled data at the same time. Differently, Welsch+ works well in largely reducing the memorization of mislabeled data.

4.2.7 Visualization of experimental results

We use 2D t-SNE [97] to visualize the experimental results which are presented in Fig. 9. We can see that the proposed methods work well, and can distinguish different classes clearly when there are noisy labels.

4.2.8 Combination with semi-supervised learning

Recall that we discussed the comparison between the proposed methods and some methods comprising multiple techniques is unfair. Therefore, to make it fair, here we boost our methods with semi-supervised learning. Specifically, we develop the framework of DivideMix [59]. Different from the original DivideMix which just uses the cross-entropy loss for follow-up sample selection and semi-supervised learning, we employ our regularly truncated M-estimators for warm-up. Results are provided in Table 9. As can be

Methods	Sym.-30%				Pair.-30%				Ins.-30%			
	$\epsilon = 1.5$	$\epsilon = 2$	$\epsilon = 2.5$	$\epsilon = 3$	$\epsilon = 1.5$	$\epsilon = 2$	$\epsilon = 2.5$	$\epsilon = 3$	$\epsilon = 1.5$	$\epsilon = 2$	$\epsilon = 2.5$	$\epsilon = 3$
RT-Log-sum	96.54 ± 0.35	96.47 ± 0.40	96.49 ± 0.40	96.52 ± 0.38	96.50 ± 0.23	96.40 ± 0.25	96.35 ± 0.24	96.31 ± 0.25	96.26 ± 0.32	96.22 ± 0.36	96.17 ± 0.32	96.14 ± 0.41
RT-Welsch+	96.55 ± 0.19	96.68 ± 0.33	96.71 ± 0.29	96.72 ± 0.33	96.52 ± 0.27	96.50 ± 0.24	96.32 ± 0.18	96.31 ± 0.30	96.42 ± 0.33	96.41 ± 0.33	96.35 ± 0.36	96.35 ± 0.42

TABLE 7: The sensitivity analyses on the intrinsic parameters of exploited M-estimators. The experiments are conducted on synthetic *MNIST* with 30% noise rates.

Methods	Sym.-30%				Pair.-30%				Ins.-30%			
	$\alpha = 1.5$	$\alpha = 2$	$\alpha = 2.5$	$\alpha = 3$	$\alpha = 1.5$	$\alpha = 2$	$\alpha = 2.5$	$\alpha = 3$	$\alpha = 1.5$	$\alpha = 2$	$\alpha = 2.5$	$\alpha = 3$
RT-Log-sum	89.46 ± 0.11	88.93 ± 0.19	88.94 ± 0.21	88.94 ± 0.09	90.73 ± 0.15	90.70 ± 0.15	90.66 ± 0.10	90.71 ± 0.21	89.35 ± 0.39	89.39 ± 0.49	88.86 ± 0.33	88.96 ± 0.49
RT-Welsch+	90.81 ± 0.15	90.64 ± 0.11	90.15 ± 0.14	89.77 ± 0.19	89.76 ± 0.19	90.46 ± 0.20	90.40 ± 0.06	90.61 ± 0.28	89.48 ± 0.41	89.44 ± 0.72	89.63 ± 0.36	89.53 ± 0.21

TABLE 8: The sensitivity analyses on the intrinsic parameters of exploited M-estimators. The experiments are conducted on synthetic *CIFAR-10* with 30% noise rates.

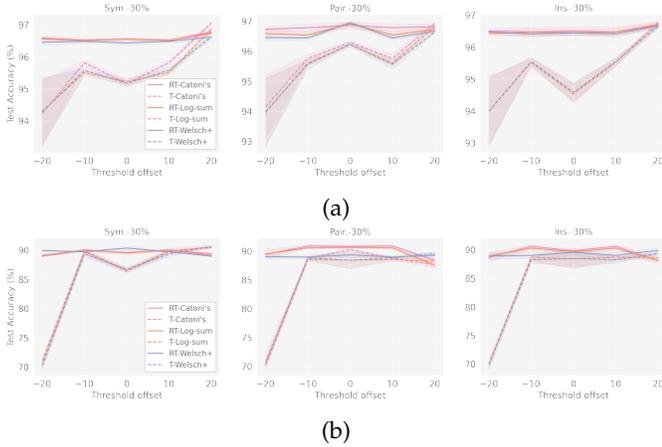


Fig. 6: Illustrations of the test accuracy with different disturbances. The experimental results reveal that regularly truncated M-estimators are more stable. The experiments are conducted on synthetic *MNIST* (subfigure (a)) and synthetic *CIFAR-10* (subfigure (b)) with 30% noise rates.

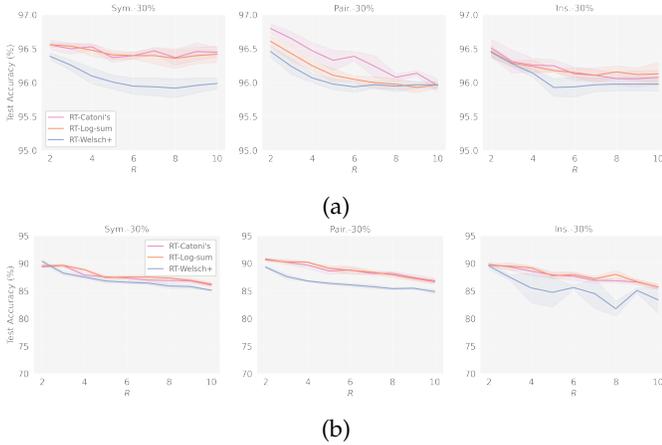


Fig. 7: Illustrations of the test accuracy with different values of R . These experiments reveal a smaller R , which means that introducing large-loss examples frequently can lead to better classification performance in general. The experiments are conducted on synthetic *MNIST* (subfigure (a)) and synthetic *CIFAR-10* (subfigure (b)) with 30% noise rates.

seen, in almost all cases, the proposed methods can bring performance improvements. Especially in the cases of Pair-

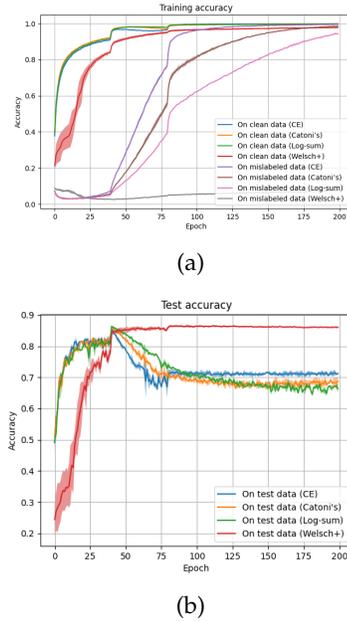


Fig. 8: Illustrations of the training and test accuracy achieved by different methods with the increase of epochs. Experiments are conducted on synthetic *CIFAR-10* with Sym.-30% noise. (a): Training accuracy vs Epoch. (b): Test accuracy vs Epoch.

Methods	Sym.-50%	Pair.-45%	Ins.-50%
DivideMix	95.00 ± 1.12	86.55 ± 2.74	<u>92.90 ± 2.26</u>
DivideMix+RT-Catoni's	<u>95.01 ± 1.01</u>	<u>94.88 ± 1.63</u>	<u>94.56 ± 2.84</u>
DivideMix+RT-Log-sum	<u>95.18 ± 0.84</u>	<u>94.84 ± 1.36</u>	<u>94.78 ± 1.29</u>
DivideMix+RT-Welsch	<u>95.02 ± 1.28</u>	<u>95.21 ± 2.77</u>	91.51 ± 2.29

TABLE 9: Mean and standard deviations of test accuracy (%) on *CIFAR-10* compared DivideMix with the methods boosted by the proposed algorithms. The best 3 experimental results are in bold while the best is underlined.

Methods	Sym.-50%	Pair.-45%	Ins.-50%
CL	82.56 ± 1.14	58.97 ± 1.52	56.51 ± 2.82
CL+RT-Catoni's	<u>87.75 ± 1.90</u>	<u>90.16 ± 2.64</u>	<u>75.45 ± 3.05</u>
CL+RT-Log-sum	<u>87.73 ± 1.59</u>	<u>88.96 ± 2.21</u>	<u>72.88 ± 3.15</u>
CL+RT-Welsch	<u>88.26 ± 1.41</u>	<u>80.99 ± 2.31</u>	<u>72.91 ± 4.43</u>

TABLE 10: Mean and standard deviations of test accuracy (%) on *CIFAR-10* compared contrastive learning (CL) with the methods boosted by the proposed algorithms. The best 3 experimental results are in bold while the best is underlined.

Methods	<i>Food-101</i>	<i>Clothing1M</i>	<i>CIFAR-10N-1</i>	<i>CIFAR-10N-2</i>	<i>CIFAR-10N-3</i>	<i>CIFAR-10N-W</i>
CE	85.15	68.88	85.41 ± 0.24	86.79 ± 0.15	85.41 ± 0.24	80.77 ± 0.24
APL	80.37	54.46	84.40 ± 0.26	84.45 ± 0.50	84.35 ± 0.43	78.16 ± 0.17
PCE	85.72	69.48	63.06 ± 0.37	62.26 ± 0.36	35.47 ± 0.36	33.80 ± 0.33
AUL	82.77	66.25	76.26 ± 0.28	75.24 ± 0.20	75.48 ± 0.40	63.61 ± 1.62
CELC	86.38	69.05	89.77 ± 0.39	89.19 ± 0.46	90.06 ± 0.33	81.16 ± 1.86
Revision	85.70	70.97	90.39 ± 0.12	90.15 ± 0.11	90.07 ± 0.08	83.47 ± 0.27
Identifiability	82.21	67.07	82.52 ± 0.87	81.97 ± 0.85	82.09 ± 0.73	71.62 ± 1.16
Joint	84.74	70.26	88.20 ± 0.29	87.54 ± 0.33	87.67 ± 0.22	84.29 ± 0.40
Co-teaching	83.73	67.94	90.26 ± 0.22	89.82 ± 0.63	90.64 ± 0.47	75.64 ± 4.04
SIGUA	79.68	65.37	87.67 ± 1.18	89.01 ± 0.34	88.40 ± 0.42	80.65 ± 1.29
Co-Dis	86.13	71.60	90.77 ± 0.35	90.22 ± 0.30	90.35 ± 1.12	76.12 ± 3.19
RT-Catoni's	86.13	72.69	91.31 ± 0.25	91.22 ± 0.40	91.23 ± 0.41	84.46 ± 0.41
RT-Log-sum	86.15	72.64	91.37 ± 0.07	91.38 ± 0.23	91.19 ± 0.10	85.03 ± 0.54
RT-Welsch+	85.86	70.81	91.49 ± 0.11	91.26 ± 0.17	91.09 ± 0.17	85.96 ± 1.56

TABLE 11: Test accuracy (%) on three real-world noisy datasets, i.e., *Food-101*, *Clothing1M*, and *CIFAR-10N*. The best 3 experimental results are in bold while the best is underlined.

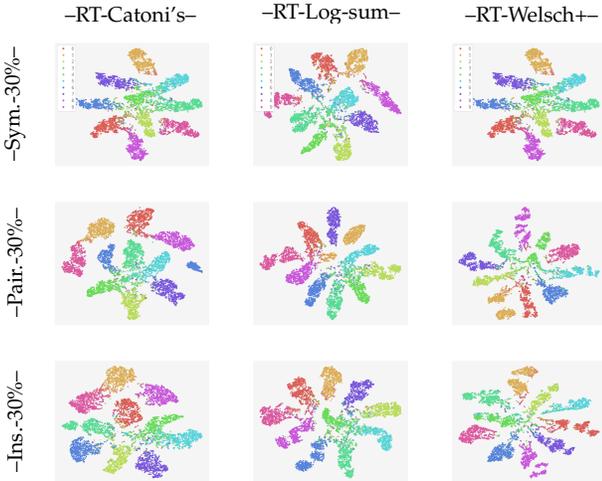


Fig. 9: Visualizations of experimental results using 2D t-SNE [97]. The experiments are conducted on synthetic *CIFAR-10*.

45%, the improvement is significant.

4.2.9 Combination with self-supervised learning

There are some works that employ self-supervised learning [98] to enhance network robustness [60], [99]. Hence, here we follow them and show that our methods can be combined with self-supervised learning to enhance network robustness. Specifically, we use MOCO V2 [100]. After the self-supervised representation learning, the baseline employs the cross-entropy loss and noisily labeled data to fine-tune the linear head. In contrast, our strategies utilize the regularly truncated M-estimators for fine-tuning. Experimental comparisons are provided in Table 10, which demonstrates the utility of our methods. Note that compared with the results in Table 4, we claim that one of the advantages of our methods is plug-and-play for robustness improvement.

4.3 Experiments on real-world noisy datasets

4.3.1 Experimental setup

Datasets. We exploit three real-world noisy datasets to justify our claims, i.e., *Food-101* [101], *Clothing1M* [55], and *CIFAR-10N* [102]¹, which consist of heterogeneous noisy labels. *Food-101* consists of 101 food categories, with 101,000 images.

1. <http://competition.noisylabels.com/>

For each class, 250 manually reviewed clean test images are provided as well as 750 training images. *Clothing1M* has 1M images with real-world noisy labels, and 50k, 14k, 10k images with clean labels for training, validating, and testing, but with 14 classes. Note that we do not use the 50k and 14k clean data in all the experiments, since it is more practical that there is no available clean data. For preprocessing, we resize the image to 256×256 , crop the middle 224×224 as input, and perform normalization. *CIFAR-10N* provides *CIFAR-10* images with human-annotated noisy labels obtained from Amazon Mechanical Turk. Four versions of *CIFAR-10N* label sets are employed here, three of which are labeled by three independent workers (named *CIFAR-10N-1/2/3*) and one of which is negatively aggregated from the above three sets (named *CIFAR-10N-W*). We leave 10% noisy training data as a validation set for model selection.

Network structure and optimizer. We exploit the ResNet-50 network pretrained on ImageNet for *Food-101* and *Clothing1M*. For *Food-101*, we use SGD with momentum 0.9, weight decay 10^{-4} , batch size 128, and an initial learning rate 10^{-2} to train the networks. The learning rate is also divided by 10 after the 40th epoch and 80th epoch. The maximum number of epochs is set to 200. For *Clothing1M*, we also use SGD with momentum 0.9. The batch size and weight decay are adjusted to 32 and 5×10^{-3} . The learning rate is initially set to 10^{-3} and then divided by 10 after the 5th epoch. The maximum number of epochs is set to 20. The experiments on *Food-101* and *Clothing1M* are performed once due to the huge computational cost. For *CIFAR-10N*, a PreAct-ResNet-18 network is exploited. We use SGD with momentum 0.9, weight decay 10^{-3} , batch size 128, and an initial learning rate 10^{-2} . The learning rate is divided by 10 after the 100th epoch. The maximum number of epochs is set to 200. Experiments on *CIFAR-10N* are repeated five times.

4.3.2 Discussions of experimental results

Experimental results on real-world noisy datasets are shown in Table 11. For *Food-101*, the proposed methods achieve great performance. Although the baseline CELC achieves the best performance, the proposed RT-Catoni's and RT-Log-sum achieve competitive performance. For *Clothing1M*, the proposed methods, e.g., RT-Catoni's and RT-Log-sum, achieve clear leads over baselines. For the proposed RT-Welsch+, although it does not outperform the best baseline Co-Dis, it still receives competitive performance. Moreover, since the training procedure of Co-Dis consists of two stages

(estimating the noise rate and performing sample selection), our method can keep an end-to-end manner, and is thus arguably easier to implement. At last, for *CIFAR-10N*, our methods consistently outperform baselines.

5 CONCLUSION

In this paper, we focus on exploiting the sample selection approach to handle noisy labels. We discuss that the prior sample selection procedure has some weaknesses, i.e., ignoring the concerns of noisy labels in selected small-loss examples and neglecting the values of discarded large-loss examples. To relieve two issues at the same time, we propose regularly truncated M-estimators, which can assign different weights to selected small-loss examples and enable large-loss examples to periodically participate in optimization. Theoretically, we discuss the noise-tolerant of truncated M-estimators. Empirically, we conduct a series of experiments to verify the effectiveness of the proposed methods. Extensive experimental results support our claims well. In the future, we are interested in applying our method to data cleaning and robustness enhancement of large-scale pre-trained models [103], [104].

REFERENCES

- [1] Zhi-Hua Zhou. A brief introduction to weakly supervised learning. *National science review*, 5(1):44–53, 2018.
- [2] Yan Yan and Yuhong Guo. Mutual partial label learning with competitive label noise. In *ICLR*, 2023.
- [3] Daniel Paleka and Amartya Sanyal. A law of adversarial risk, interpolation, and label noise. In *ICLR*, 2023.
- [4] Amanda Olmin and Fredrik Lindsten. Robustness and reliability when training with noisy labels. In *AISTATS*, pages 922–942, 2022.
- [5] Lu Jiang, Di Huang, Mason Liu, and Weilong Yang. Beyond synthetic noise: Deep learning on controlled noisy labels. In *ICML*, pages 4804–4815, 2020.
- [6] Ran He, Wei-Shi Zheng, and Bao-Gang Hu. Maximum correntropy criterion for robust face recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(8):1561–1576, 2010.
- [7] Ran He, Wei-Shi Zheng, Tieniu Tan, and Zhenan Sun. Half-quadratic-based iterative minimization for robust sparse representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(2):261–275, 2013.
- [8] Jingchen Ke, Chen Gong, Tongliang Liu, Lin Zhao, Jian Yang, and Dacheng Tao. Laplacian welsch regularization for robust semisupervised learning. *IEEE transactions on cybernetics*, 2020.
- [9] Deep Patel and PS Sastry. Adaptive sample selection for robust learning under label noise. In *WACV*, pages 3932–3942, 2023.
- [10] Ahmet Iscen, Jack Valmadre, Anurag Arnab, and Cordelia Schmid. Learning with neighbor consistency for noisy labels. In *CVPR*, pages 4672–4681, 2022.
- [11] HeeSun Bae, Seungjae Shin, Byeonghu Na, JoonHo Jang, Kyungwoo Song, and Il-Chul Moon. From noisy prediction to true label: Noisy prediction calibration via generative model. In *ICML*, pages 1277–1297, 2022.
- [12] Kevin J Liang, Samrudhdi B Rangrej, Vladan Petrovic, and Tal Hassner. Few-shot learning with noisy labels. In *CVPR*, pages 9089–9098, 2022.
- [13] Shuo Yang, Songhua Wu, Erkun Yang, Bo Han, Yang Liu, Min Xu, Gang Niu, and Tongliang Liu. A parametrical model for instance-dependent label noise. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [14] Amila Silva, Ling Luo, Shanika Karunasekera, and Christopher Leckie. Noise-robust learning from multiple unsupervised sources of inferred labels. In *AAAI*, volume 36, pages 8315–8323, 2022.
- [15] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *ICLR*, 2017.
- [16] Xuefeng Li, Tongliang Liu, Bo Han, Gang Niu, and Masashi Sugiyama. Provably end-to-end label-noise learning without anchor points. In *ICML*, 2021.
- [17] Xiaobo Xia, Tongliang Liu, Bo Han, Nannan Wang, Mingming Gong, Haifeng Liu, Gang Niu, Dacheng Tao, and Masashi Sugiyama. Part-dependent label noise: Towards instance-dependent label noise. In *NeurIPS*, 2020.
- [18] Yingbin Bai and Tongliang Liu. Me-momentum: Extracting hard confident examples from noisily labeled data. In *ICCV*, 2021.
- [19] Xiang Wu, Ran He, Zhenan Sun, and Tieniu Tan. A light cnn for deep face representation with noisy labels. *IEEE Transactions on Information Forensics and Security*, 13(11):2884–2896, 2018.
- [20] Ming-Kun Xie and Sheng-Jun Huang. Ccmn: A general framework for learning with class-conditional multi-label noise. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1):154–166, 2022.
- [21] Ming-Kun Xie and Sheng-Jun Huang. Partial multi-label learning with noisy label identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(7):3676–3687, 2021.
- [22] Enyan Dai, Wei Jin, Hui Liu, and Suhang Wang. Towards robust graph neural networks for noisy graphs with sparse labels. In *WSDM*, pages 181–191, 2022.
- [23] Haobo Wang, Ruixuan Xiao, Yiwen Dong, Lei Feng, and Junbo Zhao. Promix: combating label noise via maximizing clean sample utility. In *IJCAI*, 2023.
- [24] Maria Sofia Bucarelli, Lucas Cassano, Federico Siciliano, Amin Mantrach, and Fabrizio Silvestri. Leveraging inter-rater agreement for classification in the presence of noisy labels. In *CVPR*, pages 3439–3448, 2023.
- [25] Quanming Yao, Hansi Yang, Bo Han, Gang Niu, and James Tin-Yau Kwok. Searching to exploit memorization effect in learning with noisy labels. In *ICML*, pages 10789–10798, 2020.
- [26] Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor W Tsang, and Masashi Sugiyama. How does disagreement benefit co-teaching? In *ICML*, 2019.
- [27] Chuanwen Feng, Yilong Ren, and Xike Xie. Ot-filter: An optimal transport filter for learning with noisy labels. In *CVPR*, 2023.
- [28] 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 2309–2318, 2018.
- [29] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *NeurIPS*, pages 8527–8537, 2018.
- [30] Xiaobo Xia, Tongliang Liu, Bo Han, Mingming Gong, Jun Yu, Gang Niu, and Masashi Sugiyama. Instance correction for learning with open-set noisy labels. *arXiv preprint arXiv:2106.00455*, 2021.
- [31] 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, 2020.
- [32] Eric Arazo, Diego Ortego, Paul Albert, Noel O’Connor, and Kevin McGuinness. Unsupervised label noise modeling and loss correction. In *ICML*, pages 312–321, 2019.
- [33] Bo Han, Gang Niu, Xingrui Yu, Quanming Yao, Miao Xu, Ivor Tsang, and Masashi Sugiyama. Sigua: Forgetting may make learning with noisy labels more robust. In *ICML*, 2020.
- [34] Teng Zhang and Gilad Lerman. A novel m-estimator for robust pca. *The Journal of Machine Learning Research*, 15(1):749–808, 2014.
- [35] Zhilu Zhang and Mert Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. In *NeurIPS*, pages 8778–8788, 2018.
- [36] Xingjun Ma, Hanxun Huang, Yisen Wang, Simone Romano, Sarah Erfani, and James Bailey. Normalized loss functions for deep learning with noisy labels. In *ICML*, pages 6543–6553, 2020.
- [37] Yueming Lyu and Ivor W Tsang. Curriculum loss: Robust learning and generalization against label corruption. In *ICLR*, 2020.
- [38] 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.
- [39] Nontawat Charoenphakdee, Jongyeong Lee, and Masashi Sugiyama. On symmetric losses for learning from corrupted labels. In *ICML*, pages 961–970, 2019.
- [40] Youngdong Kim, Junho Yim, Juseung Yun, and Junmo Kim. Nlnl: Negative learning for noisy labels. In *ICCV*, pages 101–110, 2019.
- [41] Yang Liu and Hongyi Guo. Peer loss functions: Learning from noisy labels without knowing noise rates. In *ICML*, 2020.
- [42] Yilun Xu, Peng Cao, Yuqing Kong, and Yizhou Wang. L_{dmi}: A novel information-theoretic loss function for training deep nets robust to label noise. In *NeurIPS*, pages 6222–6233, 2019.

- [43] Tongliang Liu and Dacheng Tao. Classification with noisy labels by importance reweighting. *IEEE Transactions on pattern analysis and machine intelligence*, 38(3):447–461, 2016.
- [44] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *ICML*, pages 4331–4340, 2018.
- [45] Songhua Wu, Xiaobo Xia, Tongliang Liu, Bo Han, Mingming Gong, Nannan Wang, Haifeng Liu, and Gang Niu. Class2simi: A noise reduction perspective on learning with noisy labels. In *ICML*, 2021.
- [46] Xiaobo Xia, Bo Han, Nannan Wang, Jiankang Deng, Jiatong Li, Yinian Mao, and Tongliang Liu. Extended T: Learning with mixed closed-set and open-set noisy labels. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [47] Yu Yao, Tongliang Liu, Bo Han, Mingming Gong, Jiankang Deng, Gang Niu, and Masashi Sugiyama. Dual t: Reducing estimation error for transition matrix in label-noise learning. In *NeurIPS*, 2020.
- [48] Zhaowei Zhu, Yiwen Song, and Yang Liu. Clusterability as an alternative to anchor points when learning with noisy labels. In *ICML*, pages 12912–12923, 2021.
- [49] Zhaowei Zhu, Tongliang Liu, and Yang Liu. A second-order approach to learning with instance-dependent label noise. In *CVPR*, 2021.
- [50] Yang Liu. Identifiability of label noise transition matrix. *arXiv preprint arXiv:2202.02016*, 2022.
- [51] Seong Min Kye, Kwanghee Choi, Joonyoung Yi, and Buru Chang. Learning with noisy labels by efficient transition matrix estimation to combat label miscorrection. In *ECCV*, pages 717–738, 2022.
- [52] Jacob Goldberger and Ehud Ben-Reuven. Training deep neural networks using a noise adaptation layer. In *ICLR*, 2017.
- [53] Kun Yi and Jianxin Wu. Probabilistic end-to-end noise correction for learning with noisy labels. In *CVPR*, pages 7017–7025, 2019.
- [54] Yikai Zhang, Songzhu Zheng, Pengxiang Wu, Mayank Goswami, and Chao Chen. Learning with feature-dependent label noise: A progressive approach. In *ICLR*, 2021.
- [55] 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, 2015.
- [56] Arash Vahdat. Toward robustness against label noise in training deep discriminative neural networks. In *NeurIPS*, 2017.
- [57] Yuncheng Li, Jianchao Yang, Yale Song, Liangliang Cao, Jiebo Luo, and Li-Jia Li. Learning from noisy labels with distillation. In *ICCV*, pages 1910–1918, 2017.
- [58] Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa. Joint optimization framework for learning with noisy labels. In *CVPR*, 2018.
- [59] Junnan Li, Richard Socher, and Steven C.H. Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. In *ICLR*, 2020.
- [60] Shikun Li, Xiaobo Xia, Shiming Ge, and Tongliang Liu. Selective-supervised contrastive learning with noisy labels. In *CVPR*, pages 316–325, 2022.
- [61] Zhizhong Huang, Junping Zhang, and Hongming Shan. Twin contrastive learning with noisy labels. In *CVPR*, pages 11661–11670, 2023.
- [62] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- [63] 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*, 2015.
- [64] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin Raffel. Mixmatch: A holistic approach to semi-supervised learning. *arXiv preprint arXiv:1905.02249*, 2019.
- [65] Zhaoqing Wang, Ziyu Chen, Yaqian Li, Yandong Guo, Jun Yu, Mingming Gong, and Tongliang Liu. Mosaic representation learning for self-supervised visual pre-training. In *ICLR*, 2023.
- [66] Bo Han, Quanming Yao, Tongliang Liu, Gang Niu, Ivor W Tsang, James T Kwok, and Masashi Sugiyama. A survey of label-noise representation learning: Past, present and future. *arXiv preprint arXiv:2011.04406*, 2020.
- [67] Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy labels with deep neural networks: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [68] 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*, 2017.
- [69] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. MIT Press, 2018.
- [70] Xiaobo Xia, Tongliang Liu, Bo Han, Mingming Gong, Jun Yu, Gang Niu, and Masashi Sugiyama. Sample selection with uncertainty of losses for learning with noisy labels. In *ICLR*, 2022.
- [71] Yisen Wang, Weiyang Liu, Xingjun Ma, James Bailey, Hongyuan Zha, Le Song, and Shu-Tao Xia. Iterative learning with open-set noisy labels. In *CVPR*, pages 8688–8696, 2018.
- [72] Jiacheng Cheng, Tongliang Liu, Kotagiri Ramamohanarao, and Dacheng Tao. Learning with bounded instance-and label-dependent label noise. In *ICML*, 2020.
- [73] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to statistical learning*, volume 112. Springer, 2013.
- [74] Ran He, Tieniu Tan, and Liang Wang. Robust recovery of corrupted low-rankmatrix by implicit regularizers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(4):770–783, 2013.
- [75] Olivier Catoni. Challenging the empirical mean and empirical variance: a deviation study. In *Annales de l’IHP Probabilités et statistiques*, volume 48, pages 1148–1185, 2012.
- [76] Emmanuel J Candes, Michael B Wakin, and Stephen P Boyd. Enhancing sparsity by reweighted l1 minimization. *Journal of Fourier analysis and applications*, 14(5-6):877–905, 2008.
- [77] Weifeng Liu, Puskal P Pokharel, and Jose C Principe. Correntropy: Properties and applications in non-gaussian signal processing. *IEEE Transactions on signal processing*, 55(11):5286–5298, 2007.
- [78] Aritra Ghosh, Himanshu Kumar, and P Shanti Sastry. Robust loss functions under label noise for deep neural networks. In *AAAI*, 2017.
- [79] Rick Durrett. *Probability: theory and examples*, volume 49. Cambridge university press, 2019.
- [80] Naiyang Guan, Tongliang Liu, Yangmuzi Zhang, Dacheng Tao, and Larry S Davis. Truncated cauchy non-negative matrix factorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(1):246–259, 2017.
- [81] Ferenc Nagy. Parameter estimation of the cauchy distribution in information theory approach. *J. UCS*, 12(9):1332–1344, 2006.
- [82] Aditya Krishna Menon, Ankit Singh Rawat, Sashank J. Reddi, and Sanjiv Kumar. Can gradient clipping mitigate label noise? In *ICLR*, 2020.
- [83] Xiong Zhou, Xianming Liu, Deming Zhai, Junjun Jiang, and Xiangyang Ji. Asymmetric loss functions for noise-tolerant learning: Theory and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [84] Hongxin Wei, Huiqing Zhuang, Renchunzi Xie, Lei Feng, Gang Niu, Bo An, and Yixuan Li. Mitigating memorization of noisy labels by clipping the model prediction. In *ICML*, 2023.
- [85] Xiaobo Xia, Tongliang Liu, Nannan Wang, Bo Han, Chen Gong, Gang Niu, and Masashi Sugiyama. Are anchor points really indispensable in label-noise learning? In *NeurIPS*, 2019.
- [86] Yang Liu, Hao Cheng, and Kun Zhang. Identifiability of label noise transition matrix. In *ICML*, pages 21475–21496, 2023.
- [87] Xiaobo Xia, Bo Han, Yibing Zhan, Jun Yu, Mingming Gong, Chen Gong, and Tongliang Liu. Combating noisy labels with sample selection by mining high-discrepancy examples. In *ICCV*, 2023.
- [88] Duc Tam Nguyen, Chaithanya Kumar Mummadi, Thi Phuong Nhung Ngo, Thi Hoai Phuong Nguyen, Laura Beggel, and Thomas Brox. Self: Learning to filter noisy labels with self-ensembling. In *ICLR*, 2020.
- [89] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. The MNIST database of handwritten digits.
- [90] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [91] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- [92] Ken Lang. Newsweeder: Learning to filter netnews. In *Machine Learning Proceedings*, pages 331–339, 1995.
- [93] Kimin Lee, Sukmin Yun, Kibok Lee, Honglak Lee, Bo Li, and Jinwoo Shin. Robust inference via generative classifiers for handling noisy labels. In *ICML*, pages 3763–3772, 2019.
- [94] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014.

- [95] Xiaobo Xia, Tongliang Liu, Bo Han, Chen Gong, Nannan Wang, Zongyuan Ge, and Yi Chang. Robust early-learning: Hindering the memorization of noisy labels. In *ICLR*, 2021.
- [96] 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, 2017.
- [97] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [98] Zhaoqing Wang, Qiang Li, Guoxin Zhang, Pengfei Wan, Wen Zheng, Nannan Wang, Mingming Gong, and Tongliang Liu. Exploring set similarity for dense self-supervised representation learning. In *CVPR*, pages 16590–16599, 2022.
- [99] Tong Wei, Jiang-Xin Shi, Wei-Wei Tu, and Yu-Feng Li. Robust long-tailed learning under label noise. *arXiv preprint arXiv:2108.11569*, 2021.
- [100] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- [101] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *ECCV*, pages 446–461, 2014.
- [102] Jiaheng Wei, Zhaowei Zhu, Hao Cheng, Tongliang Liu, Gang Niu, and Yang Liu. Learning with noisy labels revisited: A study using real-world human annotations. In *ICLR*, 2022.
- [103] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023.
- [104] Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *arXiv preprint arXiv:2306.09927*, 2023.

A SUPPLEMENTARY THEORETICAL ANALYSIS

A.1 Preliminary knowledge

We denote the underlying clean dataset corresponding to the noisy dataset S , as $S^* = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, where y_i is the unobserved clean label of \tilde{y}_i . Given any loss function, ψ , and a classifier, f , we define the ψ -risk

$$R_\psi(f) := \mathbb{E}_{(\mathbf{x}, y) \sim S^*} [\psi(f(\mathbf{x}), y)]. \quad (6)$$

Under the risk minimization framework, the objective is to learn a classifier, f , which is a global minimizer of R_ψ . Note that the ψ -risk, R_ψ , depends on ψ , the loss function. When ψ happens to be the 0-1 loss, R_ψ would be the usual Bayes risk. Let f^* be the global minimizer (over the chosen function class) of $R_\psi(f)$. In this paper, ψ will be the loss function composed of the cross-entropy loss and truncated M-estimators.

Then the notion and notations about the label noise model are introduced. We have

$$\tilde{y}_i = \begin{cases} y_i & \text{with probability } (1 - \eta_{\mathbf{x}_i}) \\ j, j \in [k], j \neq y_i & \text{with probability } \bar{\eta}_{\mathbf{x}, j}. \end{cases} \quad (7)$$

Note that for all \mathbf{x} , conditioned on $y = i$, we have $\sum_{j \neq i} \bar{\eta}_{\mathbf{x}, j} = \eta_{\mathbf{x}}$. The label noise model is termed *symmetric* or *uniform* if $\eta_{\mathbf{x}} = \eta$, and $\bar{\eta}_{\mathbf{x}, j} = \frac{\eta}{k-1}$, $\forall j \neq y, \forall \mathbf{x}$, where η is a constant. Noise is said to be *simple non-uniform* when the noise rate $\eta_{\mathbf{x}}$ is a function of \mathbf{x} . A simple special case is when $\bar{\eta}_{\mathbf{x}, j} = \frac{\eta_{\mathbf{x}}}{k-1}$, $\forall j \neq y$. We define it as simple non-uniform noise. Then ψ -risk of a classifier f under noisy data is defined

$$R_\psi^\eta(f) := \mathbb{E}_{(\mathbf{x}, \tilde{y}) \sim S} [\psi(f(\mathbf{x}), \tilde{y})]. \quad (8)$$

A.2 Proof of Lemma 1

Our proofs are inspired by [78]. Recall that $R_\psi(f) := \mathbb{E}_{(\mathbf{x}, y) \sim S^*} [\psi(f(\mathbf{x}), y)]$. For symmetric noise, we have, for any f ,

$$\begin{aligned} R_\psi^\eta(f) &= \mathbb{E}_{(\mathbf{x}, \tilde{y}) \sim S} [\psi(f(\mathbf{x}), \tilde{y})] \\ &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{y|\mathbf{x}} \mathbb{E}_{\tilde{y}|\mathbf{x}, y} \psi(f(\mathbf{x}), \tilde{y}) \\ &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{y|\mathbf{x}} [(1 - \eta) \psi(f(\mathbf{x}), y) + \frac{\eta}{k-1} \sum_{i \neq y} \psi(f(\mathbf{x}), i)]. \end{aligned} \quad (9)$$

Note that in this paper, ψ will be the loss function that is composed of the cross-entropy loss and truncated M-estimators. Although the cross-entropy loss is not upper-bounded, with our truncation mechanism, ψ will be upper-bounded, since the largest value of ψ is limited. Therefore, for any ψ , we denote its lower and upper bounds of the sum of loss values as c_1 and c_2 respectively, i.e., $c_1 \leq \sum_i \psi(f(\mathbf{x}), i) \leq c_2$. In this way,

$$\begin{aligned} &R_\psi^\eta(f^*) - R_\psi^\eta(f) \\ &\leq \frac{c_2 \eta}{k-1} + (1 - \frac{\eta k}{k-1}) R_\psi(f^*) - \frac{c_1 \eta}{k-1} - (1 - \frac{\eta k}{k-1}) R_\psi(f) \\ &\leq \frac{(c_2 - c_1) \eta}{k-1} + \frac{k-1 - \eta k}{k-1} (R_\psi(f^*) - R_\psi(f)) \\ &= \frac{(c_2 - c_1) \eta}{k-1} + \frac{k-1 - \eta k}{k-1} \Delta(\psi, f) \\ &= \frac{(c_2 - c_1 - k \Delta(\psi, f)) \eta + (k-1) \Delta(\psi, f)}{k-1} \end{aligned} \quad (10)$$

If the noise rate $\eta < \frac{(1-k)\Delta(\psi, f)}{c_2 - c_1 - k\Delta(\psi, f)}$, we will have

$$\begin{aligned} &R_\psi^\eta(f^*) - R_\psi^\eta(f) \\ &\leq \frac{(1-k)\Delta(\psi, f) + (k-1)\Delta(\psi, f)}{k-1} \\ &= 0. \end{aligned} \quad (11)$$

This proves f^* is also a minimizer of the risk under symmetric noise. Proof completed.

A.3 Proof of Corollary 1

Corollary 1. *In a multi-class classification problem, the truncated M-estimators are noise-tolerant under the simple non-uniform noise, if $c_2 - c_1 - k\Delta(\psi, f, \mathbf{x}) > 0$ and the noise rate $\eta_{\mathbf{x}} < \frac{(1-k)\Delta(\psi, f, \mathbf{x})}{c_2 - c_1 - k\Delta(\psi, f, \mathbf{x})}$. Here c_1 and c_2 denote the lower and upper bounds of the sum of the losses obtained by predictions on all classes, and $\Delta(\psi, f, \mathbf{x}) = \sup(\psi(f^*(\mathbf{x}), y) - \psi(f(\mathbf{x}), y))$.*

The proof of Corollary 1 is as follows. For the simple non-uniform noise, we derive that

$$\begin{aligned}
R_\psi^\eta(f) &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{y|\mathbf{x}} \left[(1 - \eta_{\mathbf{x}}) \psi(f(\mathbf{x}), y) + \frac{\eta_{\mathbf{x}}}{k-1} \sum_{i \neq y} \psi(f(\mathbf{x}), i) \right] \\
&= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{y|\mathbf{x}} \left[(1 - \eta_{\mathbf{x}}) \psi(f(\mathbf{x}), y) \right. \\
&\quad \left. + \frac{\eta_{\mathbf{x}}}{k-1} \left(\sum_i \psi(f(\mathbf{x}), i) - \psi(f(\mathbf{x}), y) \right) \right] \\
&= \mathbb{E} \left[\frac{k-1-\eta_{\mathbf{x}}k}{k-1} \psi(f(\mathbf{x}), y) + \frac{\eta_{\mathbf{x}}}{k-1} \sum_i \psi(f(\mathbf{x}), i) \right].
\end{aligned} \tag{12}$$

Therefore,

$$\begin{aligned}
&R_\psi^\eta(f^*) - R_\psi^\eta(f) \\
&\leq \frac{(c_2 - c_1 - k\Delta(\psi, f, \mathbf{x})\eta + (k-1)\Delta(\psi, f, \mathbf{x}))}{k-1},
\end{aligned} \tag{13}$$

where $\Delta(\psi, f, \mathbf{x}) = \sup(\psi(f^*(\mathbf{x}), y) - \psi(f(\mathbf{x}), y))$. Therefore, if the noise rate $\eta_{\mathbf{x}} < \frac{(1-k)\Delta(\psi, f, \mathbf{x})}{c_2 - c_1 - k\Delta(\psi, f, \mathbf{x})}$, we have $R_\psi^\eta(f^*) - R_\psi^\eta(f) \leq 0$. This proves f^* is also a minimizer of the risk under the simple non-uniform noise.