
FINE Samples for Learning with Noisy Labels

Taehyeon Kim *
KAIST AI
KAIST
Daejeon, South Korea
potter32@kaist.ac.kr

Jongwoo Ko *
KAIST AI
KAIST
Daejeon, South Korea
jongwoo.ko@kaist.ac.kr

Sangwook Cho
KAIST AI
KAIST
Daejeon, South Korea
sangwookcho@kaist.ac.kr

Jinhwan Choi
KAIST AI
KAIST
Daejeon, South Korea
jinhwanchoi@kaist.ac.kr

Se-Young Yun
KAIST AI
KAIST
Daejeon, South Korea
yunseyoung@kaist.ac.kr

Abstract

Modern deep neural networks (DNNs) become weak when the datasets contain noisy (incorrect) class labels. Robust techniques in the presence of noisy labels can be categorized into two types: developing *noise-robust* functions or using *noise-cleansing* methods by detecting the noisy data. Recently, *noise-cleansing* methods have been considered as the most competitive noisy-label learning algorithms. Despite their success, their noisy label detectors are often based on heuristics more than a theory, requiring a robust classifier to predict the noisy data with loss values. In this paper, we propose a novel detector for filtering label noise. Unlike most existing methods, we focus on each data point’s latent representation dynamics and measure the alignment between the latent distribution and each representation using the eigen decomposition of the data gram matrix. Our framework, coined as *filtering noisy instances via their eigenvectors* (FINE), provides a robust detector using derivative-free simple methods with theoretical guarantees. Under our framework, we propose three applications of the FINE: sample-selection approach, semi-supervised learning (SSL) approach, and collaboration with *noise-robust* loss functions. Experimental results show that the proposed methods consistently outperform corresponding baselines for all three applications on various benchmark datasets¹.

1 Introduction

Deep neural networks (DNNs) have achieved remarkable success in numerous tasks as the amount of accessible data has dramatically increased [21, 15]. On the other hand, accumulated datasets are typically labeled by a human, a labor-intensive job or through web crawling [48] so that they may be easily corrupted (*label noise*) in real-world situations. Recent studies have shown that deep neural networks have the capacity to memorize essentially any labeling of the data [49]. Even a small amount of such noisy data can hinder the generalization of DNNs owing to their strong memorization of noisy labels [49, 29]. Hence, it becomes crucial to train DNNs that are robust to corrupted labels. As label noise problems may appear anywhere, such robustness increases reliability in many applications such as the e-commerce market [9], medical fields [45], on-device AI [46], and autonomous driving systems [11].

To improve the robustness against noisy data, the methods for learning with noisy labels (LNL) have been evolving in two main directions [18]: (1) designing noise-robust objective functions or regular-

*Equal contribution

¹Code available at https://github.com/Kthyeon/FINE_official

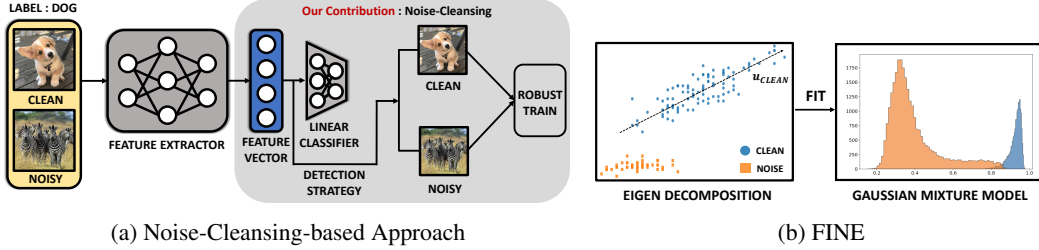


Figure 1: Illustration of (a) basic concept of this work and (b) proposed detection framework, FINE. Noise-cleansing learning generally separates clean data from the original dataset by using prediction outputs. We propose a novel derivative-free detector based on an unsupervised clustering algorithm on the high-order topological space. FINE measures the alignment of pre-logits (i.e., penultimate layer representation vectors) toward the class-representative vector that is extracted through the eigen decomposition of the gram matrix of data representations.

izations and (2) detecting and cleansing the noisy data. In general, the former *noise-robust* direction uses explicit regularization techniques [6, 52, 50] or robust loss functions [38, 13, 40, 51], but their performance is far from state-of-the-art [49, 26] on datasets with severe noise rates. Recently, researchers have designed *noise-cleansing* algorithms focused on segregating the clean data (i.e., samples with uncorrupted labels) from the corrupted data [19, 14, 47, 18, 32, 42]. One of the popular criteria for the segregation process is the loss value between the prediction of the noisy classifier and its noisy label, where it is generally assumed that the noisy data have a large loss [19, 14, 47, 18] or the magnitude of the gradient during training [51, 40]. However, these methods may still be biased by the corrupted linear classifier towards label noise because their criterion (e.g., loss values or weight gradient) uses the posterior information of such a linear classifier [24]. Maennel et al. [31] analytically showed that the principal components of the weights of a neural network align with the randomly labeled data; this phenomenon can yield more negative effects on the classifier as the number of randomly labeled classes increases. Recently, Wu et al. [42] used an inherent geometric structure induced by nearest neighbors (NN) in latent space and filtered out isolated data in such topology, and its quality was sensitive to its hyperparameters regarding NN clustering in the presence of severe noise rates.

To mitigate such issues for label noise detectors, we provide a novel yet simple detector framework, *filtering noisy labels via their eigenvectors* (FINE) with theoretical guarantees to provide a high-quality splitting of clean and corrupted examples (without the need to estimate noise rates). Instead of using the neural network’s linear classifier, FINE utilizes the principal components of latent representations made by eigen decomposition which is one of the most widely used unsupervised learning algorithms and separates clean data and noisy data by these components (Figure 1a). To motivate our approach, as Figure 1b shows, we find that the clean data (blue points) are mainly aligned on the principal component (black dotted line), whereas the noisy data (orange points) are not; thus, the dataset is well clustered with the alignment of representations toward the principal component by fitting them into Gaussian mixture models (GMM). We apply our framework to various *LNL* methods: the sample selection approach, a semi-supervised learning (SSL) approach, and collaboration with noise-robust loss functions. The key contributions of this work are summarized as follows:

- We propose a novel framework, termed FINE (*filtering noisy labels via their eigenvectors*), for detecting clean instances from noisy datasets. FINE makes robust decision boundary for the high-order topological information of data in latent space by using eigen decomposition of their gram matrix.
- We provide provable evidence that FINE allows a meaningful decision boundary made by eigenvectors in latent space. We support our theoretical analysis with various experimental results regarding the characteristics of the principal components extracted by our FINE detector.
- We develop a simple sample-selection method by replacing the existing detector method with FINE. We empirically validate that a sample-selection learning with FINE provides consistently superior detection quality and higher test accuracy than other existing alternative methods such as the Co-teaching family [14, 47], TopoFilter [42], and CRUST [32].

- We experimentally show that our detection framework can be applied in various ways to existing *LNL* methods and validate that ours consistently improves the generalization in the presence of noisy data: sample-selection approach [14, 47], SSL approach [25], and collaboration with noise-robust loss functions [51, 40, 29].

Organization. The remainder of this paper is organized as follows. In Section 2, we discuss the recent literature on LNL solutions and meaningful detectors. In Section 3, we address our motivation for creating a noisy label detector with theoretical insights and provide our main method, filtering the noisy labels via their eigenvectors (FINE). In Section 4, we present the experimental results. Finally, Section 5 concludes the paper.

2 Related Works

Zhang et al. [49] empirically showed that any convolutional neural networks trained using stochastic gradient methods easily fit a random labeling of the training data. To tackle this issue, numerous works have examined the classification task with noisy labels. We do not consider the works that assumed the availability of small subsets of training data with clean labels [17, 36, 39, 53, 3].

Noise-Cleansing-based Approaches. Noise-cleansing methods have evolved following the improvement of noisy detectors. Han et al. [14] suggested a noisy detection approach, named co-teaching, that utilizes two networks, extracts subsets of instances with small losses from each network, and trains each network with subsets of instances filtered by another network. Yu et al. [47] combined a disagreement training procedure with co-teaching, which only selects instances predicted differently by two networks. Huang et al. [18] provided a simple noise-cleansing framework, training-filtering-training; the empirical efficacy was improved by first finding label errors, then training the model only on data predicted as clean. Recently, new noisy detectors with theoretical support have been developed. Wu et al. [42] proposed a method called TopoFilter that filters noisy data by utilizing the k-nearest neighborhood algorithm and Euclidean distance between pre-logits. Mirzasoleiman et al. [32] introduced an algorithm that selects subsets of clean instances that provide an approximately low-rank Jacobian matrix and proved that gradient descent applied to the subsets prevents overfitting to noisy labels. Pleiss et al. [34] proposed an area under margin (AUM) statistic that measures the average difference between the logit values of the assigned class and its highest non-assigned class to divide clean and noisy samples. Cheng et al. [8] progressively filtered out corrupted instances using a novel confidence regularization term. The noise-cleansing method was also developed in a semi-supervised learning (SSL) manner. Li et al. [25] modeled the per-sample loss distribution and divide it into a labeled set with clean samples and an unlabeled set with noisy samples, and they leverage the noisy samples through the well-known SSL technique MixMatch [4].

Noise-Robust Models. Noise-robust models have been studied in the following directions: robust-loss functions, regularizations, and strategies. First, for robust-loss functions, Ghosh et al. [13] showed that the mean absolute error (MAE) might be robust against noisy labels. Zhang & Sabuncu et al. [51] argued that MAE performed poorly with DNNs and proposed a GCE loss function, which can be seen as a generalization of MAE and cross-entropy (CE). Wang et al. [40] introduced the reverse version of the cross-entropy term (RCE) and suggested that the SCE loss function is a weighted sum of the CE and RCE. Some studies have stated that the early-stopped model can prevent the memorization phenomenon for noisy labels [2, 49] and theoretically analyzed it [26]. Based on this hypothesis, Liu et al. [29] proposed an early-learning regularization (ELR) loss function to prohibit memorizing noisy data by leveraging the semi-supervised learning techniques. Xia et al. [43] clarified which neural network parameters cause memorization and proposed a robust training strategy for these parameters. Efforts have been made to develop regularizations on the prediction level by smoothing the one-hot vector [30], using linear interpolation between data instances [50], and distilling the rescaled prediction of other models [20]. However, these works have limitations in terms of performance as the noise rate of the dataset increases.

Dataset Resampling. Label-noise detection may be a category of data resampling which is a common technique in the machine learning community that extracts a “*helpful*” dataset from the distribution of the original dataset to remove the dataset bias. In class-imbalance tasks, numerous studies have conducted over-sampling of minority classes [7, 1] or undersampling the majority classes [5] to balance the amount of data per class. Li & Vasconcelos et al. [27] proposed a resampling procedure to reduce the representation bias of the data by learning a weight distribution that favors difficult instances for a given feature representation. Le Bras et al. [22] suggested an adversarial filtering-based approach to remove spurious artifacts in a dataset. Analogously, in anomaly detection and

Algorithm 1: FINE Algorithm for Sample Selection

INPUT : Noisy training data \mathcal{D} , feature extractor g , number of classes K , clean probability threshold ζ , set of FINE scores for class k \mathcal{F}_k

OUTPUT : Collected clean data \mathcal{C}

- 1: Initialize $\mathcal{C} \leftarrow \emptyset$, $\hat{\mathcal{D}} \leftarrow \mathcal{D}$, $\Sigma_k \leftarrow \mathbf{0}$ for all $k = 1, \dots, K$
/* Update the covariance matrices for all classes */
 - 2: **for** $(\mathbf{x}_i, y_i) \in \mathcal{D}$ **do**
 - 3: $\mathbf{z}_i \leftarrow g(\mathbf{x}_i)$
 - 4: Update the gram matrix $\Sigma_{y_i} \leftarrow \Sigma_{y_i} + \mathbf{z}_i \mathbf{z}_i^\top$
 - 5: **end for**
/* Generate the principal component with eigen decomposition */
 - 6: **for** $k = 1, \dots, K$ **do**
 - 7: $\mathbf{U}_k, \Lambda_k \leftarrow$ EIGEN DECOMPOSITION OF Σ_k
 - 8: $\mathbf{u}_k \leftarrow$ THE FIRST COLUMN OF \mathbf{U}_k
 - 9: **end for**
/* Compute the alignment score and get clean subset \mathcal{C} */
 - 10: **for** $(\mathbf{x}_i, y_i) \in \mathcal{D}$ **do**
 - 11: Compute the FINE score $f_i = \langle \mathbf{u}_{y_i}, \mathbf{z}_i \rangle^2$ and $\mathcal{F}_{y_i} \leftarrow \mathcal{F}_{y_i} \cup \{f_i\}$
 - 12: **end for**
/* Finding the samples whose clean probability is larger than ζ */
 - 13: $\mathcal{C} \leftarrow \mathcal{C} \cup \text{GMM}(\mathcal{F}_k, \zeta)$ for all $k = 1, \dots, K$
-

out-of-distribution detection problems [16, 28, 23], the malicious data are usually detected by examining the loss value or negative behavior in the feature representation space. While our research is motivated by such previous works, this paper focuses on the noisy image classification task.

3 Method

In this section, we present our detector framework and the theoretical motivation behind using the detector in high-dimensional classification. To segregate the clean data, we utilize the degree of alignment between the representations and the eigenvector of the representations' gram matrices for all classes, called **FINE** (*F*iltering *N*oisy *i*nstances *v*ia *t*heir *E*igenvectors). Our algorithm is as follows (Algorithm 1). FINE first creates a gram matrix of the representation in the noisy training dataset for each class and conducts the eigen decomposition for those gram matrices. Then, FINE finds clean and noisy instances using the square of inner product values between the representations and the first eigenvector having the largest eigenvalue. In this manner, we treat the data as clean if aligned onto the first eigenvector, while most of the noisy instances are not. Here, we formally define 'alignment' and 'alignment clusterability' in Definition 1 and Definition 2, respectively.

Definition 1. (*Alignment*) D

Definition 2. (*Alignment Clusterability*) For all features labeled as class k in dataset \mathcal{D} , let fit a Gaussian Mixture Model (GMM) on their alignment (Definition 1) distribution to divide current samples into a clean set and a noisy set; the set having larger mean value is treated as a clean set, and another one is a noisy set. Then, we say a dataset \mathcal{D} satisfies alignment clusterability if the representation \mathbf{z} labeled as the same true class belongs to the clean set.

As an empirical evaluation, the quality of our detector for noisy data is measured with the **F-score**, a widely used criterion in noisy label detection, anomaly detection and out-of-distribution detection [8, 16, 28, 23]. We treat the selected clean samples as the positive class and the noisy samples as negative class. The **F-score** is the harmonic mean of the precision and the recall; the *precision* indicates the fraction of clean samples among all samples that are predicted as clean, and the *recall* indicates the portion of clean samples that are identified correctly.

3.1 Alignment Analysis for Noisy Label Detector

To design a robust label noise filtering framework, we explore the linear nature of the topological space of feature vectors for data resampling techniques and deal with the classifier contamination due to random labels. Recent studies on the distribution of latent representations in DNNs provide

insight regarding how correctly the outlier samples can be filtered with the hidden space’s geometrical information. For instance, in [23, 24], the authors proposed frameworks for novelty detection using the topological information of pre-logit based on the Mahalanobis distance, and, in [42], the authors filtered the noisy data based on the Euclidean distance between pre-logits. Maennel et al. [31] analytically showed that an alignment between the principal components of network parameters and those of data takes place when training with random labels. This finding points out that random labels can corrupt a classifier, and thus building a robust classifier is required.

Motivated by these works, we aim to design a novel detector using the principal components of latent features to satisfy Definition 2. However, it is intractable to find the optimal classifier to maximize the separation of *alignment clusterability* because clean data distribution and noisy data distribution are inaccessible. To handle this issue, we attempt to approximate the clean eigenvector to maximize the alignment values of clean data rather than to maximize the separation; the algorithm utilizes the eigenvector of the data for each class (Figure 2). Below, we provide the upper bound for the perturbation toward the clean data’s eigenvector under simple problem settings with noisy labels referred to in other studies [29, 42]. We first introduce notations. Next, we establish the theoretical evidence that our FINE algorithm approximates the clean data’s eigenvectors under some assumptions for its analytic tractability (Theorem 1). We mainly present the theorem and its interpretation; details of the proofs can be found in the Appendix.

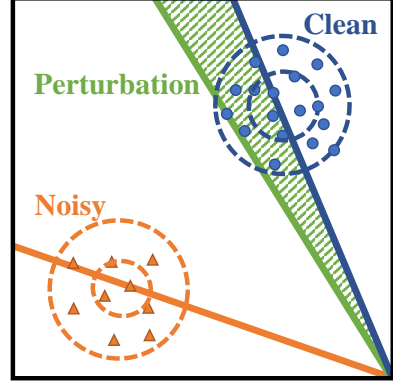


Figure 2: Illustration for the problem settings and Theorem 1. The perturbation (green shade) is the angle between the first eigenvector of clean instances (blue line) and the estimated first eigenvector (green line) which is perturbed by that of noisy instances (orange line). Note that blue and orange points are clean instances and noisy instances, respectively.

Notations. Consider a binary classification task. Assume that the data points and labels lie in $\mathcal{X} \times \mathcal{Y}$, where the feature space $\mathcal{X} \subset \mathbb{R}^d$ and label space $\mathcal{Y} = \{-1, +1\}$. A single data point \mathbf{x} and its true label y follow a distribution $(\mathbf{x}, y) \sim P_{\mathcal{X} \times \mathcal{Y}}$. Denote by \tilde{y} the observed label (potentially corrupted). Without loss of generality, we focus on the set of data points whose observed label is $\tilde{y} = +1$.

Let $\mathbf{X} \subset \mathcal{X}$ be the finite set of features with clean instances whose true label is $y = +1$. Similarly, let $\tilde{\mathbf{X}} \subset \mathcal{X}$ be the set of noisy instances whose true label is $y = -1$. To establish our theorem, we assume the following reasonable conditions referred to other works using linear discriminant analysis (LDA) assumptions [24, 12]:

Assumption 1. *The feature distribution is comprised of two Gaussians, each identified as a clean cluster and a noisy cluster.*

Assumption 2. *The features of all instances with $y = +1$ are aligned on the unit vector \mathbf{v} with the white noise, i.e., $\mathbb{E}_{\mathbf{x} \in \mathbf{X}} [\mathbf{x}] = \mathbf{v}$. Similarly, features of all instances with $y = -1$ are aligned on the unit vector \mathbf{w} , i.e., $\mathbb{E}_{\mathbf{x} \in \tilde{\mathbf{X}}} [\mathbf{x}] = \mathbf{w}$.*

Theorem 1. (Upper bound for the perturbation towards the clean data’s eigenvector \mathbf{v}) Let N_+ and N_- be the number of clean instances and noisy instances, respectively, and \mathbf{u} be the FINE’s eigenvector which is the first column of \mathbf{U} from the eigen decomposition of the whole data’s matrix Σ . For any $\delta \in (0, 1)$, its perturbation towards the \mathbf{v} in assumption 2 (i.e., 2-norm for difference of projection matrices; left hand side of Eq. (1)) holds the following with probability $1 - \delta$:

$$\|\mathbf{u}\mathbf{u}^\top - \mathbf{v}\mathbf{v}^\top\|_2 \leq \frac{3\tau \cos \theta + \mathcal{O}(\sigma^2 \sqrt{\frac{d+\log(4/\delta)}{N_+}})}{1 - \tau(\sin \theta + 3 \cos \theta) - \mathcal{O}(\sigma^2 \sqrt{\frac{d+\log(4/\delta)}{N_+}})} \quad (1)$$

where \mathbf{w} is the first eigenvector of noisy instances, τ is the fraction between noisy and clean instances ($\frac{N_-}{N_+}$), θ is $\angle(\mathbf{w}, \mathbf{v})$, and σ^2 is a variance of white noise.

Theorem 1 states that the upper bound for the perturbation toward \mathbf{v} are dependent on both the ratio τ and the angle θ between \mathbf{w} and \mathbf{v} ; small upper bound can be guaranteed as the number of clean data

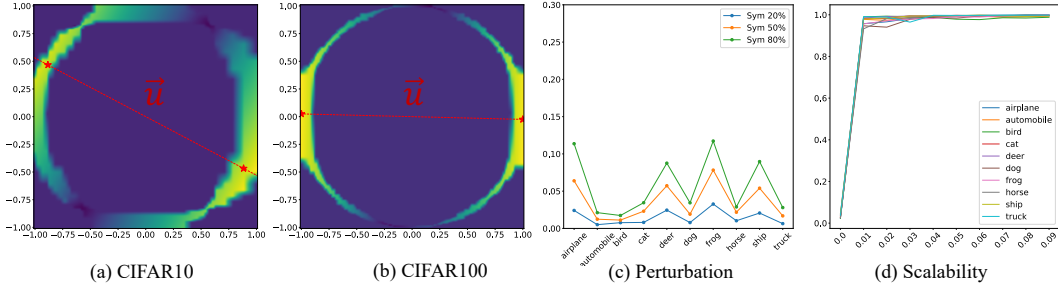


Figure 3: (a), (b): Heatmaps of Eq. (2) values on unit circle in random hyperplane. We evaluate this visualization on the ResNet34 model trained with common cross-entropy loss on CIFAR-10 with asymmetric noise 40% and CIFAR-100 with symmetric noise 80%, respectively. Colors closer to yellow indicate larger the values; (c): comparison of perturbations of Eq. (1) on CIFAR-10 with symmetric noise 20%; (d): comparison of cosine similarity values between FINE’s principal components and approximated principal components using fraction of data on CIFAR-10 with symmetric noise 80%.

increases, τ decreases, and θ approaches $\frac{\pi}{2}$. We also derive the lower bound for the precision and the recall when using the eigenvector \mathbf{u} in Appendix. In this theoretical analysis, we can ensure that such lower bound values become larger as \mathbf{v} and \mathbf{w} become orthogonal to each other. To verify these assumptions, we provide various experimental results for the separation of *alignment clusterability*, the perturbation values, the scalability to the number of samples, the quality of our detector in the application of sample selection approach, and the comparison with an alternative clustering-based estimator [24].

Validation for our Estimated Eigenvector. To validate our FINE’s principal components, we first propose a simple visualization scheme based on the following steps: (1) Pick the first eigenvector (\mathbf{u}) extracted by FINE algorithm, (2) Generate a random hyperplane spanned by such eigenvector (\mathbf{u}) and a random vector, (3) Calculate the value of the following Eq. (2) on any unit vectors (\mathbf{a}) in such hyperplane and plot a heatmap with them:

$$\frac{1}{|\mathbf{X}|} \sum_{\mathbf{x}_i \in \mathbf{X}} \langle \mathbf{a}, \mathbf{x}_i \rangle^2 - \frac{1}{|\tilde{\mathbf{X}}|} \sum_{\mathbf{x}_j \in \tilde{\mathbf{X}}} \langle \mathbf{a}, \mathbf{x}_j \rangle^2 \quad (2)$$

Eq. (2) is maximized when the unit vector \mathbf{a} not only maximizes the FINE scores of clean data for the first term in Eq. (2), but also minimizes those of noisy data for the second term in Eq. (2). This visualization shows in 2-D how FINE’s first eigenvector (\mathbf{u}) optimizes such values in the presence of noisy instances (Figure 3a and 3b). As the figures show, the FINE’s eigenvector \mathbf{u} (red dotted line) has almost maximum value of Eq. (2). Furthermore, we empirically evaluate the perturbation values in Theorem 1 as the noise rate changes (Figure 3c); FINE has small perturbation values even in a severe noise rate.

Scalability to Number of Samples. Despite FINE’s superiority, it may require high computational costs if the whole dataset is used for eigen decomposition. To address this issue, we approximate the eigenvector with a small portion of the dataset and measure the cosine similarity values between the approximated term and the original one (\mathbf{u}) (Figure 3d). Interestingly, we verify that far accurate eigenvector is computed even using 1% data (i.e., a cosine similarity value is 0.99), and thus the eigenvector can be accurately estimated with little computation time.

Validation for Dynamics of Sample-selection Approach. We evaluate the F-score dynamics of every training epoch on the symmetric and the asymmetric label noise in Figure 4. We compare FINE with the following sample-selection approaches: Co-teaching [14] and TopoFilter [42]. In Figure 4, during the training process, F-scores of FINE becomes consistently higher on both symmetric noise and asymmetric noise settings, while Co-teaching and TopoFilter achieve lower quality. Unlike TopoFilter and FINE, Co-teaching even performs the sample-selection with the access of noise rate. This evidence show that FINE is also applicable to the naive sample-selection approaches.

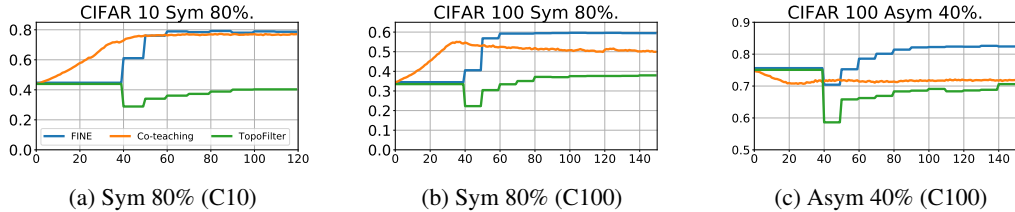


Figure 4: Comparisons of F-scores on CIFAR10 and CIFAR100 under symmetric and asymmetric label noise. C10 and C100 denote CIFAR-10 and CIFAR-100, respectively.

Comparison for Mahalanobis Distance Estimator Under similar conditions, Lee et al. [24] measured the Mahalanobis distance of pre-logits using the minimum covariance determinant (MCD) estimator and selected clean samples based on this distance. While they also utilized the LDA assumptions on pre-logits, FINE consistently outperforms MCD in both precision and recall, thus yielding better F-score (Figure 5). The experimental results justify our proposed detector, in comparison with a similar alternative.

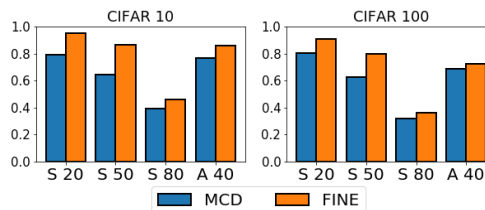


Figure 5: Comparisons of F-scores on CIFAR-10 and CIFAR-100 under symmetric (S) and asymmetric noise (A) settings.

4 Experiment

In this section, we demonstrate the effectiveness of our FINE detector for three applications: sample selection approach, SSL, and collaboration with noise-robust loss functions.

4.1 Experimental Settings

Noisy Benchmark Dataset. Following the previous setups [25, 29], we artificially generate two types of random noisy labels: injecting uniform randomness into a fraction of labels (symmetric) and corrupting a label only to a specific class (asymmetric). For example, we generate noise by mapping TRUCK \rightarrow AUTOMOBILE, BIRD \rightarrow AIRPLANE, DEER \rightarrow HORSE, CAT \leftrightarrow DOG to make asymmetric noise for CIFAR-10. For CIFAR-100, we create 20 five-size super-classes and generate asymmetric noise by changing each class to the next class within super-classes circularly. For a real-world dataset, Clothing1M [44] containing inherent noisy labels is used. This dataset contains 1 million clothing images obtained from online shopping websites with 14 classes². The dataset provides 50k, 14k, and 10k verified as clean data for training, validation, and testing. Instead of using the 50k clean training data, we use a randomly sampled pseudo-balanced subset as a training set with 120k images. For evaluation, we compute the classification accuracy on the 10k clean dataset.

Networks and Hyperparameter Settings. We use the architectures and hyperparameter settings for all baseline experiments following the setup of Liu et al. [29] except with SSL approaches. For SSL approaches, we follow the setup of Li et al. [25]. We set the threshold ζ as 0.5.

4.2 Application of FINE

4.2.1 Sample Selection-Based Approaches

We apply our FINE detector for various sample selection algorithms. In detail, after warmup training, at every epoch, FINE selects the clean data with the eigenvectors generated from the gram matrices of data predicted to be clean in the previous round, and then the neural networks are trained with them. We compare our proposed method with the following sample selection approaches: (1) *Bootstrap* [35], (2) *Forward* [33], (3) *Co-teaching* [14]; (4) *Co-teaching+* [47]; (5) *TopoFilter* [42]; (6) *CRUST* [32]. We evaluate these algorithms three times and report error bars.

²T-shirt, Shirt, Knitwear, Chiffon, Sweater, Hoodie, Windbreaker, Jacket, Down Coat, Suit, Shawl, Dress, Vest, and Underwear. The labels in this dataset are extremely noisy (estimated noisy level is 38.5%) [37].

Table 1: Test accuracies (%) on CIFAR-10 and CIFAR-100 under different noisy types and fractions. All comparison methods are reproduced with publicly available code, while the results for Bootstrap [35] and Forward [33] are taken from [29]. For CRUST [32], we experiment without mix-up to compare the intrinsic sample selection effect of each method. The average accuracies and standard deviations over three trials are reported. Here, we substitute the sample selection method of Co-teaching [14, 47] with FINE (i.e., F-Co-teaching). The best results sharing the noisy fraction and method are highlighted in bold.

Dataset	CIFAR-10				CIFAR-100			
	Noisy Type	Sym	Asym		Sym	Asym		
Noise Ratio	20	50	80	40	20	50	80	40
Standard	87.0 ± 0.1	78.2 ± 0.8	53.8 ± 1.0	85.0 ± 0.0	58.7 ± 0.3	42.5 ± 0.3	18.1 ± 0.8	42.7 ± 0.6
Bootstrap [35]	86.2 ± 0.2	-	54.1 ± 1.3	81.2 ± 1.5	58.3 ± 0.2	-	21.6 ± 1.0	45.1 ± 0.6
Forward [33]	88.0 ± 0.4	-	54.6 ± 0.4	83.6 ± 0.6	39.2 ± 2.6	-	9.0 ± 0.6	34.4 ± 1.9
Co-teaching [14]	89.3 ± 0.3	83.3 ± 0.6	66.3 ± 1.5	88.4 ± 2.8	63.4 ± 0.0	49.1 ± 0.4	20.5 ± 1.3	47.7 ± 1.2
Co-teaching+ [47]	89.1 ± 0.5	84.9 ± 0.4	63.8 ± 2.3	86.5 ± 1.2	59.2 ± 0.4	47.1 ± 0.3	20.2 ± 0.9	44.7 ± 0.6
TopoFilter [42]	90.4 ± 0.2	86.8 ± 0.3	46.8 ± 1.0	87.5 ± 0.4	66.9 ± 0.4	53.4 ± 1.8	18.3 ± 1.7	56.6 ± 0.5
CRUST [32]	89.4 ± 0.2	87.0 ± 0.1	64.8 ± 1.5	82.4 ± 0.0	69.3 ± 0.2	62.3 ± 0.2	21.7 ± 0.7	56.1 ± 0.5
FINE	91.0 ± 0.1	87.3 ± 0.2	69.4 ± 1.1	89.5 ± 0.1	70.3 ± 0.2	64.2 ± 0.5	25.6 ± 1.2	61.7 ± 1.0
F-Coteaching	92.0 ± 0.1	87.5 ± 0.1	74.2 ± 0.8	90.5 ± 0.2	71.1 ± 0.2	64.7 ± 0.3	31.6 ± 1.0	64.8 ± 0.7

Table 1 summarizes the performances of different sample selection approaches on various noise distribution and datasets. We observe that our FINE method consistently outperforms the competitive methods over the various noise rates. Our FINE methods can filter the clean instances without losing essential information, leading to training the robust network.

To go further, we improve the performance of Co-teaching [14] by substituting its sample selection state with our FINE algorithm. To combine FINE and the Co-teaching family, unlike the original methods that utilize the small loss instances to train with clean labels, we train one model with extracted samples by conducting FINE on another model. The results of the experiments are shown in the eighth and ninth rows of Table 1.

4.2.2 SSL-Based Approaches

SSL approaches [25, 10, 41] divide the training data into clean instances as labeled instances and noisy instances as unlabeled instances and use both the labeled and unlabeled samples to train the networks in SSL. Recently, methods belonging to this category have shown the best performance among the various *LNL* methods, and these methods can train robust networks for even extremely high noise rates. We compare the performances of the existing semi-supervised approaches and that in which the sample selection state of DivideMix [25] is substituted with our FINE algorithm (i.e., F-DivideMix). The results of the experiments are shown in Table 2. We achieve consistently higher performance than DivideMix by utilizing FINE instead of its loss-based filtering method and show comparable performance to the state-of-the-art SSL methods such as DST [41] and LongReMix [10]. Interestingly, as Figure 6 shows, clean and noisy data are well classified in F-DivideMix under extreme noise cases.

4.2.3 Collaboration with Noise-Robust Loss Functions

The goal of the noise-robust loss function is to achieve a small risk for unseen clean data even when noisy labels exist in the training data. There have been few collaboration studies of the noise-robust loss function methodology and dynamic sample selection. Most studies have selected clean and noisy data based on cross-entropy loss.

Here, we state the collaboration effects of FINE with various noise-robust loss functions: generalized cross entropy (GCE) [51], symmetric cross entropy (SCE) [40], and early-learning regularization (ELR) [29]. Figure 7 shows that FINE facilitates generalization in the application of noise-robust

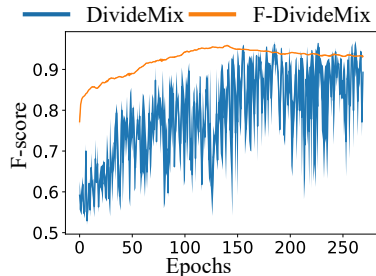


Figure 6: Comparisons of F-scores on CIFAR-10 under symmetric 90% noise. Blue line indicates the error bar of two networks’ F-score used in Dividemix [25], and Orange line indicates those replaced by our FINE detector.

Table 2: Comparison of test accuracies (%) for FINE collaborating with DivideMix and existing semi-supervised approaches on CIFAR-10 and CIFAR-100 under different noisy types and fractions. The results for all comparison methods are taken from their original works.

Dataset		CIFAR-10					CIFAR-100			
Noisy Type		Sym			Asym		Sym			
Noise Ratio		20	50	80	90	40	20	50	80	90
DivideMix [25]	Best	96.1	94.6	93.2	76.0	93.4	77.3	74.6	60.2	31.5
	Last	95.7	94.4	92.9	75.4	92.1	76.9	74.2	59.6	31.0
DST [41]	Best	96.1	95.2	92.9	-	94.3	78.0	74.3	57.8	-
	Last	95.9	94.7	92.6	-	92.3	77.4	73.6	55.3	-
LongReMix [10]	Best	96.2	95.0	93.9	82.0	94.7	77.8	75.6	62.9	33.8
	Last	96.0	94.7	93.4	81.3	94.3	77.5	75.1	62.3	33.2
F-DivideMix	Best	96.1	94.9	93.5	90.5	93.8	79.1	74.6	61.0	34.3
	Last	96.0	94.5	93.2	89.6	92.8	78.8	74.3	60.1	31.2

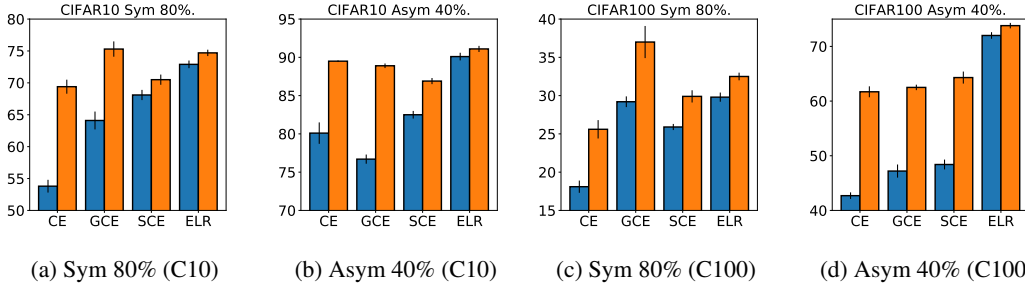


Figure 7: Test accuracies (%) on CIFAR-10 and CIFAR-100 under different noisy types and fractions for noise-robust loss approaches. Note that the blue and orange bars are results for without and with FINE, respectively. The average accuracies and standard deviations over three trials are reported.

loss functions on severe noise rate settings. The detailed results are reported in the Appendix. Unlike other methods, it is still theoretically supported because FINE extracts clean data with a robust classifier using representation.

4.3 Experiments on Real-World Dataset

As Table 3 shows, FINE and F-DivideMix work fairly well on the Clothing 1M dataset compared to other approaches when we reproduce the experimental results under the same settings.

Table 3: Test accuracy on Clothing1M dataset

Method	Standard	GCE [51]	SCE [40]	ELR [29]	DivideMix [25]	CORES ² [8]	FINE	F-DivideMix
Accuracy	68.94	69.75	71.02	72.87	74.30	73.24	72.91	74.37

5 Conclusion

This paper introduces FINE for detecting label noise by designing a robust noise detector. Our main idea is utilizing the principal components of latent representations made by eigen decomposition. Most existing detection methods are dependent on the loss values, while such losses may be biased by corrupted classifier [24, 31]. Our methodology alleviates this issue by extracting key information from representations without using explicit knowledge of the noise rates. We show that the FINE detector has an excellent ability to detect noisy labels in theoretical and experimental results. We propose three applications of the FINE detector: sample-selection approach, SSL approach, and collaboration with noise-robust loss functions. FINE yields strong results on standard benchmarks and a real-world dataset for various *LNL* approaches.

We believe that our work opens the door to detecting samples having noisy labels with explainable results. It is a non-trivial task and of social significance, and thus, our work will have a substantial social impact on DL practitioners because it avoids the a labor-intensive job of checking data label quality. As future work, we hope that our work will trigger interest in the design of new label-

noise detectors and bring a fresh perspective for other data-resampling approaches (e.g., anomaly detection and novelty detection). The development of robustness against label noise even leads to an improvement in the performance of network trained with data collected through web crawling. We believe that our contribution will lower the barriers to entry for developing robust models for DL practitioners and greatly impact the internet industry. On the other hand, we are concerned that it can be exploited to train robust models using data collected illegally and indiscriminately on the dark web (e.g., web crawling), and thus it may raise privacy concerns (e.g., copyright).

Acknowledgments and Disclosure of Funding

This work was supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) [No.2019-0-00075, Artificial Intelligence Graduate School Program (KAIST)] and [No. 2021-0-00907, Development of Adaptive and Lightweight Edge-Collaborative Analysis Technology for Enabling Proactively Immediate Response and Rapid Learning]. We thank Seongyoon Kim for discussing about the concept of perturbation.

References

- [1] Shin Ando and Chun Yuan Huang. Deep over-sampling framework for classifying imbalanced data. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 770–785. Springer, 2017.
- [2] Devansh Arpit, Stanisław Jastrzębski, 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 *International Conference on Machine Learning*, pages 233–242. PMLR, 2017.
- [3] Dara Bahri, Heinrich Jiang, and Maya Gupta. Deep k-nn for noisy labels. In *International Conference on Machine Learning*, pages 540–550. PMLR, 2020.
- [4] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *Advances in Neural Information Processing Systems*, pages 5049–5059, 2019.
- [5] Mateusz Buda, Atsuto Maki, and Maciej A Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 106:249–259, 2018.
- [6] Kaidi Cao, Yining Chen, Junwei Lu, Nikos Arechiga, Adrien Gaidon, and Tengyu Ma. Heteroskedastic and imbalanced deep learning with adaptive regularization. *arXiv preprint arXiv:2006.15766*, 2020.
- [7] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.
- [8] Hao Cheng, Zhaowei Zhu, Xingyu Li, Yifei Gong, Xing Sun, and Yang Liu. Learning with instance-dependent label noise: A sample sieve approach. *arXiv preprint arXiv:2010.02347*, 2020.
- [9] Charles Corbier, Hedi Ben-Younes, Alexandre Ramé, and Charles Ollion. Leveraging weakly annotated data for fashion image retrieval and label prediction. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 2268–2274, 2017.
- [10] Filipe R Cordeiro, Ragav Sachdeva, Vasileios Belagiannis, Ian Reid, and Gustavo Carneiro. Longremix: Robust learning with high confidence samples in a noisy label environment. *arXiv preprint arXiv:2103.04173*, 2021.
- [11] Di Feng, Christian Haase-Schutz, Lars Rosenbaum, Heinz Hertlein, Claudius Glaser, Fabian Timm, Werner Wiesbeck, and Klaus Dietmayer. Deep multi-modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges. *IEEE Transactions on Intelligent Transportation Systems*, 22(3):1341–1360, Mar 2021.

- [12] Ronald A Fisher. The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2):179–188, 1936.
- [13] Aritra Ghosh, Himanshu Kumar, and PS Sastry. Robust loss functions under label noise for deep neural networks. In *AAAI*, 2017.
- [14] 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 *Advances in neural information processing systems*, pages 8527–8537, 2018.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [16] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- [17] Dan Hendrycks, Mantas Mazeika, Duncan Wilson, and Kevin Gimpel. Using trusted data to train deep networks on labels corrupted by severe noise. *arXiv preprint arXiv:1802.05300*, 2018.
- [18] Jinchi Huang, Lie Qu, Rongfei Jia, and Binqiang Zhao. O2u-net: A simple noisy label detection approach for deep neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3326–3334, 2019.
- [19] 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, 2018.
- [20] Taehyeon Kim, Jaehoon Oh, NakYil Kim, Sangwook Cho, and Se-Young Yun. Comparing kullback-leibler divergence and mean squared error loss in knowledge distillation, 2021.
- [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25:1097–1105, 2012.
- [22] Ronan Le Bras, Swabha Swayamdipta, Chandra Bhagavatula, Rowan Zellers, Matthew Peters, Ashish Sabharwal, and Yejin Choi. Adversarial filters of dataset biases. In *International Conference on Machine Learning*, pages 1078–1088. PMLR, 2020.
- [23] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *Advances in Neural Information Processing Systems*, pages 7167–7177, 2018.
- [24] Kimin Lee, Sukmin Yun, Kibok Lee, Honglak Lee, Bo Li, and Jinwoo Shin. Robust inference via generative classifiers for handling noisy labels. In *International Conference on Machine Learning*, pages 3763–3772. PMLR, 2019.
- [25] Junnan Li, Richard Socher, and Steven C.H. Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. In *International Conference on Learning Representations*, 2020.
- [26] Mingchen Li, Mahdi Soltanolkotabi, and Samet Oymak. Gradient descent with early stopping is provably robust to label noise for overparameterized neural networks. In *International Conference on Artificial Intelligence and Statistics*, pages 4313–4324. PMLR, 2020.
- [27] Yi Li and Nuno Vasconcelos. Repair: Removing representation bias by dataset resampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9572–9581, 2019.
- [28] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017.
- [29] Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-learning regularization prevents memorization of noisy labels. *arXiv preprint arXiv:2007.00151*, 2020.

- [30] Michal Lukasik, Srinadh Bhojanapalli, Aditya Menon, and Sanjiv Kumar. Does label smoothing mitigate label noise? In *International Conference on Machine Learning*, pages 6448–6458. PMLR, 2020.
- [31] Hartmut Maennel, Ibrahim Alabdulmohsin, Ilya Tolstikhin, Robert JN Baldock, Olivier Bousquet, Sylvain Gelly, and Daniel Keysers. What do neural networks learn when trained with random labels? *arXiv preprint arXiv:2006.10455*, 2020.
- [32] Baharan Mirzasoleiman, Kaidi Cao, and Jure Leskovec. Coresets for robust training of neural networks against noisy labels. *arXiv preprint arXiv:2011.07451*, 2020.
- [33] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1944–1952, 2017.
- [34] Geoff Pleiss, Tianyi Zhang, Ethan R Elenberg, and Kilian Q Weinberger. Identifying mislabeled data using the area under the margin ranking. *arXiv preprint arXiv:2001.10528*, 2020.
- [35] Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping. *arXiv preprint arXiv:1412.6596*, 2014.
- [36] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *International Conference on Machine Learning*, pages 4334–4343. PMLR, 2018.
- [37] Hwanjun Song, Minseok Kim, Dongmin Park, and Jae-Gil Lee. Prestopping: How does early stopping help generalization against label noise?, 2020.
- [38] Brendan Van Rooyen, Aditya Krishna Menon, and Robert C Williamson. Learning with symmetric label noise: The importance of being unhinged. *arXiv preprint arXiv:1505.07634*, 2015.
- [39] Andreas Veit, Neil Alldrin, Gal Chechik, Ivan Krasin, Abhinav Gupta, and Serge Belongie. Learning from noisy large-scale datasets with minimal supervision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 839–847, 2017.
- [40] Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross entropy for robust learning with noisy labels. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 322–330, 2019.
- [41] Yi Wei, Xue Mei, Xin Liu, and Pengxiang Xu. Dst: Data selection and joint training for learning with noisy labels. *arXiv preprint arXiv:2103.00813*, 2021.
- [42] Pengxiang Wu, Songzhu Zheng, Mayank Goswami, Dimitris Metaxas, and Chao Chen. A topological filter for learning with label noise. *arXiv preprint arXiv:2012.04835*, 2020.
- [43] 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 *International Conference on Learning Representations*, 2021.
- [44] Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. Learning from massive noisy labeled data for image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2691–2699, 2015.
- [45] 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.
- [46] Seunghan Yang, Hyoungseob Park, Junyoung Byun, and Changick Kim. Robust federated learning with noisy labels, 2020.
- [47] Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor W Tsang, and Masashi Sugiyama. How does disagreement help generalization against label corruption? *arXiv preprint arXiv:1901.04215*, 2019.

- [48] Xiyu Yu, Tongliang Liu, Mingming Gong, and Dacheng Tao. Learning with biased complementary labels. In *Proceedings of the European conference on computer vision (ECCV)*, pages 68–83, 2018.
- [49] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- [50] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- [51] Zhilu Zhang and Mert R Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. *arXiv preprint arXiv:1805.07836*, 2018.
- [52] Zizhao Zhang, Han Zhang, Sercan O Arik, Honglak Lee, and Tomas Pfister. Distilling effective supervision from severe label noise. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9294–9303, 2020.
- [53] Zizhao Zhang, Han Zhang, Sercan O. Arik, Honglak Lee, and Tomas Pfister. Distilling effective supervision from severe label noise. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.