
Sample Selection with Uncertainty of Losses for Learning with Noisy Labels

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In learning with noisy labels, the *sample selection* approach is very popular, which
2 regards *small-loss* data as correctly labeled during training. However, losses are
3 generated on-the-fly based on the model being trained with noisy labels, and thus
4 *large-loss* data are *likely but not certainly* to be incorrect. There are actually
5 two possibilities of a large-loss data point: (a) it is mislabeled, and then its loss
6 *decreases slower* than other data, since deep neural networks “learn patterns first”;
7 (b) it belongs to an underrepresented group of data and *has not been selected yet*. In
8 this paper, we incorporate the uncertainty of losses by adopting *interval estimation*
9 instead of *point estimation* of losses, where lower bounds of the *confidence intervals*
10 of losses derived from *distribution-free concentration inequalities*, but not losses
11 themselves, are used for sample selection. In this way, we also give large-loss but
12 less selected data a try; then, we can better distinguish between the cases (a) and
13 (b) by seeing if the losses *effectively decrease* with the uncertainty after the try. As
14 a result, we can better explore underrepresented data that are correctly labeled but
15 seem to be mislabeled *at first glance*. Experiments demonstrate that the proposed
16 method is superior to baselines and robust to a broad range of label noise types.

17 1 Introduction

18 Learning with noisy labels is one of the most challenging problems in weakly-supervised learning,
19 since noisy labels are ubiquitous in the real world [36, 65, 40, 1, 61]. For instance, both crowdsourcing
20 and web crawling yield large numbers of noisy labels everyday [12]. Noisy labels can severely impair
21 the performance of deep neural networks with strong memorization capacities [67, 69, 42, 30].

22 To reduce the influence of noisy labels, a lot of approaches have been recently proposed [38, 29, 31,
23 68, 71, 55, 56, 46, 33, 25, 34, 47, 60, 49, 19, 17, 14]. They can be generally divided into two main
24 categories. The first one is to estimate the noise transition matrix [41, 44, 15, 11], which denotes the
25 probabilities that clean labels flip into noisy labels. However, the noise transition matrix is hard to be
26 estimated accurately, especially when the number of classes is large [65]. The second approach is
27 sample selection, which is *our focus* in this paper. This approach is based on selecting possibly clean
28 examples from a mini-batch for training [12, 62, 50, 65, 23, 50, 51]. Intuitively, if we can exploit less
29 noisy data for network parameter updates, the network will be more robust.

30 A major question in sample selection is what *criteria* can be used to select possibly clean examples.
31 At the present stage, the selection based on the *small-loss* criteria is the most common method, and
32 has been verified to be effective in many circumstances [12, 16, 65, 52, 62]. Specifically, since
33 deep networks *learn patterns first* [2], they would first memorize training data of clean labels and
34 then those of noisy labels with the assumption that clean labels are of the majority in a noisy class.
35 Small-loss examples can thus be regarded as clean examples *with high probability*. Therefore, in
36 each iteration, prior methods [12, 52] select the small-loss examples based on *the predictions of the*
37 *current network* for robust training.

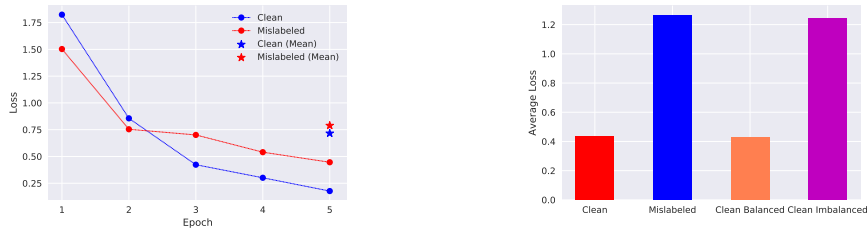


Figure 1: Illustrations of *uncertainty of losses*. Experiments are conducted on the imbalanced noisy *MNIST* dataset. **Left**: uncertainty of *small-loss* examples. At the beginning of training (Epochs 1 and 2), due to the instability of the current prediction, the network gives a larger loss to the clean example and does not select it for updates. If we consider the *mean* of training losses at different epochs, the clean example can be equipped with a smaller loss and then selected for updates. **Right**: uncertainty of *large-loss* examples. Since the deep network learns easy examples at the beginning of training, it gives a large loss to *clean imbalanced* data with non-dominant labels, which causes such data unable to be selected and severely influence generalization.

38 However, such a selection procedure is *debatable*, since it arguably does *not consider uncertainty*
 39 in selection. The uncertainty comes from two aspects. First, this procedure has *uncertainty about*
 40 *small-loss examples*. Specifically, the procedure uses *limited time intervals* and only exploits the
 41 losses provided by the *current predictions*. For this reason, the estimation for the noisy class posterior
 42 is *unstable* [63], which causes the network predictions to be equally unstable. It thus *takes huge risks*
 43 to only use losses provided by the current predictions (Figure 1, left). Once wrong selection is made,
 44 the inferiority of accumulated errors will arise [65]. Second, this procedure has *uncertainty about*
 45 *large-loss examples*. To be specific, deep networks learn easy examples at the beginning of training,
 46 but ignore some clean examples with large losses. Nevertheless, such examples are always critical for
 47 generalization. For instance, when learning with *imbalanced* data, distinguishing the examples with
 48 *non-dominant labels* are more pivotal during training [35]. Deep networks often give large losses to
 49 such examples (Figure 1, right). Therefore, when learning under the realistic scenes, e.g., learning
 50 with noisy imbalanced data, prior sample selection methods cannot address such an issue well.

51 To relieve the above issues, we study the uncertainty of losses in the sample selection procedure to
 52 combat noisy labels. To reduce the uncertainty of small-loss examples, we extend time intervals and
 53 utilize the *mean* of training losses at different training iterations. In consideration of the bad influence
 54 of mislabeled data on training losses, we build two *robust mean estimators* from the perspectives of
 55 *soft truncation* and *hard truncation* w.r.t. the truncation level, respectively. Soft truncation makes the
 56 mean estimation more robust by *holistically* changing the behavior of losses. Hard truncation makes
 57 the mean estimation more robust by *locally* removing outliers from losses. To reduce the uncertainty
 58 of large-loss examples, we encourage networks to pick the sample that has not been selected in a
 59 conservative way. Furthermore, to address the two issues *simultaneously*, we derive *concentration*
 60 *inequalities* [5] for robust mean estimation and further employ statistical *confidence bounds* [3] to
 61 consider the number of times an example was selected during training.

62 The study of uncertainty of losses in learning with noisy labels can be justified as follows. In statistical
 63 learning, it is known that uncertainty is related to the quality of data [48]. Philosophically, we need
 64 *variety decrease* for selected data and *variety search* for unselected data, which share a common
 65 objective, i.e., *reduce the uncertainty of data to improve generalization* [37]. This is our original
 66 intention, since noisy labels could bring more uncertainty because of the low quality of noisy data.
 67 Nevertheless, due to the harm of noisy labels for generalization, we need to strike a good balance
 68 between variety decrease and search. Technically, our method is specially designed for handling
 69 noisy labels, which robustly uses network predictions and conservatively seeks less selected examples
 70 meanwhile to reduce the uncertainty of losses and then generalize well.

71 Before delving into details, we clearly emphasize our contributions in two folds. First, we reveal prior
 72 sample selection criteria in learning with noisy labels have some potential weaknesses and discuss
 73 them in detail. The new selection criteria are then proposed with detailed theoretical analyses. Second,
 74 we experimentally validate the proposed method on both synthetic noisy balanced/imbalanced datasets
 75 and real-world noisy datasets, on which it achieves superior robustness compared with the state-
 76 of-the-art methods in learning with noisy labels. The rest of the paper is organized as follows. In
 77 Section 2, we propose our robust learning paradigm step by step. Experimental results are discussed
 78 in Section 3. The conclusion is given in Section 4.

79 2 Method

80 In this section, we first introduce the problem setting and some background (Section 2.1). Then we
81 discuss how to exploit training losses at different iterations (Section 2.2). Finally, we introduce the
82 proposed method, which exploits training losses at different iterations more robustly and encourages
83 networks to pick the sample that is less selected but could be correctly labeled (Section 2.3).

84 2.1 Preliminaries

85 Let \mathcal{X} and \mathcal{Y} be the input and output spaces. Consider a k -class classification problem, i.e., $\mathcal{Y} = [k]$,
86 where $[k] = \{1, \dots, k\}$. In learning with noisy labels, the training data are all sampled from a
87 corrupted distribution on $\mathcal{X} \times \mathcal{Y}$. We are given a sample with noisy labels, i.e., $\tilde{S} = \{(\mathbf{x}, \tilde{y})\}$, where
88 \tilde{y} is the noisy label. The aim is to learn a robust classifier that could assign clean labels to test data by
89 only exploiting a training sample with noisy labels.

90 Let $f : \mathcal{X} \rightarrow \mathbb{R}^k$ be the classifier with learnable parameters \mathbf{w} . At the i -th iteration during training,
91 the parameters of the classifier f can be denoted as \mathbf{w}_i . Let $\ell : \mathbb{R}^k \times \mathcal{Y} \rightarrow \mathbb{R}$ be a *surrogate loss*
92 *function* for k -class classification. We exploit the *softmax cross entropy loss* in this paper. Given an
93 arbitrary training example (\mathbf{x}, \tilde{y}) , at the i -th iteration, we can obtain a loss ℓ_i , i.e., $\ell_i = \ell(f(\mathbf{w}_i; \mathbf{x}), \tilde{y})$.
94 Hence, until the t -th iteration, we can obtain a training loss set L_t about the example (\mathbf{x}, \tilde{y}) , i.e.,
95 $L_t = \{\ell_1, \dots, \ell_t\}$.

96 In this paper, we assume that the training losses in L_t conform to a *Markov process*, which is to
97 represent a changing system under the assumption that future states only depend on the current state
98 (the Markov property) [43]. More specifically, at the i -th iteration, if we exploit an optimization
99 algorithm for parameter updates (e.g., the stochastic gradient descent algorithm [4]) and omit other
100 dependencies (e.g., \tilde{S}), we will have $P(\mathbf{w}_i | \mathbf{w}_{i-1}, \dots, \mathbf{w}_0) = P(\mathbf{w}_i | \mathbf{w}_{i-1})$, which means that the
101 future state of the classifier f only depends on the current state. Furthermore, given a training example
102 and the parameters of the classifier f , we can determine the loss of the training example as discussed.
103 Therefore, the training losses in L_t will also conform to a Markov process.

104 2.2 Extended Time Intervals

105 As limited time interval cannot address the instability issue of the estimation for the noisy class
106 posterior well [42], we extend time intervals and exploit the training losses at different training
107 iterations for sample selection. One straightforward idea is to use the *mean* of training losses at
108 different training iterations. Hence, the selection criterion could be

$$\tilde{\mu} = \frac{1}{t} \sum_{i=1}^t \ell_i. \quad (1)$$

109 It is intuitive and reasonable to use such a selection criterion for sample selection, since the operation
110 of averaging can mitigate the risks caused by the unstable estimation for the noisy class posterior,
111 following better generalization. Nevertheless, such a method could arguably achieve suboptimal
112 classification performance for learning with noisy labels. The main reason is that, due to the great
113 harm of mislabeled data, part of training losses are with too large uncertainty and could be seen as
114 outliers. Therefore, it could be biased to use the mean of training losses consisting of such outliers
115 [10], which further influences sample selection. More evaluations for our claims are provided in
116 Section 3.

117 2.3 Robust Mean Estimation and Conservative Search

118 We extend time intervals and meanwhile exploit the training losses at different training iterations more
119 robustly. Specifically, we build two robust mean estimators from the perspectives of *soft truncation*
120 and *hard truncation* [7]. Note that for specific tasks, it is feasible to decide the types of robust mean
121 estimation with statistical tests based on some assumptions [8]. We leave the analysis as future work.
122 Two *distribution-free* robust mean estimators are introduced as follows.

123 **Soft truncation.** We extend a classical M-estimator from [7] and exploit the *widest* possible choice of
124 the *influence function*. More specifically, give a random variable X , let us consider a non-decreasing

125 influence function $\psi : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$\psi(X) = \log(1 + X + X^2/2), X \geq 0. \quad (2)$$

126 The choice of ψ is inspired by the *Taylor expansion of the exponential function*, which can make the
 127 estimation results more robust by reducing the side effect of extremum *holistically*. The illustration
 128 for this influence function is provided in Appendix A.1. For our task, given the observations on
 129 training losses, i.e., $L_t = \{\ell_1, \dots, \ell_t\}$, we estimate the mean robustly as follows:

$$\tilde{\mu}_s = \frac{1}{t} \sum_{i=1}^t \psi(\ell_i). \quad (3)$$

130 We term the above robust mean estimator (3) the *soft estimator*.

131 **Hard truncation.** We propose a new robust mean estimator based on hard truncation. Specifically,
 132 given the observations on training losses L_t , we first exploit the K-nearest neighbor (KNN) algorithm
 133 [27] to remove some underlying outliers in L_t . The number of outliers is denoted by t_o ($t_o < t$), which
 134 can be *adaptively determined* as discussed in [70]. Note that we can also employ other algorithms,
 135 e.g., principal component analysis [45] and the local outlier factor [6], to identify underlying outliers
 136 in L_t . The main reason we employ KNN is because of its relatively low computation costs [70].

137 The truncated loss observations on training losses are denoted by L_{t-t_o} . We then utilize L_{t-t_o} for
 138 the mean estimation. As the potential outliers are removed with high probability, the robustness of
 139 the estimation results will be enhanced. We denote such an estimated mean as $\tilde{\mu}_h$. We have

$$\tilde{\mu}_h = \frac{1}{t-t_o} \sum_{\ell_i \in L_{t-t_o}} \ell_i. \quad (4)$$

140 The corresponding estimator (4) is termed the *hard estimator*.

141 We derive concentration inequalities for the soft and hard estimators respectively. The search strategy
 142 for less selected examples and overall selection criterion are then provided. Note that we do not need
 143 to explicitly quantify the mean of training losses. We only need to sort the training examples based
 144 on the proposed selection criterion and then use the selected examples for robust training.

145 **Theorem 1.** Let $Z_n = \{z_1, \dots, z_n\}$ be an observation set with mean μ_z and variance σ^2 . By
 146 exploiting the non-decreasing influence function $\psi(z) = \log(1 + z + z^2/2)$. For any $\epsilon > 0$, we have

$$\left| \frac{1}{n} \sum_{i=1}^n \psi(z_i) - \mu_z \right| \leq \frac{\sigma^2 (n + \frac{\sigma^2 \log(\epsilon^{-1})}{n^2})}{n - \sigma^2}, \quad (5)$$

147 with probability at least $1 - 2\epsilon$.

148 Proof can be found in Appendix A.1.

149 **Theorem 2.** Let $Z_n = \{z_1, \dots, z_n\}$ be a (not necessarily time homogeneous) Markov chain with
 150 mean μ_z , taking values in a Polish state space $\Lambda_1 \times \dots \times \Lambda_n$, and with a minimal mixing time τ_{\min} .
 151 The truncated set with hard truncation is denoted by Z_{n_o} , with $n_o < n$. If $|z_i|$ is upper bounded by Z .
 152 For any $\epsilon_1 > 0$ and $\epsilon_2 > 0$, we have

$$\left| \frac{1}{n - n_o} \sum_{z_i \in Z_n \setminus Z_{n_o}} - \mu_z \right| \leq \frac{1}{n - n_o} \left(2Z \sqrt{2\tau_{\min} \log \frac{2}{\epsilon_1}} + \frac{2Zn_o}{n} \sqrt{2\tau_{\min} \log \frac{2n}{\epsilon_2}} \right), \quad (6)$$

153 with probability at least $1 - \epsilon_1 - \epsilon_2$.

154 Proof can be found in Appendix A.2. For our task, let the training loss be upper-bounded by L . The
 155 value of L can be determined easily by training networks on noisy datasets and observing the loss
 156 distribution [1].

157 **Conservative search and selection criteria.** In this paper, we will use the concentration inequalities
 158 (5) and (6) to present conservative search and the overall sample selection criterion. Specifically,
 159 we exploit their *lower bounds* and consider the selected number of examples during training. The
 160 selection of the examples that are less selected is encouraged.

Algorithm 1 CNLCU Algorithm.

1: **Input** θ_1 and θ_2 , learning rate η , fixed τ , epoch T_k and T_{\max} , iteration t_{\max} ;
for $T = 1, 2, \dots, T_{\max}$ **do**
 2: **Shuffle** training dataset \tilde{S} ;
 for $t = 1, \dots, t_{\max}$ **do**
 3: **Fetch** mini-batch \tilde{S} from \tilde{S} ;
 4: **Obtain** $\tilde{S}_1 = \arg \min_{S': |S'| \geq R(T)|\tilde{S}|} \ell^*(\theta_1, S')$; // calculated with Eq. (7) or Eq. (8)
 5: **Obtain** $\tilde{S}_2 = \arg \min_{S': |S'| \geq R(T)|\tilde{S}|} \ell^*(\theta_2, S')$; // calculated with Eq. (7) or Eq. (8)
 6: **Update** $\theta_1 = \theta_1 - \eta \nabla \ell(\theta_1, \tilde{S}_2)$;
 7: **Update** $\theta_2 = \theta_2 - \eta \nabla \ell(\theta_2, \tilde{S}_1)$;
 end
 8: **Update** $R(T) = 1 - \min \left\{ \frac{T}{T_k} \tau, \tau \right\}$;
end
9: **Output** θ_1 and θ_2 .

161 Denote the number of times one example was selected by $n_t (n_t \leq t)$. Let $\epsilon = \frac{1}{2t}$. For the
162 circumstance with soft truncation, the selection criterion is

$$\ell_s^* = \tilde{\mu}_s - \frac{\sigma^2(t + \frac{\sigma^2 \log(2t)}{t^2})}{n_t - \sigma^2}. \quad (7)$$

163 Let $\epsilon_1 = \epsilon_2 = \frac{1}{2t}$, for the situation with hard truncation, by rewriting (6), the selection criterion is

$$\ell_h^* = \tilde{\mu}_h - \frac{2\sqrt{2\tau_{\min}}L(t + \sqrt{2}t_o)}{(t - t_o)\sqrt{t}} \sqrt{\frac{\log(4t)}{n_t}}. \quad (8)$$

164 Note that we directly replace t with n_t . If an example is rarely selected during training, n_t will be far
165 less than n , which causes the lower bounds to change drastically. Hence, we do not use the mean of
166 all training losses, but use the mean of training losses in fixed-length time intervals. More details
167 about this can be checked in Section 3.

168 For the selection criteria (7) and (8), we can see that they consist of two terms and have one term
169 with a minus sign. The first term in Eq. (7) (or Eq. (8)) is to reduce the uncertainty of small-loss
170 examples, where we use robust mean estimation on training losses. The second term, i.e., the
171 statistical confidence bound, is to encourage the network to choose the less selected examples (with a
172 small n_t). The two terms are constraining and balanced with σ^2 or τ_{\min} . To avoid introducing strong
173 assumptions on the underlying distribution of losses [8], we tune σ and τ_{\min} with a noisy validation
174 set. For the mislabeled data, although the model has high uncertainties on them (i.e., a small n_t)
175 and tends to pick them, the overfitting to the mislabeled data is harmful. Also, the mislabeled data
176 and clean data are rather hard to distinguish in some cases as discussed. Thus, we should search
177 underlying clean data in a conservative way. In this paper, we initialize σ and τ_{\min} with small values.
178 This way can reduce the adverse effects of mislabeled data and meanwhile select the clean examples
179 with large losses, which helps generalize. More evaluations will be presented in Section 3.

180 The overall procedure of the proposed method, which combats noisy labels by concerning uncertainty
181 (CNLCU), is provided in Algorithm 1. CNLCU works in a mini-batch manner since all deep learning
182 training methods are based on stochastic gradient descent. Following [12], we exploit two networks
183 with parameters θ_1 and θ_2 respectively to teach each other. Specifically, when a mini-batch \tilde{S} is
184 formed (Step 3), we let two networks select a small proportion of examples in this mini-batch with
185 Eq. (7) or (8) (Step 4 and Step 5). The number of instances is controlled by the function $R(T)$, and
186 two networks only select $R(T)$ percentage of examples out of the mini-batch. The value of $R(T)$
187 should be larger at the beginning of training, and be smaller when the number of epochs goes large,
188 which can make better use of memorization effects of deep networks [12] for sample selection. Then,
189 the selected instances are fed into its peer network for parameter updates (Step 6 and Step 7).

190 3 Experiments

191 In this section, we evaluate the robustness of our proposed method to noisy labels with comprehensive
192 experiments on the synthetic balanced noisy datasets (Section 3.1), synthetic imbalanced noisy
193 datasets (Section 3.2), and real-world noisy dataset (Section 3.3).

194 3.1 Experiments on Synthetic Balanced Noisy Datasets

195 **Datasets.** We verify the effectiveness of our method on the manually corrupted version of the
196 following datasets: *MNIST* [22], *F-MNIST* [58], *CIFAR-10* [21], and *CIFAR-100* [21], because
197 these datasets are popularly used for the evaluation of learning with noisy labels in the literature
198 [12, 65, 54, 23]. The four datasets are class-balanced. The important statistics of the used synthetic
199 datasets are summarized in Appendix B.1.

200 **Generating noisy labels.** We consider broad types of label noise: (1). Symmetric noise (abbreviated
201 as Sym.) [53, 31, 26]. (2) Asymmetric noise (abbreviated as Asym.) [32, 57, 52]. (3) Pairflip noise
202 (abbreviated as Pair.) [12, 65, 71]. (4). Tridiagonal noise (abbreviated as Trid.) [68]. (5). Instance
203 noise (abbreviated as Ins.) [9, 56]. The noise rate is set to 20% and 40% to ensure clean labels are
204 diagonally dominant [32]. More details about above noise are provided in Appendix B.1. We leave
205 out 10% of noisy training examples as a validation set.

206 **Baselines.** We compare the proposed method (Algorithm 1) with following methods which focus on
207 sample selection, and implement all methods with default parameters by PyTorch, and conduct all the
208 experiments on NVIDIA Titan Xp GPUs. (1). S2E [62], which properly controls the sample selection
209 process so that deep networks can better benefit from the memorization effects. (2). MentorNet [16],
210 which learns a curriculum to filter out noisy data. We use self-paced MentorNet in this paper. (3).
211 Co-teaching [12], which trains two networks simultaneously and cross-updates parameters of peer
212 networks. (4). SIGUA [13], which exploits stochastic integrated gradient underweighted ascent to
213 handle noisy labels. We use self-teaching SIGUA in this paper. (5). JoCor [52], which reduces the
214 diversity of networks to improve robustness. Other types of baselines such as *adding regularization*
215 are provided in Appendix B.2. Note that we do not compare the proposed method with some state-
216 of-the-art methods, e.g., SELF [39] and DivideMix [24]. It is because their proposed methods are
217 aggregations of multiple techniques. We mainly focus on sample selection in learning with noisy
218 labels. Therefore, the comparison is not fair. Here, we term our methods with soft truncation and
219 hard truncation as CNLCU-S and CNLCU-H respectively.

220 **Network structure and optimizer.** For *MNIST*, *F-MNIST*, and *CIFAR-10*, we use a 9-layer CNN
221 structure from [12]. Due to the limited space, the experimental details on *CIFAR-100* are provided
222 in Appendix B.3. All network structures we used here are standard test beds for weakly-supervised
223 learning. For all experiments, the Adam optimizer [20] (momentum=0.9) is used with an initial
224 learning rate of 0.001, and the batch size is set to 128 and we run 200 epochs. We linearly decay
225 learning rate to zero from 80 to 200 epochs as did in [12]. We take two networks with the same
226 architecture but different initializations as two classifiers as did in [12, 65, 52], since even with the
227 same network and optimization method, different initializations can lead to different local optimal
228 [12]. The details of network structures can be checked in Appendix C.

229 For the hyper-parameters σ^2 and τ_{\min} , we determine them in the range $\{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}\}$
230 with a noisy validation set. Here, we assume the noise level τ is known and set $R(T) = 1 -$
231 $\min\{\frac{T}{T_k}\tau, \tau\}$ with $T_k=10$. If τ is not known in advanced, it can be inferred using validation sets
232 [29, 66]. As for performance measurement, we use test accuracy, i.e., *test accuracy* = (# of correct
233 prediction) / (# of testing). All experiments are repeated five times. We report the mean and standard
234 deviation of experimental results.

235 **Experimental results.** The experimental results about test accuracy are provided in Table 1, 2, and
236 3. Specifically, for *MNIST*, as can be seen, our proposed methods, i.e., CNLCU-S and CNLCU-H,
237 produce the best results in the vast majority of cases. In some cases such as asymmetric noise, the
238 baseline S2E outperforms ours, which benefits the accurate estimation for the number of selected
239 small-loss examples. For *F-MNIST*, the training data becomes complicated. S2E cannot achieve the
240 accurate estimation in such situation and thus has no great performance like it got on *MNIST*. Our
241 methods achieve varying degrees of lead over baselines. For *CIFAR-10*, our methods once again
242 outperforms all the baseline methods. Although some baseline, e.g., Co-teaching, can work well
243 in some cases, experimental results show that it cannot handle various noise types. In contrast, the
244 proposed methods achieve superior robustness against broad noise types. The results mean that our
245 methods can be better applied to actual scenarios, where the noise is diversiform.

246 **Ablation study.** We first conduct the ablation study to analyze the sensitivity of the length of time
247 intervals. In order to *avoid too dense figures*, we exploit *MNIST* and *F-MNIST* with the mentioned
248 noise settings as representative examples. For CNLCU-S, the length of time intervals is chosen in

Noise type	Sym.		Asym.		Pair.		Trid.		Ins.	
	20%	40%	20%	40%	20%	40%	20%	40%	20%	40%
Method/Noise ratio	20%	40%	20%	40%	20%	40%	20%	40%	20%	40%
S2E	98.46	95.62	99.05	98.45	98.56	94.22	99.02	97.23	97.93	94.02
	± 0.06	± 0.91	± 0.02	± 0.26	± 0.32	± 0.79	± 0.09	± 1.26	± 1.26	± 2.39
MentorNet	95.04	92.08	96.32	90.86	93.19	90.93	96.42	93.28	94.65	90.11
	± 0.03	± 0.42	± 0.17	± 0.97	± 0.17	± 1.54	± 0.09	± 1.37	± 0.73	± 1.26
Co-teaching	97.53	95.62	98.25	95.08	96.05	94.16	98.05	96.18	97.96	95.02
	± 0.12	± 0.30	± 0.08	± 0.43	± 0.96	± 1.37	± 0.06	± 0.85	± 0.09	± 0.39
SIGUA	92.31	91.88	93.96	62.59	93.77	86.22	94.92	83.46	92.90	86.34
	± 1.10	± 0.92	± 0.82	± 0.15	± 1.40	± 1.75	± 0.83	± 2.98	± 1.82	± 3.51
JoCor	98.42	98.04	98.05	94.55	98.01	96.85	98.45	96.98	98.62	96.07
	± 0.14	± 0.07	± 0.37	± 1.08	± 0.19	± 0.43	± 0.17	± 0.25	± 0.06	± 0.31
CNLCU-S	98.82	98.31	98.93	97.67	98.86	97.71	99.09	98.02	98.77	97.78
	± 0.03	± 0.05	± 0.06	± 0.22	± 0.06	± 0.64	± 0.04	± 0.17	± 0.08	± 0.25
CNLCU-H	98.70	98.24	99.01	98.01	98.44	97.37	98.89	97.92	98.74	97.42
	± 0.06	± 0.06	± 0.04	± 0.03	± 0.19	± 0.32	± 0.15	± 0.05	± 0.16	± 0.39

Table 1: Test accuracy (%) on *MNIST* over the last ten epochs. The best two results are in bold.

Noise type	Sym.		Asym.		Pair.		Trid.		Ins.	
	20%	40%	20%	40%	20%	40%	20%	40%	20%	40%
Method/Noise ratio	20%	40%	20%	40%	20%	40%	20%	40%	20%	40%
S2E	89.99	75.32	89.00	81.03	88.66	67.09	89.53	77.29	88.65	79.35
	± 2.07	± 5.84	± 0.95	± 1.93	± 1.32	± 4.03	± 2.63	± 3.97	± 2.12	± 3.04
MentorNet	90.37	86.53	89.69	67.21	87.92	83.70	88.74	85.63	87.52	83.27
	± 0.17	± 0.65	± 0.19	± 2.94	± 1.08	± 0.49	± 0.33	± 0.59	± 0.15	± 1.42
Co-teaching	91.48	88.80	91.03	68.07	90.77	86.91	91.24	89.18	90.60	87.90
	± 0.10	± 0.29	± 0.14	± 4.58	± 0.23	± 0.71	± 0.11	± 0.36	± 0.12	± 0.45
SIGUA	87.64	87.23	76.97	45.96	69.59	68.93	79.97	76.14	76.92	74.89
	± 1.29	± 0.72	± 2.59	± 3.40	± 5.75	± 2.80	± 3.23	± 4.24	± 5.09	± 4.84
JoCor	91.97	89.96	90.95	79.79	91.52	87.40	92.01	89.42	91.43	87.59
	± 0.13	± 0.19	± 0.21	± 2.39	± 0.24	± 0.58	± 0.17	± 0.33	± 0.71	± 0.94
CNLCU-S	92.37	91.45	92.57	83.14	92.04	88.20	92.24	90.08	91.69	89.02
	± 0.15	± 0.28	± 0.15	± 1.77	± 0.26	± 0.44	± 0.17	± 0.34	± 0.10	± 1.02
CNLCU-H	92.42	91.60	92.60	82.69	91.70	87.70	92.33	90.22	91.50	88.79
	± 0.21	± 0.19	± 0.18	± 0.43	± 0.18	± 0.69	± 0.26	± 0.71	± 0.21	± 1.22

Table 2: Test accuracy on *F-MNIST* over the last ten epochs. The best two results are in bold.

249 the range from 3 to 8. For CNLCU-H, the length of time intervals is chosen in the range from 10 to
250 15. Note that the reason for their different lengths is that their different mechanisms. Specifically,
251 CNLCU-S holistically changes the behavior of losses, but does not remove any loss from the loss set.
252 We thus do not need too long length of time intervals. As a comparison, CNLCU-H needs to remove
253 some outliers from the loss set as discussed. The length should be longer to guarantee the number of
254 examples available for robust mean estimation. The experimental results are provided in Appendix
255 B.4, which show the proposed CNLCU-S and CNLCU-H are robust to the choices of the length of
256 time intervals. Such robustness to hyperparameters means our methods can be applied in practice and
257 does not need too much effort to tune the hyperparameters.

258 Furthermore, since our methods concern uncertainty from two aspects, i.e., the uncertainty from both
259 small-loss and large-loss examples, we conduct experiments to analyze each part of our methods.
260 Also, as mentioned, we compare robust mean estimation with non-robust mean estimation when
261 learning with noisy labels. More details are provided in Appendix B.4.

262 3.2 Experiments on Synthetic Imbalanced Noisy Datasets

263 **Experimental setup.** We exploit *MNIST* and *F-MNIST*. For these two datasets, we reduce the number
264 of training examples along with the labels from “0” to “4” to 1% of previous numbers. We term
265 such synthetic imbalanced noisy datasets as *IM-MNIST* and *IM-F-MNIST* respectively. This setting
266 aims to simulate the extremely imbalanced circumstance, which is common in practice. Moreover,
267 we exploit asymmetric noise, since these types of noise can produce more imbalanced case [41, 32].
268 Other settings such as the network structure and optimizer are the same as those in experiments on
269 synthetic balanced noisy datasets.

Noise type	Sym.		Asym.		Pair.		Trid.		Ins.	
	20%	40%	20%	40%	20%	40%	20%	40%	20%	40%
Method/Noise ratio	20%	40%	20%	40%	20%	40%	20%	40%	20%	40%
S2E	80.78 ±0.88	69.72 ±3.94	84.03 ±1.01	75.04 ±1.24	81.72 ±0.93	61.50 ±4.63	81.44 ±0.59	64.39 ±2.82	79.89 ±0.26	62.42 ±3.11
MentorNet	80.92 ±0.48	74.67 ±1.17	80.37 ±0.26	71.69 ±1.06	77.98 ±0.31	69.39 ±1.73	78.02 ±0.29	71.56 ±0.93	77.02 ±0.71	68.17 ±2.52
Co-teaching	82.35 ±0.16	77.96 ±0.39	83.87 ±0.24	73.43 ±0.62	80.94 ±0.46	72.81 ±0.92	81.17 ±0.60	74.37 ±0.64	79.92 ±0.57	73.29 ±1.62
SIGUA	78.19 ±0.22	77.67 ±0.41	75.14 ±0.36	52.76 ±0.68	74.41 ±0.81	61.91 ±5.27	75.75 ±0.53	74.05 ±0.41	74.34 ±0.39	67.98 ±1.34
JoCor	80.96 ±0.25	76.65 ±0.43	81.39 ±0.74	69.92 ±1.63	80.33 ±0.20	71.62 ±1.05	79.03 ±0.13	74.33 ±1.09	78.21 ±0.34	71.46 ±1.27
CNLCU-S	83.03 ±0.21	78.25 ±0.70	85.06 ±0.17	75.34 ±0.32	83.16 ±0.25	73.19 ±1.25	82.77 ±0.32	74.37 ±1.37	82.03 ±0.37	73.67 ±1.09
CNLCU-H	83.03 ±0.47	78.33 ±0.50	84.95 ±0.27	75.29 ±0.80	83.39 ±0.68	73.40 ±1.53	82.52 ±0.71	74.79 ±1.13	81.93 ±0.25	73.58 ±1.39

Table 3: Test accuracy (%) on *CIFAR-10* over the last ten epochs. The best two results are in bold.

As for performance measurements, we use test accuracy. In addition, we exploit the selected ratio of training examples with the imbalanced classes, i.e., $selected\ ratio = (\#\ of\ selected\ imbalanced\ labels / \#\ of\ all\ selected\ labels)$. Intuitively, a higher selected ratio means the proposed method can make better use of training examples with the imbalanced classes, following better generalization [18].

Experimental results. The test accuracy achieved on *IM-MNIST* and *IM-F-MNIST* is presented in Figure 2. Recall the experimental results in Table 1 and 2, we can see that the imbalanced issue is *catastrophic* to the sample selection approach when learning with noisy labels. For *IM-MNIST*, as can be seen, all the baselines have serious overfitting in the early stages of training. The curves of test accuracy drop dramatically. As a comparison, the proposed CNLCU-S and CNLCU-H can give a try to large-loss but less selected data which are possible to be clean but equipped with imbalanced labels. Therefore, our methods always outperform baselines clearly. In the case of Asym. 10%, our methods achieve nearly 30% lead over baselines. For *IM-F-MNIST*, we can also see that our methods perform well and always achieve about 5% lead over all the baselines. Note that due to the huge challenge of this task, some baseline, e.g., S2E, has a large error bar. In addition, the baseline SIGUA performs badly. It is because SIGUA exploits stochastic integrated gradient underweighted ascent on large-loss examples, which makes the examples with imbalanced classes more difficult to be selected than them in other sample selection methods.

The selected ratio achieved on *IM-MNIST* and *IM-F-MNIST* is presented in Table 4. The results explain well why our methods perform better on synthetic imbalanced noisy datasets, i.e., our methods can make better use of training examples with the imbalanced classes. Note that since we give a try to large-loss but less selected data in a conservative way, the selected ratio is still far away from the class prior probability on the test set, i.e., 10%. However, a little improvement of the selection ratio can bring a considerable improvement of test accuracy. These results tell us that, in the sample selection approach when learning with noisy labels, improving the selected ratio of training examples with the imbalanced classes is challenging but promising for generalization. This practical problem deserves to be studied in depth.

3.3 Experiments on Real-world Noisy Datasets

Experimental setup. To verify the efficacy of our methods in the real-world scenario, we conduct experiments on the noisy dataset *ClothingIM* [59]. Specifically, for experiments on *ClothingIM*, we use the 1M images with noisy labels for training and 10k clean data for test respectively. Note that we do not use the 50k clean training data in all the experiments. For preprocessing, we resize the image to 256×256 , crop the middle 224×224 as input, and perform normalization. The experiments on *ClothingIM* are performed once due to the huge computational cost. We leave 10% noisy training data as a validation set for model selection. Note that we do not exploit the resampling trick during training [24]. Here, *Best* denotes the test accuracy of the epoch where the validation accuracy was optimal. *Last* denotes test accuracy of the last epoch. For the experiments on *ClothingIM*, we use a ResNet-18 pretrained on ImageNet as did in [52]. We also use the Adam optimizer and set the batch size to 64. During the training stage, we run 15 epochs in total and set the learning rate 8×10^{-4} , 5×10^{-4} , and 5×10^{-5} for 5 epochs each.

Dataset	<i>IM-MNIST</i>				<i>IM-F-MNIST</i>			
	Method/Noise ratio	10%	20%	30%	40%	10%	20%	30%
S2E	0.13 ± 0.12	0.11 ± 0.05	0.09 ± 0.02	0.05 ± 0.01	0.13 ± 0.04	0.17 ± 0.03	0.16 ± 0.02	0.12 ± 0.04
MentorNet	0.10 ± 0.02	0.15 ± 0.02	0.12 ± 0.03	0.13 ± 0.02	0.12 ± 0.01	0.15 ± 0.03	0.09 ± 0.01	0.14 ± 0.02
Co-teaching	0.09 ± 0.03	0.07 ± 0.02	0.05 ± 0.01	0.12 ± 0.01	0.17 ± 0.05	0.04 ± 0.00	0.13 ± 0.04	0.07 ± 0.01
SIGUA	0.04 ± 0.00	0.04 ± 0.00	0.01 ± 0.00	0.02 ± 0.00	0.03 ± 0.00	0.02 ± 0.00	0.04 ± 0.00	0.00 ± 0.00
JoCor	0.11 ± 0.04	0.08 ± 0.01	0.07 ± 0.03	0.06 ± 0.02	0.05 ± 0.01	0.13 ± 0.04	0.13 ± 0.03	0.07 ± 0.02
CNLCU-S	0.60 ± 0.11	0.37 ± 0.09	0.39 ± 0.04	0.38 ± 0.06	0.35 ± 0.03	0.39 ± 0.04	0.36 ± 0.03	0.30 ± 0.02
CNLCU-H	0.57 ± 0.13	0.32 ± 0.01	0.37 ± 0.07	0.32 ± 0.05	0.34 ± 0.02	0.35 ± 0.06	0.32 ± 0.04	0.28 ± 0.03

Table 4: Selected ratio (%) on *IM-MNIST* and *IM-F-MNIST*. The best two results are in bold.

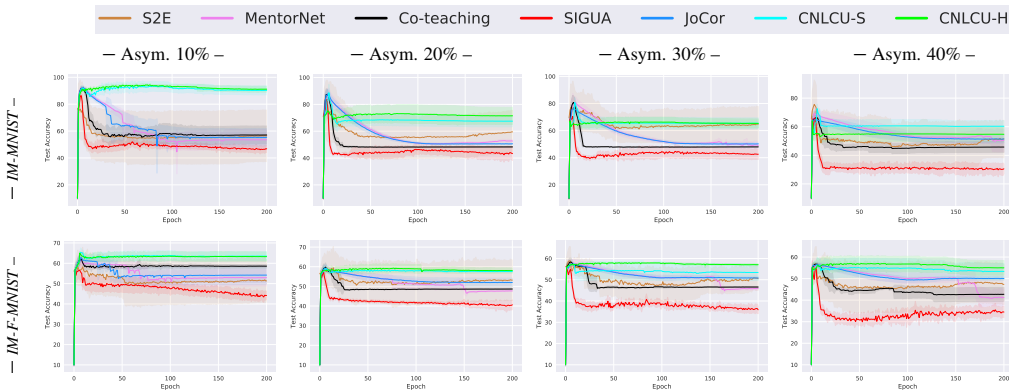


Figure 2: Test accuracy vs. number of epochs on *IM-MNIST* and *IM-F-MNIST*. The error bar for standard deviation in each figure has been shaded.

309 **Experimental results.** The results on *ClothingIM* are provided in Table 5. Specifically, the proposed
310 methods get better results than state-of-the-art methods on *Best*, which achieve an improvement of
311 +1.28% and +0.99% over the best baseline JoCor. Likewise, the proposed methods outperform all the
312 baselines on *Last*. We achieve an improvement of +1.01% and +0.54% over JoCor. Note that the
313 results are a bit lower than some state-of-art methods, e.g., [64] and [46], because of the following
314 reasons. (1). We follow [52] and use ResNet-18 as a backbone. The state-of-art methods [64, 46]
315 use ResNet-50 as a backbone. Our aim is to make the experimental results directly comparable with
316 previous papers [52] in the same area. (2). We only focus on the sample selection approach and do
317 not employ other advanced techniques, e.g., introducing the prior distribution [46] and combining
semi-supervised learning [24, 39, 28].

Methods	S2E	MentorNet	Co-teaching	SIGUA	JoCor	CNLCU-S	CNLCU-H
<i>Best</i>	67.34	68.36	69.37	62.89	70.09	71.37	71.08
<i>Last</i>	65.90	67.42	68.62	58.73	69.75	70.76	70.29

Table 5: Test accuracy (%) on *ClothingIM*. The best two results are in bold.

318 4 Conclusion

319
320 In this paper, we focus on promoting the prior sample selection in learning with noisy labels, which
321 starts from concerning the uncertainty of losses during training. We robustly use the training losses at
322 different iterations to reduce the uncertainty of small-loss examples, and adopt confidence interval
323 estimation to reduce the uncertainty of large-loss examples. Experiments are conducted on benchmark
324 datasets, demonstrating the effectiveness of our method. We believe that this paper opens up new
325 possibilities in the topics of using sample selection to handle noisy labels, especially in improving
326 the robustness of models on imbalanced noisy datasets.

327 **References**

- 328 [1] Eric Arazo, Diego Ortego, Paul Albert, Noel O’Connor, and Kevin McGuinness. Unsupervised
329 label noise modeling and loss correction. In *ICML*, pages 312–321, 2019.
- 330 [2] Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio,
331 Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A
332 closer look at memorization in deep networks. In *ICML*, pages 233–242, 2017.
- 333 [3] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine*
334 *Learning Research*, 3(Nov):397–422, 2002.
- 335 [4] Léon Bottou. Stochastic gradient descent tricks. In *Neural networks: Tricks of the trade*, pages
336 421–436. Springer, 2012.
- 337 [5] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A*
338 *nonasymptotic theory of independence*. Oxford university press, 2013.
- 339 [6] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. Lof: identifying
340 density-based local outliers. In *SIGMOD*, pages 93–104, 2000.
- 341 [7] Olivier Catoni. Challenging the empirical mean and empirical variance: a deviation study. In
342 *Annales de l’IHP Probabilités et statistiques*, volume 48, pages 1148–1185, 2012.
- 343 [8] Arijit Chakrabarty and Gennady Samorodnitsky. Understanding heavy tails in a bounded world
344 or, is a truncated heavy tail heavy or not? *Stochastic models*, 28(1):109–143, 2012.
- 345 [9] Jiacheng Cheng, Tongliang Liu, Kotagiri Ramamohanarao, and Dacheng Tao. Learning with
346 bounded instance-and label-dependent label noise. In *ICML*, 2020.
- 347 [10] Ilias Diakonikolas, Daniel M Kane, and Ankit Pensia. Outlier robust mean estimation with
348 subgaussian rates via stability. *arXiv preprint arXiv:2007.15618*, 2020.
- 349 [11] Bo Han, Jiangchao Yao, Gang Niu, Mingyuan Zhou, Ivor Tsang, Ya Zhang, and Masashi
350 Sugiyama. Masking: A new perspective of noisy supervision. In *NeurIPS*, pages 5836–5846,
351 2018.
- 352 [12] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi
353 Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels.
354 In *NeurIPS*, pages 8527–8537, 2018.
- 355 [13] Bo Han, Gang Niu, Xingrui Yu, Quanming Yao, Miao Xu, Ivor Tsang, and Masashi Sugiyama.
356 Sigua: Forgetting may make learning with noisy labels more robust. In *ICML*, pages 4006–4016,
357 2020.
- 358 [14] Hrayr Harutyunyan, Kyle Reing, Greg Ver Steeg, and Aram Galstyan. Improving generalization
359 by controlling label-noise information in neural network weights. In *ICML*, pages 4071–4081,
360 2020.
- 361 [15] Dan Hendrycks, Mantas Mazeika, Duncan Wilson, and Kevin Gimpel. Using trusted data to
362 train deep networks on labels corrupted by severe noise. In *NeurIPS*, 2018.
- 363 [16] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. MentorNet: Learning
364 data-driven curriculum for very deep neural networks on corrupted labels. In *ICML*, pages
365 2309–2318, 2018.
- 366 [17] Lu Jiang, Di Huang, Mason Liu, and Weilong Yang. Beyond synthetic noise: Deep learning on
367 controlled noisy labels. In *ICML*, pages 4804–4815, 2020.
- 368 [18] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and
369 Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition. In *ICLR*,
370 2020.
- 371 [19] Youngdong Kim, Junho Yim, Juseung Yun, and Junmo Kim. Nlnl: Negative learning for noisy
372 labels. In *ICCV*, pages 101–110, 2019.

- 373 [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
374 *arXiv:1412.6980*, 2014.
- 375 [21] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- 376 [22] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. The MNIST database of handwritten
377 digits. <http://yann.lecun.com/exdb/mnist/>.
- 378 [23] Kimin Lee, Sukmin Yun, Kibok Lee, Honglak Lee, Bo Li, and Jinwoo Shin. Robust inference
379 via generative classifiers for handling noisy labels. In *ICML*, pages 3763–3772, 2019.
- 380 [24] Junnan Li, Richard Socher, and Steven C.H. Hoi. Dividemix: Learning with noisy labels as
381 semi-supervised learning. In *ICLR*, 2020.
- 382 [25] Mingchen Li, Mahdi Soltanolkotabi, and Samet Oymak. Gradient descent with early stopping
383 is provably robust to label noise for overparameterized neural networks. In *AISTATS*, 2020.
- 384 [26] Xuefeng Li, Tongliang Liu, Bo Han, Gang Niu, and Masashi Sugiyama. Provably end-to-end
385 label-noise learning without anchor points. 2021.
- 386 [27] Yihua Liao and V Rao Vemuri. Use of k-nearest neighbor classifier for intrusion detection.
387 *Computers & security*, 21(5):439–448, 2002.
- 388 [28] Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-
389 learning regularization prevents memorization of noisy labels. In *NeurIPS*, 2020.
- 390 [29] Tongliang Liu and Dacheng Tao. Classification with noisy labels by importance reweighting.
391 *IEEE Transactions on pattern analysis and machine intelligence*, 38(3):447–461, 2016.
- 392 [30] Michal Lukasik, Srinadh Bhojanapalli, Aditya Menon, and Sanjiv Kumar. Does label smoothing
393 mitigate label noise? In *ICML*, pages 6448–6458, 2020.
- 394 [31] Xingjun Ma, Yisen Wang, Michael E Houle, Shuo Zhou, Sarah M Erfani, Shu-Tao Xia, Sudanthi
395 Wijewickrema, and James Bailey. Dimensionality-driven learning with noisy labels. In *ICML*,
396 pages 3361–3370, 2018.
- 397 [32] Xingjun Ma, Hanxun Huang, Yisen Wang, Simone Romano, Sarah Erfani, and James Bailey.
398 Normalized loss functions for deep learning with noisy labels. In *ICML*, pages 6543–6553,
399 2020.
- 400 [33] Eran Malach and Shai Shalev-Shwartz. Decoupling "when to update" from "how to update". In
401 *NeurIPS*, pages 960–970, 2017.
- 402 [34] Aditya Krishna Menon, Brendan Van Rooyen, and Nagarajan Natarajan. Learning from binary
403 labels with instance-dependent noise. *Machine Learning*, 107(8-10):1561–1595, 2018.
- 404 [35] Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit,
405 and Sanjiv Kumar. Long-tail learning via logit adjustment. *arXiv preprint arXiv:2007.07314*,
406 2020.
- 407 [36] Baharan Mirzasoleiman, Kaidi Cao, and Jure Leskovec. Coresets for robust training of neural
408 networks against noisy labels. In *NeurIPS*, 2020.
- 409 [37] David S Moore. Uncertainty. *On the shoulders of giants: New approaches to numeracy*, pages
410 95–137, 1990.
- 411 [38] Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning
412 with noisy labels. In *NeurIPS*, pages 1196–1204, 2013.
- 413 [39] Duc Tam Nguyen, Chaithanya Kumar Mummadi, Thi Phuong Nhung Ngo, Thi Hoai Phuong
414 Nguyen, Laura Beggel, and Thomas Brox. Self: Learning to filter noisy labels with self-
415 ensembling. In *ICLR*, 2020.
- 416 [40] Kento Nishi, Yi Ding, Alex Rich, and Tobias Höllerer. Augmentation strategies for learning
417 with noisy labels. *arXiv preprint arXiv:2103.02130*, 2021.

- 418 [41] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu.
419 Making deep neural networks robust to label noise: A loss correction approach. In *CVPR*, pages
420 1944–1952, 2017.
- 421 [42] Geoff Pleiss, Tianyi Zhang, Ethan R Elenberg, and Kilian Q Weinberger. Identifying mislabeled
422 data using the area under the margin ranking. In *NeurIPS*, 2020.
- 423 [43] Jeffrey S Rosenthal. Faithful couplings of markov chains: now equals forever. *Advances in*
424 *Applied Mathematics*, 18(3):372–381, 1997.
- 425 [44] Jun Shu, Qian Zhao, Zengben Xu, and Deyu Meng. Meta transition adaptation for robust deep
426 learning with noisy labels. *arXiv preprint arXiv:2006.05697*, 2020.
- 427 [45] Mei-Ling Shyu, Shu-Ching Chen, Kanoksri Sarinnapakorn, and LiWu Chang. A novel anomaly
428 detection scheme based on principal component classifier. Technical report, 2003.
- 429 [46] Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa. Joint optimization
430 framework for learning with noisy labels. In *CVPR*, 2018.
- 431 [47] Kiran K Thekumparampil, Ashish Khetan, Zinan Lin, and Sewoong Oh. Robustness of
432 conditional gans to noisy labels. In *NeurIPS*, pages 10271–10282, 2018.
- 433 [48] Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media,
434 2013.
- 435 [49] Qizhou Wang, Jiangchao Yao, Chen Gong, Tongliang Liu, Mingming Gong, Hongxia Yang,
436 and Bo Han. Learning with group noise. In *AAAI*, 2021.
- 437 [50] Xiaobo Wang, Shuo Wang, Jun Wang, Hailin Shi, and Tao Mei. Co-mining: Deep face
438 recognition with noisy labels. In *ICCV*, pages 9358–9367, 2019.
- 439 [51] Yisen Wang, Weiyang Liu, Xingjun Ma, James Bailey, Hongyuan Zha, Le Song, and Shu-Tao
440 Xia. Iterative learning with open-set noisy labels. In *CVPR*, pages 8688–8696, 2018.
- 441 [52] Hongxin Wei, Lei Feng, Xiangyu Chen, and Bo An. Combating noisy labels by agreement: A
442 joint training method with co-regularization. In *CVPR*, pages 13726–13735, 2020.
- 443 [53] Pengxiang Wu, Songzhu Zheng, Mayank Goswami, Dimitris Metaxas, and Chao Chen. A
444 topological filter for learning with label noise. In *NeurIPS*, 2020.
- 445 [54] Songhua Wu, Xiaobo Xia, Tongliang Liu, Bo Han, Mingming Gong, Nannan Wang, Haifeng
446 Liu, and Gang Niu. Class2simi: A noise reduction perspective on learning with noisy labels. In
447 *ICML*, 2021.
- 448 [55] Xiaobo Xia, Tongliang Liu, Nannan Wang, Bo Han, Chen Gong, Gang Niu, and Masashi
449 Sugiyama. Are anchor points really indispensable in label-noise learning? In *NeurIPS*, pages
450 6835–6846, 2019.
- 451 [56] Xiaobo Xia, Tongliang Liu, Bo Han, Nannan Wang, Mingming Gong, Haifeng Liu, Gang Niu,
452 Dacheng Tao, and Masashi Sugiyama. Part-dependent label noise: Towards instance-dependent
453 label noise. In *NeurIPS*, 2020.
- 454 [57] Xiaobo Xia, Tongliang Liu, Bo Han, Chen Gong, Nannan Wang, Zongyuan Ge, and Yi Chang.
455 Robust early-learning: Hindering the memorization of noisy labels. In *ICLR*, 2021.
- 456 [58] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for
457 benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 458 [59] Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. Learning from massive
459 noisy labeled data for image classification. In *CVPR*, pages 2691–2699, 2015.
- 460 [60] Yilun Xu, Peng Cao, Yuqing Kong, and Yizhou Wang. L_dmi: A novel information-theoretic
461 loss function for training deep nets robust to label noise. In *NeurIPS*, pages 6222–6233, 2019.

- 462 [61] Shuo Yang, Lu Liu, and Min Xu. Free lunch for few-shot learning: Distribution calibration. In
463 *ICLR*, 2021.
- 464 [62] Quanming Yao, Hansi Yang, Bo Han, Gang Niu, and James Tin-Yau Kwok. Searching to exploit
465 memorization effect in learning with noisy labels. In *ICML*, pages 10789–10798, 2020.
- 466 [63] Yu Yao, Tongliang Liu, Bo Han, Mingming Gong, Jiankang Deng, Gang Niu, and Masashi
467 Sugiyama. Dual t: Reducing estimation error for transition matrix in label-noise learning. In
468 *NeurIPS*, 2020.
- 469 [64] Kun Yi and Jianxin Wu. Probabilistic end-to-end noise correction for learning with noisy labels.
470 In *CVPR*, pages 7017–7025, 2019.
- 471 [65] Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor W Tsang, and Masashi Sugiyama. How
472 does disagreement benefit co-teaching? In *ICML*, 2019.
- 473 [66] Xiyu Yu, Tongliang Liu, Mingming Gong, Kayhan Batmanghelich, and Dacheng Tao. An
474 efficient and provable approach for mixture proportion estimation using linear independence
475 assumption. In *CVPR*, pages 4480–4489, 2018.
- 476 [67] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
477 deep learning requires rethinking generalization. In *ICLR*, 2017.
- 478 [68] Yivan Zhang, Gang Niu, and Masashi Sugiyama. Learning noise transition matrix from only
479 noisy labels via total variation regularization. In *ICML*, 2021.
- 480 [69] Zhilu Zhang and Mert Sabuncu. Generalized cross entropy loss for training deep neural networks
481 with noisy labels. In *NeurIPS*, pages 8778–8788, 2018.
- 482 [70] Yue Zhao, Zain Nasrullah, and Zheng Li. Pyod: A python toolbox for scalable outlier detection.
483 *Journal of Machine Learning Research*, 20(96):1–7, 2019.
- 484 [71] Songzhu Zheng, Pengxiang Wu, Aman Goswami, Mayank Goswami, Dimitris Metaxas, and
485 Chao Chen. Error-bounded correction of noisy labels. In *ICML*, pages 11447–11457, 2020.

486 Checklist

487 The checklist follows the references. Please read the checklist guidelines carefully for information on
488 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
489 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
490 the appropriate section of your paper or providing a brief inline description. For example:

- 491 • Did you include the license to the code and datasets? **[No]** The code and the data are
492 proprietary.

493 Please do not modify the questions and only use the provided macros for your answers. Note that the
494 Checklist section does not count towards the page limit. In your paper, please delete this instructions
495 block and only keep the Checklist section heading above along with the questions/answers below.

496 1. For all authors...

- 497 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
498 contributions and scope? **[Yes]**
- 499 (b) Did you describe the limitations of your work? **[Yes]**
- 500 (c) Did you discuss any potential negative societal impacts of your work? **[No]**
- 501 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
502 them? **[Yes]**

503 2. If you are including theoretical results...

- 504 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
- 505 (b) Did you include complete proofs of all theoretical results? **[Yes]**

506 3. If you ran experiments...

- 507 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
508 imental results (either in the supplemental material or as a URL)? **[Yes]** The code
509 and instructions are provided in the supplemental material. The used datasets can be
510 publicly downloaded. Besides, the code for generating noisy labels is provided.
- 511 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
512 were chosen)? **[Yes]** See Section 3.
- 513 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
514 ments multiple times)? **[Yes]** See Section 3.1 and 3.2.
- 515 (d) Did you include the total amount of compute and the type of resources used (e.g., type
516 of GPUs, internal cluster, or cloud provider)? **[Yes]** See “Baselines” in Section 3.1.

517 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 518 (a) If your work uses existing assets, did you cite the creators? **[Yes]** We use *MNIST*,
519 *F-MNIST*, *CIFAR-10*, *CIFAR-100*, and *Clothing1M* in this paper. We cite the creators,
520 which can be checked in Section 3.
- 521 (b) Did you mention the license of the assets? **[N/A]**
- 522 (c) Did you include any new assets either in the supplemental material or as a URL? **[N/A]**
- 523
- 524 (d) Did you discuss whether and how consent was obtained from people whose data you’re
525 using/curating? **[N/A]**
- 526 (e) Did you discuss whether the data you are using/curating contains personally identifiable
527 information or offensive content? **[N/A]**

528 5. If you used crowdsourcing or conducted research with human subjects...

- 529 (a) Did you include the full text of instructions given to participants and screenshots, if
530 applicable? **[N/A]**
- 531 (b) Did you describe any potential participant risks, with links to Institutional Review
532 Board (IRB) approvals, if applicable? **[N/A]**
- 533 (c) Did you include the estimated hourly wage paid to participants and the total amount
534 spent on participant compensation? **[N/A]**