# Supplementary Material of "Sample Selection with Uncertainty of Losses for Learning with Noisy Labels"

**Anonymous Author(s)**
Affiliation
Address
email

## A Proof of Theoretical Results

### A.1 Proof of Theorem 1

For the circumstance with soft truncation, $\tilde{\mu}_z = \frac{1}{n} \sum_{i=1}^{n} \psi(z_i)$. As suggested in [1], we can exploit $\tilde{\mu}_z^-$ and $\tilde{\mu}_z^+$ such that

$$\tilde{\mu}_z^- \leq \tilde{\mu}_z \leq \tilde{\mu}_z^+, \tag{1}$$

to derive a bound for $\tilde{\mu}_z$. For some positive real parameter $\alpha$, we define

$$r(\tilde{\mu}_z) = \sum_{i=1}^{n} \psi\left[\alpha(z_i - \tilde{\mu}_z)\right] = 0. \tag{2}$$

Let us introduce the quantity

$$r(\theta) = \frac{1}{\alpha n} \sum_{i=1}^{n} \psi\left[\alpha(z_i - \theta)\right]. \tag{3}$$

With the exponential moment inequality [5] and the $C_r$ inequality [9], we have

$$\exp\{\alpha n r(\theta)\} \leq \left\{1 + \alpha(\mu_z - \theta) + \alpha^2[\sigma^2 + (\mu_z - \theta)^2]\right\}^n$$
$$\leq \exp\{n\alpha(\mu_z - \theta) + n\alpha^2[\sigma^2 + (\mu_z - \theta)^2]\}. \tag{4}$$

In the same way,

$$\exp\{-\alpha n r(\theta)\} \leq \exp\{-n\alpha(\mu_z - \theta) + n\alpha^2[\sigma^2 + (\mu_z - \theta)^2]\}. \tag{5}$$

If we define for any $\mu_s \in \mathbb{R}$ the bounds

$$B_-(\theta) = \mu_z - \theta - \alpha[\sigma^2 + (\mu_z - \theta)^2] - \frac{\log(\epsilon^{-1})}{\alpha n} \tag{6}$$

and

$$B_+(\theta) = \mu_z - \theta + \alpha[\sigma^2 + (\mu_z - \theta)^2] + \frac{\log(\epsilon^{-1})}{\alpha n}. \tag{7}$$

From [2] (Lemma 2.2), we obtain that

$$P(r(\theta) > B_-(\theta)) \geq 1 - \epsilon \quad \text{and} \quad P(r(\theta) < B_+(\theta)) \geq 1 - \epsilon. \tag{8}$$

Let $\tilde{\mu}_z^-$ be the largest solution of the quadratic equation $B_-(\theta)$ and $\tilde{\mu}_z^+$ be the smallest solution of the quadratic equation $B_+(\theta)$. Also, to guarantee the solution of the quadratic equation, we assume

$$4\alpha^2\sigma^2 + \frac{4\log(\epsilon^{-1})}{n} \leq 1. \tag{9}$$

From [2] (Theorem 2.6), we then have

$$\tilde{\mu}_z^- \geq \mu_z - \frac{\alpha\sigma^2 + \frac{\log(\epsilon^{-1})}{\alpha n}}{\alpha - 1}, \tag{10}$$

and

$$\tilde{\mu}_z^+ \leq \mu_z + \frac{\alpha\sigma^2 + \frac{\log(\epsilon^{-1})}{\alpha n}}{\alpha - 1}. \tag{11}$$

With probability at least 1-2$\epsilon$, we have $\tilde{\mu}_z^- \leq \tilde{\mu}_z \leq \tilde{\mu}_z^+$. We can choose $\alpha = \frac{n}{\sigma^2}$. Then we have

$$|\tilde{\mu}_z - \mu_z| \leq \frac{\sigma^2(n + \frac{\sigma^2 \log(\epsilon^{-1})}{n^2})}{n - \sigma^2}, \tag{12}$$

which holds with probability at least 1-2$\epsilon$.

We exploit the lower bound and let $\epsilon = \frac{1}{2t}$. Then we have

$$\ell_s^\star = \tilde{\mu}_s - \frac{\sigma^2(t + \frac{\sigma^2 \log(2t)}{t^2})}{n_t - \sigma^2}, \tag{13}$$

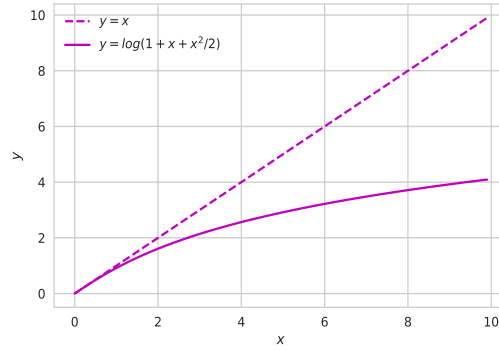where $n_t$ denotes the number of times that the example was selected in the time intervals.



Figure 1: The illustration of the influence function for the soft estimator.

Here, we provide the graph of the used influence function for the soft estimator, which explains the mechanism of the function $y = \log(1 + x + x^2/2)$ more clearly. The illustration is presented in Figure 1. As can be seen, when $x$ is large and may be an outlier, the influence function can reduce its negative impact for mean estimation. Therefore, we exploit such an influence function for robust mean estimation, which brings better classification performance.

## A.2 Proof of Theorem 2

**Lemma 1** ([11]). *Let $Z_n = \{z_1, \ldots, z_n\}$ be a (not necessarily time homogeneous) Markov chain with mean $\mu_z$, taking values in a Polish state space $\Lambda_1 \times \ldots \times \Lambda_n$, with a mixing time $\tau(\upsilon)$ (for $0 \leq \upsilon \leq 1$). Let*

$$\tau_{\min} = \inf_{0 \leq \upsilon < 1} \tau(\upsilon) \cdot \left(\frac{2 - \upsilon}{1 - \upsilon}\right)^2. \tag{14}$$

*For some $\eta \in \mathbb{R}^+$, suppose that $f : \Lambda \to \mathbb{R}$ satisfies the following inequality:*

$$f(a) - f(b) \leq \sum_{i=1}^n \eta \mathbb{1}[a_i \neq b_i], \tag{15}$$

*for every $a, b \in \Lambda$. Then for any $\epsilon \geq 0$, we have*

$$P(|f(Z_n) - \mathbb{E}f(Z_n)| \geq \epsilon) \leq 2\exp\left(\frac{-2\epsilon^2}{\eta^2 \tau_{\min}}\right). \tag{16}$$

The detailed definition of the mixing time for the Markov chain can be found in [11, 12]. Let $f$ be the mean function. Following the prior work on mean estimation [7, 4, 3, 10], without loss of generality, we assume $\mu_z = 0$ for the underlying true distribution, and $|z_i|$ is upper bounded by $Z$. Then we can set $\eta$ to $4Z/n$ for Eq. (15). Combining the above analyses, we can revise Eq. (16) as follows:

$$P\left(\left|\frac{1}{n}\sum_{i=1}^{n} z_i\right| \geq \frac{2Z}{n}\sqrt{2\tau_{\min}\log\frac{2}{\epsilon_1}}\right) \leq \epsilon_1, \tag{17}$$

and

$$P\left(\max_{i\in[n]}|z_i| \geq \frac{2Z}{n}\sqrt{2\tau_{\min}\log\frac{2n}{\epsilon_2}}\right) \leq \epsilon_2, \tag{18}$$

for $\epsilon_1 > 0$ and $\epsilon_2 > 0$. If we remove the potential outliers $Z_{n_o}$ from $Z_n$. Therefore, we have

$$
\begin{aligned}
\left|\frac{1}{n-n_o}\sum_{z_i\in Z_n\backslash Z_{n_o}} -\mu_z\right| &= \frac{1}{n-n_o}\left|\sum_{z_i\in Z_n} - \sum_{z_i\in Z_{n_o}}\right| \\
&\leq \frac{1}{n-n_o}\left(\left|\sum_{z_i\in Z_n}\right| + \left|\sum_{z_i\in Z_{n_o}}\right|\right) \\
&\leq \frac{1}{n-n_o}\left(\left|\sum_{z_i\in Z_n}\right| + n_o\max_{i\in[n]}|z_i|\right) \\
&\leq \frac{1}{n-n_o}\left(2Z\sqrt{2\tau_{\min}\log\frac{2}{\epsilon_1}} + \frac{2Zn_o}{n}\sqrt{2\tau_{\min}\log\frac{2n}{\epsilon_2}}\right),
\end{aligned}
\tag{19}
$$

which holds with probability at least $1 - \epsilon_1 - \epsilon_2$.

For our task, we exploit the concentration inequality. Let $\epsilon_1 = \epsilon_2 = \frac{1}{2t}$, and the losses be bounded by $L$. Next we can obtain

$$
\begin{aligned}
|\tilde{\mu}_h - \mu| &\leq \frac{2L}{t-t_o}\left(\sqrt{2\tau_{\min}\log(4t)} + \frac{t_o}{t}\sqrt{4\tau_{\min}\log(4t)}\right) \\
&= \frac{2\sqrt{2\tau_{\min}}L(t+\sqrt{2}t_o)}{(t-t_o)t}\sqrt{\log(4t)}
\end{aligned}
\tag{20}
$$

with the probability at least $1 - \frac{1}{t}$. In practice, it is easy to identify the value of $L$. For example, we can training deep networks on noisy datasets to observe the loss distributions. Then, we exploit the lower bound such that

$$\ell_h^\star = \tilde{\mu}_h - \frac{2\sqrt{2\tau_{\min}}L(t+\sqrt{2}t_o)}{(t-t_o)\sqrt{t}}\sqrt{\frac{\log(4t)}{n_t}} \tag{21}$$

for sample selection.

# B  Complementary Experimental Analyses

|  | # of training | # of testing | # of class | size |
|---|---|---|---|---|
| MNIST | 60,000 | 10,000 | 10 | 28×28×1 |
| F-MNIST | 60,000 | 10,000 | 10 | 28×28×1 |
| CIFAR-10 | 50,000 | 10,000 | 10 | 32×32×3 |
| CIFAR-100 | 50,000 | 10,000 | 100 | 32×32×3 |

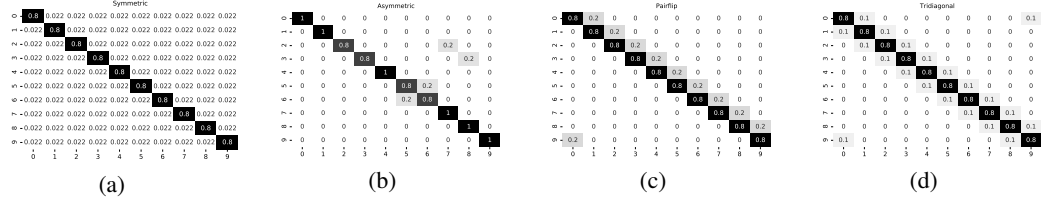Table 1: Summary of synthetic datasets used in the experiments.

Figure 2: Synthetic class-dependent transition matrices used in our experiments on *MNIST*. The noise rate is set to 20%.

## B.1 The Details of Datasets and Generating Noisy Labels

For the details of datasets, the important statistics of the used datasets are summarized in Table 1.

For the details of generating noisy labels, we exploit both *class-dependent* and *instance-dependent label noise* which include five types of synthetic label noise to verify the effectiveness of the proposed method. Here, we describe the details of the noise setting as follows:

(1). Class-dependent label noise:

• Symmetric noise: this kind of label noise is generated by flipping labels in each class uniformly to incorrect labels of other classes.

• Asymmetic noise : this kind of label noise is generated by flipping labels within a set of similar classes. In this paper, for *MNIST*, flipping 2→7, 3→8, 5↔6. For *F-MNIST*, flipping TSHIRT→SHIRT, PULLOVER→COAT, SANDALS→SNEAKER. For *CIFAR-10*, flipping TRUCK→AUTOMOBILE, BIRD→AIRPLANE, DEER→HORSE, CAT↔DOG. For *CIFAR-100*, the 100 classes are grouped into 20 super-classes, and each has 5 sub-classes. Each class is then flipped into the next within the same super-class.

• Pairflip noise: the noise flips each class to its adjacent class.

• Tridiagonal noise: the noise corresponds to a spectral of classes where adjacent classes are easier to be mutually mislabeled, unlike the unidirectional pair flipping. It can be implemented by two consecutive pair flipping transformations in the opposite direction.

(2). Instance-dependent label noise:

• Instance noise: the noise is quite realistic, where the probability that an instance is mislabeled depends on its features. We generate this type of label noise to validate the effectiveness of the proposed method as did in [13].

We use synthetic noisy *MNIST* as an example and plot the noise transition matrices in Figure 2. The noise rate is set to 20%.

## B.2 Comparison with Other Types of Baselines

As we focus on the sample selection approach in learning with noisy labels, in the main paper (Section 3.1), we fairly compare our methods with the baselines which also focus on sample selection. Here, we evaluate other types of baselines. We exploit APL [8] and CDR [14], which add implicit regularization from different perspectives. The experiments are conducted on *MNIST* and *F-MNIST*. Other experimental settings are the same as those in the main paper. The experimental results are provided in Table 2 and 3, which show that the proposed methods can outperform them with respect to classification performance.

## B.3 Experiments on Synthetic *CIFAR-100*

For *CIFAR-100*, we use a 7-layer CNN structure from [17, 16]. Other experimental settings are the same as those in the experiments on *MNIST*, *F-MNIST*, and *CIFAR-10*. The results are provided in Table 4. We can see the proposed method outperforms all the baselines.

| Noise type | Sym. | | Asym. | | Pair. | | Trid. | | Ins. | |
|---|---|---|---|---|---|---|---|---|---|---|
| Method/Noise ratio | 20% | 40% | 20% | 40% | 20% | 40% | 20% | 40% | 20% | 40% |
| APL | 98.76 | 94.92 | 98.63 | 88.65 | 98.66 | 68.44 | 98.93 | 76.44 | 97.63 | 87.90 |
| | ±0.06 | ±0.31 | ±0.05 | ±1.72 | ±0.10 | ±2.95 | ±0.04 | ±3.04 | ±0.73 | ±1.94 |
| CDR | 94.77 | 92.16 | 96.73 | 91.05 | 93.25 | 71.02 | 94.06 | 70.28 | 93.17 | 77.45 |
| | ±0.17 | ±0.73 | ±0.19 | ±0.76 | ±0.90 | ±3.89 | ±0.92 | ±4.01 | ±0.96 | ±3.04 |

Table 2: Test accuracy (%) on *MNIST* over the last ten epochs.

| Noise type | Sym. | | Asym. | | Pair. | | Trid. | | Ins. | |
|---|---|---|---|---|---|---|---|---|---|---|
| Method/Noise ratio | 20% | 40% | 20% | 40% | 20% | 40% | 20% | 40% | 20% | 40% |
| APL | 91.73 | 89.06 | 90.13 | 80.34 | 90.22 | 78.54 | 90.84 | 86.53 | 90.96 | 85.55 |
| | ±0.20 | ±0.41 | ±0.17 | ±0.63 | ±0.80 | ±4.33 | ±0.22 | ±0.76 | ±0.77 | ±2.86 |
| CDR | 85.62 | 71.83 | 89.78 | 79.05 | 85.72 | 69.07 | 86.75 | 73.63 | 85.92 | 73.14 |
| | ±0.96 | ±1.37 | ±0.41 | ±1.39 | ±0.65 | ±2.31 | ±1.19 | ±2.82 | ±1.43 | ±3.12 |

Table 3: Test accuracy on *F-MNIST* over the last ten epochs.

| Noise type | Sym. | | Asym. | | Pair. | | Trid. | | Ins. | |
|---|---|---|---|---|---|---|---|---|---|---|
| Method/Noise ratio | 20% | 40% | 20% | 40% | 20% | 40% | 20% | 40% | 20% | 40% |
| S2E | 44.59 | 25.78 | 42.18 | 26.81 | 42.99 | 26.96 | 43.16 | 27.72 | 43.13 | 27.12 |
| | ±0.32 | ±5.44 | ±1.73 | ±2.25 | ±1.54 | ±2.48 | ±0.93 | ±3.56 | ±0.67 | ±3.86 |
| MentorNet | 43.15 | 37.62 | 41.03 | 28.27 | 40.06 | 27.17 | 42.20 | 31.74 | 40.54 | 33.09 |
| | ±0.42 | ±0.89 | ±0.22 | ±0.41 | ±0.37 | ±0.92 | ±0.30 | ±0.88 | ±0.69 | ±1.53 |
| Co-teaching | 45.17 | 40.95 | 42.76 | 30.27 | 42.50 | 30.07 | 44.41 | 34.96 | 42.23 | 35.87 |
| | ±0.25 | ±0.52 | ±0.34 | ±0.33 | ±0.39 | ±0.17 | ±0.41 | ±0.35 | ±0.52 | ±1.47 |
| SIGUA | 42.03 | 40.53 | 36.67 | 26.71 | 36.48 | 26.73 | 39.21 | 32.69 | 39.19 | 33.51 |
| | ±0.33 | ±0.49 | ±0.25 | ±0.42 | ±0.37 | ±0.33 | ±0.40 | ±0.36 | ±0.32 | ±0.43 |
| JoCor | 45.93 | 41.56 | 42.89 | 29.19 | 42.12 | 30.12 | 44.98 | 34.23 | 44.28 | 35.60 |
| | ±0.21 | ±0.57 | ±0.37 | ±1.42 | ±0.35 | ±0.65 | ±0.27 | ±1.13 | ±0.59 | ±0.99 |
| CNLCU-S | **46.09** | **42.11** | **43.06** | **30.47** | **43.08** | **30.33** | **45.19** | **35.49** | **44.80** | **36.23** |
| | **±0.29** | **±0.70** | **±0.28** | **±0.37** | **±0.92** | **±0.74** | **±0.90** | **±1.30** | **±0.70** | **±0.49** |
| CNLCU-H | **46.27** | **42.05** | **43.21** | **30.55** | **43.25** | **30.79** | **45.02** | **35.24** | **45.02** | **36.17** |
| | **±0.38** | **±0.87** | **±0.93** | **±0.72** | **±0.75** | **±0.86** | **±1.06** | **±0.93** | **±1.07** | **±1.54** |

Table 4: Test accuracy (%) on *CIFAR-100* over the last ten epochs. The best two results are in bold.

## B.4 Experiments for Ablation Study

We conduct the ablation study to analyze the sensitivity of the length of time intervals. The results are shown in Figure. 3 and 4. As we can seen, the proposed method, i.e., CNLCU-S and CNLCU-H are robust to the choices of hyperparameters.
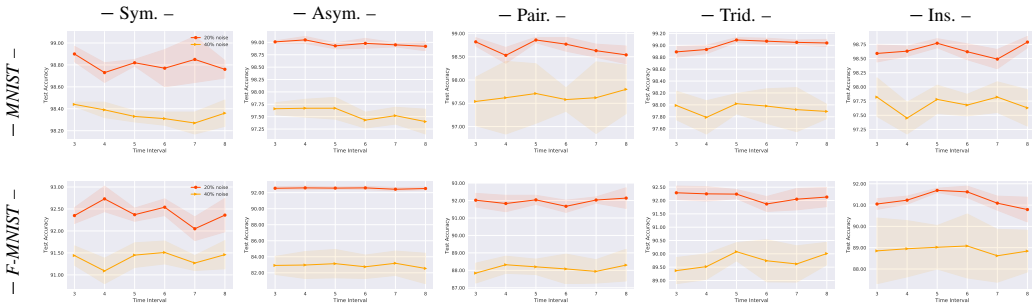


Figure 3: Illustrations of the hyperparameter sensitivity for the proposed CNLCU-S. The error bar for standard deviation in each figure has been shaded.

Note that in this paper, we concern uncertainty from *two aspects*, i.e., the uncertainty about small-loss examples and the uncertainty about large-loss examples. Here, we conduct ablation study to show the effect of removing different components to provide insights into what makes the proposed methods successful. The experiments are conducted on *MNIST* and *F-MNIST*. Other experimental settings
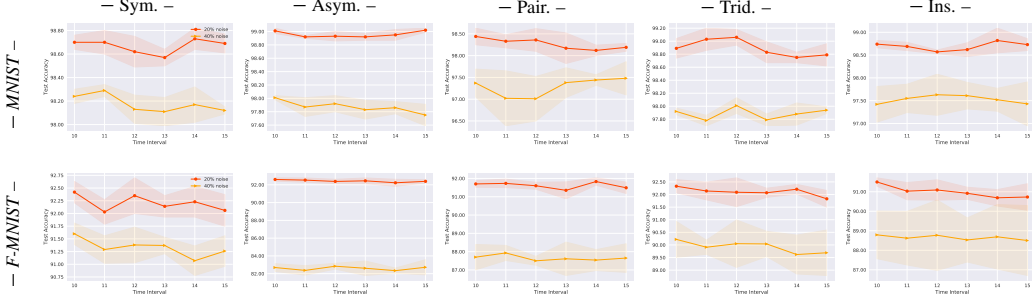
Figure 4: Illustrations of the hyperparameter sensitivity for the proposed CNLCU-H. The error bar for standard deviation in each figure has been shaded.

are the same as those in the main paper (Section 3.1). Note that we employ two networks to teach each other following [6]. Therefore, when we do not consider uncertainty in sample selection, the proposed methods will *reduce to* the baseline Co-teaching [6].

To study the effect of concerning uncertainty about small-loss examples, we remove the concerns about large-loss examples, i.e., the network is not encouraged to choose the less selected examples for updates. We express such a setting as "**with**o**ut c**oncerning about **l**arge-loss examples" (abbreviated as *w/o cl*). To study the effect of concerning uncertainty about large-loss examples, we remove the concerns about small-loss examples, i.e., we only exploit the predictions of the current network. We express such a setting as "**with**o**ut c**oncerning about **s**mall-loss examples" (abbreviated as *w/o cs*). Besides, we express the setting which directly uses non-robust mean as Co-teaching-M.

The experimental results of ablation study are provided in Table 5 and 6. As can be seen, both aspects of uncertainty concerns can improve the robustness of models. Therefore, combining two uncertainty concerns, we can better combat noisy labels. In addition, robust mean estimation is superior to the non-robust mean in learning with noisy labels.

| Noise type | Sym. | | Asym. | | Pair. | | Trid. | | Ins. | |
|---|---|---|---|---|---|---|---|---|---|---|
| Method/Noise ratio | 20% | 40% | 20% | 40% | 20% | 40% | 20% | 40% | 20% | 40% |
| 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-S *w/o cl* | 98.02 | 96.83 | 98.50 | 96.25 | 98.22 | 96.08 | 98.64 | 97.25 | 98.17 | 97.13 |
| | ±0.08 | ±0.29 | ±0.04 | ±0.13 | ±0.13 | ±0.75 | ±0.31 | ±0.24 | ±0.20 | ±0.40 |
| CNLCU-S *w/o cs* | 98.15 | 97.12 | 98.36 | 96.39 | 98.04 | 96.12 | 98.74 | 97.30 | 98.11 | 97.32 |
| | ±0.20 | ±0.22 | ±0.07 | ±0.48 | ±0.24 | ±0.68 | ±0.05 | ±0.52 | ±0.15 | ±0.43 |
| 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 |
| CNLCU-H *w/o cl* | 98.06 | 96.92 | 98.39 | 96.51 | 97.04 | 95.62 | 98.33 | 97.41 | 98.01 | 96.15 |
| | ±0.13 | ±0.23 | ±0.04 | ±0.57 | ±0.87 | ±0.93 | ±0.47 | ±0.92 | ±0.20 | ±0.28 |
| CNLCU-H *w/o cs* | 98.19 | 97.05 | 98.76 | 97.17 | 97.26 | 96.31 | 98.29 | 97.65 | 98.34 | 96.49 |
| | ±0.22 | ±0.49 | ±0.59 | ±0.60 | ±1.19 | ±0.25 | ±0.17 | ±0.92 | ±0.36 | ±0.48 |
| Co-teaching-M | 97.72 | 97.78 | 98.27 | 95.42 | 96.22 | 95.01 | 97.92 | 96.64 | 98.02 | 96.03 |
| | ±0.08 | ±0.32 | ±0.03 | ±0.42 | ±0.10 | ±0.65 | ±0.14 | ±0.77 | ±0.04 | ±0.57 |
| 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 |

Table 5: Test accuracy (%) on *MNIST* over last ten epochs.

## C  Complementary Explanation for Network Structures

Table 7 describes the 9-layer CNN [6] used on *MNIST*, *F-MNIST*, and *CIFAR-10*. Table 8 describes the 9-layer CNN [17] used on *CIFAR-100*. Here, LReLU stands for Leaky ReLU [15]. The slopes of all LReLU functions in the networks are set to 0.01. Note that that the 7/9-layer CNN is a standard and common practice in weakly supervised learning. We decided to use these CNNs, since then the experimental results are directly comparable with previous approaches in the same area, i.e., learning with noisy labels.

| Noise type | Sym. | | Asym. | | Pair. | | Trid. | | Ins. | |
|---|---|---|---|---|---|---|---|---|---|---|
| Method/Noise ratio | 20% | 40% | 20% | 40% | 20% | 40% | 20% | 40% | 20% | 40% |
| 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-S *w/o cl* | 91.77 | 89.40 | 91.25 | 72.93 | 91.53 | 87.31 | 91.31 | 89.50 | 91.09 | 88.45 |
| | ±0.35 | ±0.26 | ±0.30 | ±2.63 | ±0.17 | ±0.59 | ±0.52 | ±0.32 | ±0.13 | ±0.57 |
| CNLCU-S *w/o cs* | 91.85 | 90.76 | 91.94 | 80.99 | 91.28 | 87.31 | 91.39 | 89.29 | 90.98 | 88.73 |
| | ±0.33 | ±0.28 | ±0.09 | ±2.74 | ±0.20 | ±0.72 | ±0.07 | ±0.51 | ±0.43 | ±0.62 |
| 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 |
| CNLCU-H *w/o cl* | 91.70 | 90.05 | 91.08 | 71.35 | 91.03 | 87.22 | 91.59 | 90.01 | 90.80 | 88.31 |
| | ±0.04 | ±0.31 | ±0.06 | ±2.30 | ±0.29 | ±0.72 | ±0.07 | ±0.24 | ±0.27 | ±1.09 |
| CNLCU-H *w/o cs* | 91.82 | 90.92 | 92.45 | 80.73 | 91.21 | 87.49 | 92.08 | 89.72 | 91.21 | 88.62 |
| | ±0.13 | ±0.42 | ±0.25 | ±1.63 | ±0.17 | ±0.32 | ±0.13 | ±0.24 | ±0.38 | ±0.73 |
| Co-teaching-M | 91.33 | 89.05 | 91.14 | 71.03 | 90.85 | 86.95 | 91.50 | 89.18 | 90.74 | 88.25 |
| | ±0.18 | ±0.73 | ±0.90 | ±3.73 | ±0.61 | ±0.19 | ±0.46 | ±0.44 | ±1.06 | ±0.92 |
| 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 |

Table 6: Test accuracy (%) on *F-MNIST* over last ten epochs.

Table 7: CNN on *MNIST*, *F-MNIST*, and *CIFAR-10*.

| CNN on *MNIST* | CNN on *F-MNIST* | CNN on *CIFAR-10* |
|---|---|---|
| 28×28 Gray Image | 28×28 Gray Image | 32×32 RGB Image |
| 3×3 conv, 128 LReLU | | |
| 3×3 conv, 128 LReLU | | |
| 3×3 conv, 128 LReLU | | |
| 2×2 max-pool | | |
| dropout, $p = 0.25$ | | |
| 3×3 conv, 256 LReLU | | |
| 3×3 conv, 256 LReLU | | |
| 3×3 conv, 256 LReLU | | |
| 2×2 max-pool | | |
| dropout, $p = 0.25$ | | |
| 3×3 conv, 512 LReLU | | |
| 3×3 conv, 256 LReLU | | |
| 3×3 conv, 128 LReLU | | |
| avg-pool | | |
| dense 128→10 | dense 128→10 | dense 128→10 |

Table 8: CNN on *CIFAR-100*.

| CNN on *CIFAR-100* |
|---|
| 32×32 RGB Image |
| 3×3 conv, 64 ReLU |
| 3×3 conv, 64 ReLU |
| 2×2 max-pool |
| 3×3 conv, 128 ReLU |
| 3×3 conv, 128 ReLU |
| 2×2 max-pool |
| 3×3 conv, 196 ReLU |
| 3×3 conv, 196 ReLU |
| 2×2 max-pool |
| dense 256→100 |

# References

[1] 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.

[2] Peng Chen, Xinghu Jin, Xiang Li, and Lihu Xu. A generalized catoni's $m$-estimator under finite $\alpha$-th moment assumption with $\alpha \in (1, 2)$. *arXiv preprint arXiv:2010.05008*, 2020.

[3] Ilias Diakonikolas and Daniel M Kane. Recent advances in algorithmic high-dimensional robust statistics. *arXiv preprint arXiv:1911.05911*, 2019.

[4] Ilias Diakonikolas, Daniel M Kane, and Ankit Pensia. Outlier robust mean estimation with subgaussian rates via stability. *arXiv preprint arXiv:2007.15618*, 2020.

[5] Evarist Giné, Rafał Latała, and Joel Zinn. Exponential and moment inequalities for u-statistics. In *High Dimensional Probability II*, pages 13–38. Springer, 2000.

[6] 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.

[7] Liu Liu, Tianyang Li, and Constantine Caramanis. High dimensional robust m-estimation: Arbitrary corruption and heavy tails. *arXiv preprint arXiv:1901.08237*, 2019.

[8] 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.

[9] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. MIT Press, 2018.

[10] Laura Niss and Ambuj Tewari. What you see may not be what you get: Ucb bandit algorithms robust to $\epsilon$-contamination. In *UAI*, pages 450–459, 2020.

[11] Daniel Paulin et al. Concentration inequalities for markov chains by marton couplings and spectral methods. *Electronic Journal of Probability*, 20, 2015.

[12] Gareth O Roberts, Jeffrey S Rosenthal, et al. General state space markov chains and mcmc algorithms. *Probability surveys*, 1:20–71, 2004.

[13] 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.

[14] 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.

[15] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015.

[16] 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.

[17] Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor W Tsang, and Masashi Sugiyama. How does disagreement benefit co-teaching? In *ICML*, 2019.