

Supplementary Material:

Robust Inference via Generative Classifiers for Handling Noisy Labels

A. Preliminaries

Gaussian discriminant analysis. In this section, we describe the basic concepts of the discriminative and generative classifier (Ng & Jordan, 2002). Formally, denote the random variable of the input and label as \mathbf{x} and $y = \{1, \dots, C\}$, respectively. For the classification task, the discriminative classifier directly defines a posterior distribution $P(y|\mathbf{x})$, i.e., learning a direct mapping between input \mathbf{x} and label y . A popular model for discriminative classifier is softmax classifier which defines the posterior distribution as follows: $P(y = c|\mathbf{x}) = \frac{\exp(\mathbf{w}_c^\top \mathbf{x} + b_c)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{x} + b_{c'})}$, where \mathbf{w}_c and b_c are weights and bias for a class c , respectively. In contrast to the discriminative classifier, the generative classifier defines the class conditional distribution $P(\mathbf{x}|y)$ and class prior $P(y)$ in order to indirectly define the posterior distribution by specifying the joint distribution $P(\mathbf{x}, y) = P(y)P(\mathbf{x}|y)$. Gaussian discriminant analysis (GDA) is a popular method to define the generative classifier by assuming that the class conditional distribution follows the multivariate Gaussian distribution and the class prior follows Bernoulli distribution: $P(\mathbf{x}|y = c) = \mathcal{N}(\mathbf{x}|\mu_c, \Sigma_c)$, $P(y = c) = \frac{\beta_c}{\sum_{c'} \beta_{c'}}$, where μ_c and Σ_c are the mean and covariance of multivariate Gaussian distribution, and β_c is the unnormalized prior for class c . This classifier has been studied in various machine learning areas (e.g., semi-supervised learning (Lasserre et al., 2006) and incremental learning (Lee et al., 2018)).

In this paper, we focus on the special case of GDA, also known as the linear discriminant analysis (LDA). In addition to Gaussian assumption, LDA further assumes that all classes share the same covariance matrix, i.e., $\Sigma_c = \Sigma$. Since the quadratic term is canceled out with this assumption, the posterior distribution of generative classifier can be represented as follows:

$$P(y = c|\mathbf{x}) = \frac{P(y = c)P(\mathbf{x}|y = c)}{\sum_{c'} P(y = c')P(\mathbf{x}|y = c')} = \frac{\exp(\mu_c^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_c^\top \Sigma^{-1} \mu_c + \log \beta_c)}{\sum_{c'} \exp(\mu_{c'}^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_{c'}^\top \Sigma^{-1} \mu_{c'} + \log \beta_{c'})}.$$

One can note that the above form of posterior distribution is equivalent to the softmax classifier by considering $\mu_c^\top \Sigma^{-1}$ and $-\frac{1}{2} \mu_c^\top \Sigma^{-1} \mu_c + \log \beta_c$ as its weight and bias, respectively. This implies that \mathbf{x} might be fitted in Gaussian distribution during training a softmax classifier.

Breakdown points. The robustness of MCD estimator can be explained by the fact that it has high breakdown points (Hampel, 1971). Specifically, the breakdown point of an estimator measures the smallest fraction of observations that need to be replaced by arbitrary values to carry the estimate beyond all bounds. Formally, denote \mathcal{Y}_M as a set obtained by replacing M data points of set \mathcal{Y} by some arbitrary values. Then, for a multivariate mean estimator $\mu = \mu(\mathcal{Y})$ from \mathcal{Y} , the breakdown point is defined as follows (see Appendix A for more detailed explanations including the breakdown point of covariance estimator):

$$\varepsilon^*(\mu, \mathcal{Y}) = \frac{1}{|\mathcal{Y}|} \min \left\{ M \in \{1, \dots, |\mathcal{Y}|\} : \sup_{\mathcal{Y}_M} \|\mu(\mathcal{Y}) - \mu(\mathcal{Y}_M)\| = \infty \right\}.$$

For a multivariate estimator of covariance Σ , we have

$$\varepsilon^*(\Sigma, \mathcal{Y}) = \frac{1}{|\mathcal{Y}|} \min \{ M \in \{1, \dots, |\mathcal{Y}|\} : \sup_M \max_i \{ |\log \lambda_i(\Sigma(\mathcal{Y})) - \log \lambda_i(\Sigma(\mathcal{Y}_M))| \} \},$$

where the k -th largest eigenvalue of a general $n \times n$ matrix is denoted by $\lambda_k(\Sigma)$, $k = 1, \dots, n$ such that $\lambda_1(\Sigma) \leq \lambda_2(\Sigma) \leq \dots \leq \lambda_n(\Sigma)$. This implies that we consider a covariance estimator to be broken whenever any of the eigenvalues can become arbitrary large or arbitrary close to 0.

B. Experimental setup

We describe the detailed explanation about all the experiments in Section 4. The code is available at [anonymized].

Detailed model architecture and datasets. We consider two state-of-the-art neural network architectures: DenseNet (Huang & Liu, 2017) and ResNet (He et al., 2016). For DenseNet, our model follows the same setup as in Huang & Liu

(2017): 100 layers, growth rate $k = 12$ and dropout rate 0. Also, we use ResNet with 34 and 44 layers, filters = 64 and dropout rate 0^5 . The softmax classifier is used, and each model is trained by minimizing the cross-entropy loss. We train DenseNet and ResNet for classifying CIFAR-10 (or 100) and SVHN datasets: the former consists of 50,000 training and 10,000 test images with 10 (or 100) image classes, and the latter consists of 73,257 training and 26,032 test images with 10 digits.⁶ By following the experimental setup of Ma et al. (2018), All networks were trained using SGD with momentum 0.9, weight decay 10^{-4} and an initial learning rate of 0.1. The learning rate is divided by 10 after epochs 40 and 80 for CIFAR-10/SVHN (120 epochs in total), and after epochs 80, 120 and 160 for CIFAR-100 (200 epochs in total). For our method, we extract the hidden features at $\{79,89,99\}$ -th layers and $\{27,29,31,33\}$ -th layers for DenseNet and ResNet, respectively. We assume the uniform class prior distribution.

In the Table 5, we evaluate RoG for the NLP tasks on Twitter and Reuters dataset: the former has a task of part-of-speech (POS) tagging, and the latter has a task of text categorization. The Twitter dataset consists of 14,619 training data from 25 different classes, while some of classes only contain a small number of training data. We exclude such classes that are smaller than 100 in size, then we finally have 14,468 training data and 7,082 test data from 19 different classes. Similarly, we exclude the classes which have a size of less than 100 in the Reuters dataset. Then it consists of 5,444 training data and 2,179 test data from 7 different classes. For training, we use 2-layer FCNs with ReLU non-linearity and uniform noise on both Twitter and Reuters datasets and we extract the hidden features at both layers for 2-layer FCNs.

We also consider the open-set noisy scenarios (Wang et al., 2018). Figure 4 shows the examples of open-set noisy dataset which is built by replacing some training images in CIFAR-10 by external images in CIFAR-100 and Downsampled ImageNet (Chrabaszcz et al., 2017) which is equivalent to the ILSVRC 1,000-class ImageNet dataset (Deng et al., 2009), but with images downsampled to 32×32 resolution. In the Table 7, we maintain the original label of the CIFAR-10 dataset, while replacing 60% of the training data in CIFAR-10 with training data in CIFAR-100 and Downsampled ImageNet.

Validation. For our methods, the ensemble weights are chosen by optimizing the NLL loss over the validation set. We assume that the validation set consists of 1000 images with (same type and fraction of) noisy labels. However, one can expect that if we use all validation samples, the performance of our method can be affected by outliers. To relax this issue, we use only half of them, chosen by the MCD estimator. Specifically, we first compute the Mahalanobis distance for all validation samples using the parameters from MCD estimator, and select 500 samples with smallest distance. Then, one can expect that our ensemble method is more robust against the noisy labels in validation sets. In the case of Twitter and Reuters dataset, validation set consists of 570 and 210 samples with noisy labels, respectively.

Training method for noisy label learning. We consider the following training methods for noisy label learning:

- (a) **Hard bootstrapping** (Reed et al., 2014): Training with new labels generated by a convex combination (the hard version) of the noisy labels and their predicted labels.
- (b) **Soft bootstrapping** (Reed et al., 2014): Training with new labels generated by a convex combination (the soft version) of the noisy labels and their predictions.
- (c) **Backward** (Patrini et al., 2017): Training via loss correction by multiplying the cross-entropy loss by a noise-aware correction matrix.
- (d) **Forward** (Patrini et al., 2017): Training with label correction by multiplying the network prediction by a noise-aware correction matrix.
- (e) **Forward (Gold)** (Hendrycks et al., 2018): An augmented version of Forward method which replaces its corruption matrix estimation with the identity on trusted samples.
- (f) **GLC (Gold Loss Correction)** (Hendrycks et al., 2018): Training with the corruption matrix which is estimated by using the trusted dataset.
- (g) **D2L** (Ma et al., 2018): Training with new labels generated by a convex combination of the noisy labels and their predictions, where its weights are chosen by utilizing the Local Intrinsic Dimensionality (LID).
- (h) **Decoupling** (Malach & Shalev-Shwartz, 2017): Updating the parameters only using the samples which have different prediction from two classifier.

⁵ResNet architecture is available at <https://github.com/kuangliu/pytorch-cifar>.

⁶We do not use the extra SVHN dataset for training.

- (i) **MentorNet** (Jiang et al., 2018): An extra teacher network is pre-trained and then used to select clean samples for its student network.
- (j) **Co-teaching** (Han et al., 2018b): A simple ensemble method where each network selects its small-loss training data and back propagates the training data selected by its peer network.
- (k) **Cross-entropy**: the conventional approach of training with cross-entropy loss.

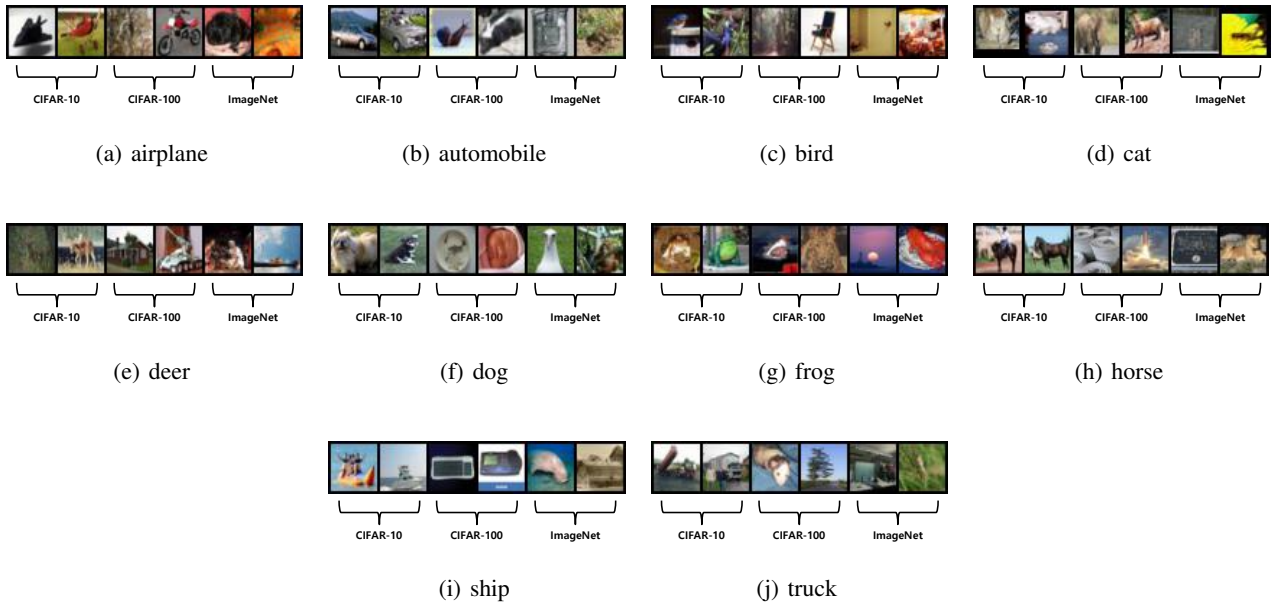


Figure 4. Examples of open-set noise in CIFAR-10 dataset.

C. Layer-wise characteristics of generative classifiers

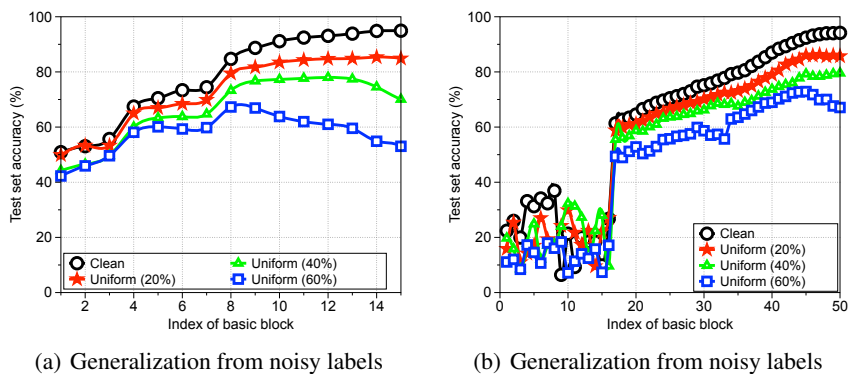


Figure 5. Layer-wise characteristics of generative classifiers from (a) ResNet-34 and (b) DenseNet-100 trained on the CIFAR-10 dataset.

Figure 5 shows the classification accuracy of the generative classifiers from different basic blocks of ResNet-34 (He et al., 2016) and DenseNet-100 (Huang & Liu, 2017). One can note that the generative classifiers from DenseNet and ResNet have different patterns due to the architecture design. In the case of DenseNet, we found that it produces meaningful features after 20-th basic blocks.

D. Proof of Theorem 1

In this section, we present a proof of Theorem 1, which consists of two statements: the limit of estimation error (3) and estimated error ratio (4). We prove both statements one by one as stated in below. For convenience, we skip to mention the Continuous Mapping Theorem⁷ and the number of training samples N goes to infinity for all convergences in the proof.

D.1. Proof of the limit of estimation error (3)

We start with a following lemma, which shows the convergences of sample and MCD estimators as the number of training samples N goes to infinity.

Lemma 1 *Suppose we have N number of d -dimensional training samples $\mathcal{X}_N = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and \mathcal{X}_N contains outlier samples with the fixed fraction $\delta_{\text{out}} < 1$. We assume the outlier samples are from an arbitrary distribution P_{out} with zero mean and finite covariance matrix $\sigma_{\text{out}}^2 \mathbf{I}$, and the clean samples are from a distribution of the hidden features of DNNs P_{clean} with mean μ and covariance matrix $\sigma^2 \mathbf{I}$. Let $\bar{\mu}$ and $\bar{\Sigma}$ be the mean and covariance matrix of sample estimator, and let $\hat{\mu}$ and $\hat{\Sigma}$ be the mean and covariance matrix of MCD estimator which selects samples from \mathcal{X}_N with the fixed fraction $\frac{d}{N} < \delta_{\text{mcd}} < 1$ to optimize its objective (2). Then the mean and covariance matrix of sample estimator converge almost surely to below as $N \rightarrow \infty$:*

$$\bar{\mu} \xrightarrow{\text{a.s.}} (1 - \delta_{\text{out}}) \mu, \quad \bar{\Sigma} \xrightarrow{\text{a.s.}} ((1 - \delta_{\text{out}}) \sigma^2 + \delta_{\text{out}} \sigma_{\text{out}}^2) \mathbf{I} + \delta_{\text{out}} (1 - \delta_{\text{out}}) \mu \mu^T.$$

In addition, if $\delta_{\text{mcd}} \leq 1 - \delta_{\text{out}}$ and $\sigma^2 < \sigma_{\text{out}}^2$, the mean and covariance matrix of MCD estimator converge almost surely to below as $N \rightarrow \infty$:

$$\hat{\mu} \xrightarrow{\text{a.s.}} \mu, \quad \hat{\Sigma} \xrightarrow{\text{a.s.}} \sigma^2 \mathbf{I}.$$

A proof of the lemma is given in appendix D.3, where it is built upon the fact that the determinant of covariance matrix with some assumptions can be expressed as the d -th degree polynomial of outlier ratio.

Lemma 1 states the convergences of sample and MCD estimators on a single distribution of hidden features of DNNs. Without loss of generality, one can assume the mean of outlier distribution is zero, i.e., $\mu_{\text{out}} = 0$ by an affine translation of hidden features. Furthermore, one can extend Lemma 1 to C number of distributions, which have the class mean μ_c and class covariance matrix Σ_c on each class label $c \in \{1, \dots, C\}$ with the assumptions $\mathcal{A}1 \sim \mathcal{A}4$. Then the class mean of MCD and sample estimators converge almost surely as follows:

$$\hat{\mu}_c \xrightarrow{\text{a.s.}} \mu_c, \quad \bar{\mu}_c \xrightarrow{\text{a.s.}} (1 - \delta_{\text{out}}) \mu_c,$$

which implies that

$$\|\mu_c - \hat{\mu}_c\|_1 \xrightarrow{\text{a.s.}} 0, \quad \|\mu_c - \bar{\mu}_c\|_1 \xrightarrow{\text{a.s.}} \delta_{\text{out}} \|\mu_c\|_1.$$

This completes the proof of the limit of estimation error (3).

D.2. Proof of the limit of estimated error ratio (4)

Recall the class mean of MCD and sample estimators converge almost surely as follows:

$$\hat{\mu}_c \xrightarrow{\text{a.s.}} \mu_c, \quad \bar{\mu}_c \xrightarrow{\text{a.s.}} (1 - \delta_{\text{out}}) \mu_c.$$

Then one can induce that the limit of mean distance of sample and MCD estimators as follow:

$$\|\bar{\mu}_c - \bar{\mu}_{c'}\|_2 \xrightarrow{\text{a.s.}} (1 - \delta_{\text{out}})^2 \|\mu_i - \mu_c\|_2, \quad (5)$$

$$\|\hat{\mu}_c - \hat{\mu}_{c'}\|_2 \xrightarrow{\text{a.s.}} \|\mu_c - \mu_{c'}\|_2. \quad (6)$$

⁷P. Billingsley, Convergence of Probability Measures, John Wiley & Sons, 1999

On the other hand, the assumptions A1 states that all class covariance matrices are the same, i.e., $\Sigma_c = \sigma^2 \mathbf{I}$. Then tied covariance matrices $\bar{\Sigma}$ and $\hat{\Sigma}$ are given by gathering $\bar{\Sigma}_c$ and $\hat{\Sigma}_c$ on each class c respectively:

$$\bar{\Sigma} = \frac{\sum_c N_c \bar{\Sigma}_c}{\sum_c N_c} = \frac{\sum_c \bar{\Sigma}_c}{C}, \quad \hat{\Sigma} = \frac{\sum_c K_c \hat{\Sigma}_c}{\sum_c K_c} = \frac{\sum_c \hat{\Sigma}_c}{C}. \quad (7)$$

From the tied covariance matrices (7) and Lemma 1, one can induce their convergences and limits as follow:

$$\bar{\Sigma} \xrightarrow{\text{a.s.}} ((1 - \delta_{\text{out}}) \sigma^2 + \delta_{\text{out}} \sigma_{\text{out}}^2) \mathbf{I} + \delta_{\text{out}} (1 - \delta_{\text{out}}) \frac{1}{C} \sum_c \mu_c \mu_c^T, \quad (8)$$

$$\hat{\Sigma} \xrightarrow{\text{a.s.}} \sigma^2 \mathbf{I}. \quad (9)$$

Next, we define a function of the tied covariance matrix as follow:

$$\phi(\hat{\Sigma}) = 4 \|\hat{\Sigma}^{-1}\|_2 \|\hat{\Sigma}\|_2 \left(1 + \|\hat{\Sigma}^{-1}\|_2 \|\hat{\Sigma}\|_2\right)^{-2}.$$

Then (9) implies that $\phi(\hat{\Sigma}) \xrightarrow{\text{a.s.}} 1$ clearly. Since the condition number t of $\hat{\Sigma}^{-1}$ is in $[1, \infty)$, i.e. $t = \|\hat{\Sigma}^{-1}\|_2 \|\hat{\Sigma}\|_2 \in [1, \infty)$, one can induces that $\phi(\hat{\Sigma}) = \phi(t) = \frac{4t}{(1+t)^2}$ for $t \in [1, \infty)$ and it is a monotonic decreasing function by using the change of variables with t . Hence $\phi(t)$ has the maximum at $t = 1$ and it implies

$$1 = \phi(1) = \lim_{N \rightarrow \infty} \phi(\hat{\Sigma}) \geq \lim_{N \rightarrow \infty} \phi(\bar{\Sigma}). \quad (10)$$

Therefore the limits of mean distance of estimators (5), (6) and ratio of the function ϕ (10) hold the statement (4),

$$\frac{\phi(\hat{\Sigma}) \|\hat{\mu}_c - \hat{\mu}_{c'}\|_2}{\phi(\bar{\Sigma}) \|\bar{\mu}_c - \bar{\mu}_{c'}\|_2} \xrightarrow{\text{a.s.}} \lim_{N \rightarrow \infty} \frac{\phi(\hat{\Sigma}) \|\hat{\mu}_c - \hat{\mu}_{c'}\|_2}{\phi(\bar{\Sigma}) \|\bar{\mu}_c - \bar{\mu}_{c'}\|_2} = \lim_{N \rightarrow \infty} \frac{1}{(1 - \delta_{\text{out}})^2 \phi(\bar{\Sigma})} \geq 1.$$

This completes the proof of Theorem 1.

D.3. Proof of Lemma 1

In this part, we present a proof of Lemma 1. We show the almost surely convergences of sample and MCD estimators as the number of training samples N goes to infinity.

Proof of the convergence of sample estimator. First of all, the set of training samples $\mathcal{X}_N = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ contains outlier samples with the fixed fraction δ_{out} . So, \mathcal{X}_N is from a mixture distribution $P_{\text{mix}} = (1 - \delta_{\text{out}}) P_{\text{clean}} + \delta_{\text{out}} P_{\text{out}}$. Then mean and covariance matrix of sample estimator, $\bar{\mu}$ and $\bar{\Sigma}$, estimate mean μ_{mix} and covariance matrix Σ_{mix} of the mixture distribution P_{mix} , respectively. One can induce μ_{mix} and Σ_{mix} directly as follow:

$$\mu_{\text{mix}} = (1 - \delta_{\text{out}}) \mu, \quad \Sigma_{\text{mix}} = (1 - \delta_{\text{out}}) \sigma^2 \mathbf{I} + \delta_{\text{out}} \sigma_{\text{out}}^2 \mathbf{I} + \delta_{\text{out}} (1 - \delta_{\text{out}}) \mu \mu^T. \quad (11)$$

Since P_{mix} has the finite covariance matrix, i.e., $\Sigma_{\text{mix}} < \infty$, one can apply the the Strong Law of Large Numbers⁸ to the sample estimator of the mixture distribution P_{mix} . Then the mean and covariance matrix of sample estimator converge almost surely to the mean and covariance matrix of P_{mix} , respectively:

$$\bar{\mu} \xrightarrow{\text{a.s.}} \mu_{\text{mix}}, \quad \bar{\Sigma} \xrightarrow{\text{a.s.}} \Sigma_{\text{mix}}.$$

This completes the proof of the convergence of sample estimator.

Proof of the convergence of MCD estimator. Consider a collection E_q of subsets $\mathcal{X}_{K,q} \subset \mathcal{X}_N$ with the size $K (= \lfloor \delta_{\text{mcd}} N \rfloor)$, and each subset $\mathcal{X}_{K,q} \in E_q$ contains the outlier samples with the fraction $q \in [0, 1]$. Then $\mathcal{X}_{K,q} \in E_q$ is from a mixture distribution $P_q = (1 - q) P_{\text{clean}} + q P_{\text{out}}$. One can induce that the mean μ_q and covariance matrix Σ_q of the mixture distribution P_q as (11):

$$\mu_q = (1 - q) \mu, \quad \Sigma_q = (1 - q) \sigma^2 \mathbf{I} + q \sigma_{\text{out}}^2 \mathbf{I} + q(1 - q) \mu \mu^T. \quad (12)$$

⁸W. Feller, An Introduction to Probability Theory and Its Applications, John Wiley & Sons, 1968

Thus sample mean estimator $\bar{\mu}_{\mathcal{X}_{K,q}}$ and covariance estimator $\bar{\Sigma}_{\mathcal{X}_{K,q}}$ of a subset $\mathcal{X}_{K,q}$ converge almost surely to μ_q and Σ_q respectively:

$$\bar{\mu}_{\mathcal{X}_{K,q}} \xrightarrow{\text{a.s.}} \mu_q, \quad \bar{\Sigma}_{\mathcal{X}_{K,q}} \xrightarrow{\text{a.s.}} \Sigma_q,$$

by the Strong Law of Large Numbers.

On the other hand, there is a subset $\mathcal{X}_{K,q}^* \subset \mathcal{X}_N$ in E_{q^*} which is selected by MCD estimator. Then the determinant of its covariance matrix is the minimum over all subset of size K in \mathcal{X}_N , and $\bar{\mu}_{\mathcal{X}_{K,q}^*} = \hat{\mu} \xrightarrow{\text{a.s.}} \mu_{q^*}$, $\bar{\Sigma}_{\mathcal{X}_{K,q}^*} = \hat{\Sigma} \xrightarrow{\text{a.s.}} \Sigma_{q^*}$ as $N \rightarrow \infty$. Since the determinant is a continuous function, the Continuous Mapping Theorem⁹ implies

$$\min_{\mathcal{X}_{K,q} \in E_q, \forall q} \det(\bar{\Sigma}_{\mathcal{X}_{K,q}}) \xrightarrow{\text{a.s.}} \min_q \det(\Sigma_q),$$

and

$$\min_{\mathcal{X}_{K,q} \in E_q, \forall q} \det(\bar{\Sigma}_{\mathcal{X}_{K,q}}) = \det(\bar{\Sigma}_{\mathcal{X}_{K,q}^*}) = \det(\hat{\Sigma}) \xrightarrow{\text{a.s.}} \det(\Sigma_{q^*}).$$

Now, we'd like to show

$$\min_q \det(\Sigma_q) = \det(\Sigma_{q^*}) = \det(\Sigma_0), \quad (13)$$

to complete the proof of Lemma 1.

By the assumption $\delta_{\text{mcd}} \leq 1 - \delta_{\text{out}}$, E_0 is non-empty. It shows the existence of Σ_0 . From the covariance matrix Σ_q (12), $\det(\Sigma_q)$ is a d -th degree polynomial of q as follow:

$$\begin{aligned} \det(\Sigma_q) &= \det((1-q)\sigma^2 \mathbf{I} + q\sigma_{\text{out}}^2 \mathbf{I} + q(1-q)\mu\mu^T) \\ &= ((1-q)\sigma^2 + q\sigma_{\text{out}}^2)^{d-1} ((1-q)\sigma^2 + q\sigma_{\text{out}}^2 + q(1-q)\mu^T \mu). \end{aligned}$$

Since the assumption gives $\sigma_{\text{out}}^2 > \sigma^2$, $\det(\Sigma_q)$ has the lower bound $\det(\Sigma_0)$ as follow:

$$\begin{aligned} \det(\Sigma_q) &= ((1-q)\sigma^2 + q\sigma_{\text{out}}^2)^{d-1} ((1-q)\sigma^2 + q\sigma_{\text{out}}^2 + q(1-q)\mu^T \mu) \\ &\geq (\sigma^2)^{d-1} (\sigma^2 + q(1-q)\mu^T \mu) \\ &\geq (\sigma^2)^{d-1} \sigma^2 = \det(\Sigma_0). \end{aligned}$$

Then $\det(\Sigma_q) \geq \det(\Sigma_0)$ for all $q \in [0, 1]$ and the equality holds for only $q = 0$. It implies $q^* = 0$ and (13) is the shown. Therefore the mean and covariance matrix of MCD estimator converge almost surely to μ and $\sigma^2 \mathbf{I}$, respectively:

$$\hat{\mu} \xrightarrow{\text{a.s.}} \mu_0 = \mu, \quad \hat{\Sigma} \xrightarrow{\text{a.s.}} \Sigma_0 = \sigma^2 \mathbf{I}.$$

This completes the proof of Lemma 1.

E. Comparison with robust clustering methods

To remove the outliers (i.e., samples with noisy labels) in hidden feature, one can consider other robust clustering methods to estimate the parameters of generative classifiers. In this section, we test trimmed K-means (TKM) (Garcia-Escudero & Gordaliza, 1999) as a new baseline, and compare the performances with MCD estimator. Specifically, we use noisy labels to initialize clusters of TKM, and assign a "majority" label to each trained cluster. For pair comparison, a tied covariance across clusters is assumed, and we run the same number of iterations for both TKM and MCD. Table 8 reports the corresponding results under ResNet and CIFAR-10 with uniform noises when we induce a generative classifier only using a penultimate layer. First, we remark that TKM outperforms Softmax, which supports our claim that the clustering property of DNN features is useful to handle noisy labels. However, the generative classifier with MCD estimator still outperforms all baselines because it utilizes the information of noisy labels carefully to derive a new decision rule at all iterations, while TKM is essentially an unsupervised method and does not utilize the information, except for initialization and termination.

⁹P. Billingsley, Convergence of Probability Measures, John Wiley & Sons, 1999

Model	Inference method	CIFAR-10			
		Clean	Uniform (20%)	Uniform (40%)	Uniform (60%)
ResNet	Softmax	94.76%	80.88%	61.98%	39.96%
	Generative + TKM	94.35%	81.41%	63.27%	41.78%
	Generative + MCD (ours)	94.76%	83.86%	68.03%	44.87%

Table 8. Test accuracy (%) of ResNet on the CIFAR-10 dataset with uniform noise. The best results are highlighted in bold if the gain is bigger than 1%.