Safeguarded Dynamic Label Regression for Generalized Noisy Supervision

Jiangchao Yao, Ya Zhang, Ivor Tsang and Jun Sun

Abstract—Learning with noisy labels, which aims to reduce expensive labors on accurate annotations, has become imperative in the Big Data era. Previous noise transition based method has achieved promising results and presented a theoretical guarantee on performance in the case of class-conditional noise. However, this type of approaches critically depend on an accurate preestimation of the noise transition, which is usually impractical. Subsequent improvement adapts the pre-estimation along with the training progress via a Softmax layer. However, the parameters in the Softmax layer are highly tweaked for the fragile performance due to the ill-posed stochastic approximation. To address these issues, we propose a Latent Class-Conditional Noise model (LCCN) that naturally embeds the noise transition under a Bayesian framework. By projecting the noise transition into a Dirichlet-distributed space, the learning is constrained on a simplex based on the whole dataset, instead of some ad-hoc parametric space. We then deduce a dynamic label regression method for LCCN to iteratively infer the latent labels, to stochastically train the classifier and to model the noise. Our approach safeguards the bounded update of the noise transition, which avoids previous arbitrarily tuning via a batch of samples. We further generalize LCCN for open-set noisy labels and the semi-supervised setting. We perform extensive experiments with the controllable noise data sets, CIFAR-10 and CIFAR-100, and the agnostic noise data sets, Clothing1M and WebVision17. The experimental results have demonstrated that the proposed model outperforms several state-of-the-art methods.

Index Terms—Bayesian learning, noisy labels, Gibbs sampling.

I. INTRODUCTION

ARGE scale datasets with editorial labels have driven the success of deep neural networks (DNNs) in computer vision [1], natural language processing [2], and speech recognition [3]. However, for many real-world applications, it is usually expensive to collect accurately annotated data in large volume. Instead, samples with noisy supervision, as an alternative to alleviate the annotation burden, can be acquired inexhaustibly on the social websites and have shown potential to many applications in the deep learning area [4–7].

It is challenging to train DNNs in the presence of noisy supervision since it can easily memorize the clean data as well as the noisy data [8]. To overcome the above issue, several methods have been explored from the perspective of model regularization and sample re-weighting, respectively. Arpit *et al.* [8] applied the dropout regularization in DNNs to limit its speed of memorizing noise, which prevents the classifier from noise pollution. Ren *et al.* [9] explored dynamically weighting noisy labels with the corresponding predictions to weaken the noise effect. However, the model regularization and sample re-weighting based methods usually require either

a careful hyperparameter setting [8], or auxiliary samples [10] or elaborate curricula [11].

The study of this paper falls into the third popular perspective, learning with noise transition, which places a noise transition on top of the classifier. Early study [12] presents a two-step solution, that is, first pre-estimate the noise transition and then fix it to train the classifier. However, it suffers from the inaccurate pre-estimation via an ideal but impractical anchor set. Subsequent improvement [13] uses the stochastic approximation to adapt the noise transition in the form of a Softmax layer along with the training progress. Although it shows promise, the optimization of the Softmax layer depends on highly tweaking and the model parameters easily fall into undesired local minimums. Essentially, such instability is due to the inconsideration of the global dependency in the stochastic approximation, yielding a "local" mini-batch of samples can unbounded update the "global" noise transition in the back-propagation.

To solve this issue, we propose a Latent Class-Conditional Noise model (LCCN) that embeds the noise transition into a Dirichlet-distributed space. Compared to the previous Softmax layer [13], LCCN constrains the learning of the noise transition as a global variable depending on the whole dataset. Namely, a "local" mini-batch of samples can only partially affect the estimation of the "global" noise transition. Besides, a new dynamic label regression method is derived to stochastically optimize LCCN. Although it iteratively infers the latent labels and applies them for the classifier training and the noise modeling, only a small amount of extra computational cost is introduced. We theoretically demonstrate our method safeguards the bounded update of the noise transition via a minibatch of samples. Fig. 1 provides a simple illustration of our safeguarded dynamic label regression for LCCN. As can be seen, images are first inputted to the classifier to have the prediction of latent labels. Noisy labels are also forwarded to Bayesian noise modeling to compute the conditional transition of latent labels. Then, the latent labels are sampled based on their product and used to supervise the classifier training and refine the noise modeling. In a nutshell, our main contributions can be summarized into the following three points.

- We propose a Latent Class-Conditional Noise model that embeds the noise transition into a Dirichlet space to emphasize its global dependency, and then deduce a scalable dynamic label regression method for its optimization.
- The theoretical analysis on the convergence of the dynamic label regression, the generalization gap as well as the complexity is provided. Importantly, we prove that our optimization of the noise transition via a batch of samples

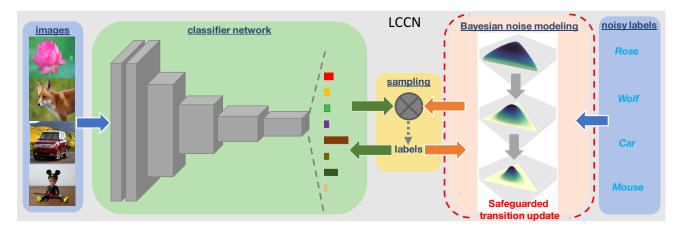


Fig. 1. Safeguarded dynamic label regression for LCCN. The images and noisy labels are respectively inputted to the classifier and the safeguarded Bayesian noise modeling to compute the prediction and the conditional transition. Then, the latent labels are sampled based on their product and then used for the classifier training and the safeguarded Bayesian noise modeling.

is bounded to avoid previous non-trivial tweaking.

- A more general variant of LCCN is further extended in order to handle the open-set noisy labels setting and the semi-supervised learning setting for the practical needs.
- We conduct a range of experiments in the popular CIFAR-10, CIFAR-100 datasets and large real-world noisy datasets, Clothing1M and WebVision17. Comprehensive results have demonstrated the superior performance of our model compared with existing state-of-the-art methods.

The rest part of this paper is organized as follows. Section II briefly reviews the related research of learning with noisy labels in deep learning. Then, we introduce our Latent Class-Conditional model and the dynamic label regression method in Section III, where the corresponding theoretical analysis and the further extension of LCCN is also included. We validate the efficiency of our method over a range of experiments in Section IV. Section V concludes the whole paper.

II. RELATED WORK

Recently, several approaches combined with deep learning have been developed for learning with noisy labels. In this section, we review these works according to noise transition, sample re-weighting and model regularization.

1) Learning with Noise Transition: This branch of research models a noise transition on top of the classifier to minimize the influence of label noise. Sukhbaatar et al. [14] introduced a noise transition matrix on top of CNN to learn with noisy supervision. With a heuristic learning procedure, they gradually make the transition matrix absorb the noise among labels. Misra et al. [15] considered the "reporting bias" phenomenon in human-centric annotations via a content-based transition, which is a special case of learning with noisy labels. Patrini et al. [12] theoretically demonstrated: the backward correction with the inverse of the noise transition is unbiased to train the classifier in the presence of noisy labels; the forward noise transition make the training share the same minimizer with that on the clean data. However, the performance quite depends on the accuracy of the pre-estimated noise transition. Subsequent improvement in [13] models the noise transition via a Softmax layer and tunes its parameters along with the training progress. Based on this research, Yao et al. [7] introduced an auxiliary variable to augment the noise transition with more uncertainty. The structure information [16] is further added to constrain the optimization. Although better performance has been achieved, these methods depend on the carefully tweaking. However, our model embeds the noise transition into a nonparametric space and naturally constrains its optimization to avoid undesired minimums via a dynamic label regression method.

- 2) Learning with Sample Re-weighting: This line of works weight the contribution of each training sample in parameter estimation to reduce the effect of label noise [10, 17]. It can be implemented by the label or the training pair re-weighting. For example, Reed et al. [18] facilitated the notion of perceptual consistency to linearly combine the label and the prediction as the new supervision, which shows the substantial robustness to label noise. Then, Li et al. [19] substituted the prediction with the refined label by the graph distillation. Wang et al. [20] leveraged the local intrinsic dimensionality to design an self-weighting strategy for Bootstrapping [18]. Recently, several works [21–23] also explore to collaboratively learn a weight or selection for each training pair and adjust their contribution to the training of the classifier. However, these methods critically depend on the elaborate sample re-weighting strategy.
- 3) Learning with Model Regularization: This type of methods attempt to regularize the training procedure in the presence of noisy supervision. Zhang et al. [24] have shown DNNs can easily memorize the random labels completely, characterizing the challenge to deep learning with noisy labels. Their further study [25] that used the convex combinations of images and noisy labels as the data augmentation, has been demonstrated as an efficient regularization to prevent DNNs from overfitting. Arpit et al. [8] investigated the memorization order of DNNs on feature patterns in noisy datasets and demonstrated dropout can efficiently limit the speed of memorization on noise in DNNs. Tanaka et al. [26] explicitly introduced a regularization term to prevent the trivial case of assigning all labels to a single class in label correction. Compared with above methods, we indirectly regularize the training by Bayesian noise modeling.

III. THE PROPOSED FRAMEWORK

A. Preliminaries

In the c-class classification setting, a collection of N noisy training pairs $\{(x_n,y_n)\}_{n=1}^N$ is given, where x_n is the raw input data or the feature vector and $y_n \in \{1,\ldots,K\}$ is the corresponding noisy label. Assume z_n denotes the latent label of x_n , which is unknown in practice. Then the goal in this task is to train a deep network classifier from the noisy dataset $\{(x_n,y_n)\}_{n=1}^N$ analogous to the one trained from the clean dataset $\{(x_n,z_n)\}_{n=1}^N$, so that a promising performance can be achieved in a clean test dataset. As shown in [24], directly minimizing the following equation will make DNNs memorize both the classification pattern and noise,

$$\hat{f}_{\theta} = \underset{f_{\theta} \in \mathcal{F}}{\operatorname{arg\,min}} - \frac{1}{N} \sum_{n=1}^{N} \ell(y_n, f_{\theta}(x_n)), \tag{1}$$

where f_{θ} is from the function class \mathcal{F} , which is parameterized by θ via DNNs, and ℓ is the loss function between y_n and the prediction $f_{\theta}(x_n)$. Eq. (1) leads to a bad performance in the clean test dataset since it does not squeeze out the noise influence from f_{θ} . Therefore, we follows one mainstream of approaches to handle this dilemma, which models a noise transition ϕ in simplex Δ when learning with noisy labels. The objective is then mathematically expressed with the following empirical risk minimization problem

$$\hat{f}_{\theta}, \phi = \underset{f_{\theta} \in \mathcal{F}, \phi \in \Delta}{\operatorname{arg\,min}} - \frac{1}{N} \sum_{n=1}^{N} \ell(y_n, \phi \circ f_{\theta}(x_n)), \tag{2}$$

Patrini et al. [12] theoretically demonstrate Eq. (2) trained with the noisy data shares the same minimizer with Eq. (1) trained with the clean data, if ϕ is accurately estimated. Unfortunately, it is usually impractical to acquire such a ϕ in advance. Thus, subsequent work [13] adapts the pre-estimation with a Softmax layer along with the training progress. Although this shows a promising performance, expensive tweaking is required due to the ill-posed stochastic approximation as a simple neural layer.

B. Latent Class-Conditional Noise model

In this section, we will present our Latent Class-Conditional Noise model (LCCN). Specifically, it avoids non-trivially tweaking for the fragile performance in [13] by modeling ϕ in a Bayesian form. The graphical notation is illustrated in Fig. 2 and the generative procedure is summarized as follows,

- The latent label $z_n \sim P(\cdot|x_n)$, where $P(\cdot|x_n)$ is a *Categorical* distribution modeled by the deep neural network f_{θ} and the given x_n is its input feature.
- The transition vector of the kth class $\phi_k \sim Dirichlet(\alpha)$, where α is the parameter of a Dirichlet distribution and $[\phi_1, \cdots, \phi_K]^T$ constitutes the noise transition matrix.
- The observed noisy label $y_n \sim P(\cdot | \phi_{z_n})$, where $P(\cdot | \phi_{z_n})$ is a *Categorical* distribution parameterized by ϕ_{z_n} .

The general way to solve such a probabilistic model combined with deep learning is amortized variational inference [27, 28]. However, this way for LCCN will require an approximate Categorical reparameterization [29, 30] and introduce an unstable Digamma function to optimize. To avoid this issue,

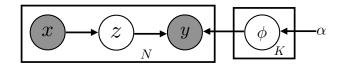


Fig. 2. Latent Class-Conditional Noise model. x and y is the observed training pair. z is the latent label. ϕ is the unknown noise transition. α is a Dirichlet parameter. N is the sample number and K is the class number.

we specifically deduce a dynamic label regression method for optimization and demonstrate its safeguarded update for ϕ .

C. Dynamic Label Regression

In the following, we will give the dynamic label regression method for LCCN, which stacks an *autoencoded Gibbs sampling* to infer the latent labels and loss minimization for parameter learning. It naturally suits LCCN and we show its deduction via a two-step formulation. Note that, in despite of the complex deduction, only a small computational cost is extra introduced, which will be explained in the following section. Simply, the first step is computing the probability of each z conditional on the others Z^{-1} , i.e., $P(z_n|Z^{-n})$. Then, with the samples from $P(z_n|Z^{-n})$, the classifier training and the noise modeling can be explicitly decoupled as the following optimization problem,

$$\begin{cases} \min -\frac{1}{n} \sum_{n=1}^{N} \ell_1(z_n, P(z_n | x_n)) \\ \min -\frac{1}{n} \sum_{n=1}^{N} \ell_2(y_n, P(y_n | z_n)). \end{cases}$$
(3)

 ℓ_1 is the ξ -clipped cross-entropy loss² and ℓ_2 is the likelihood loss. Alternating between the sampling of $P(z_n|Z^{-n})$ and the optimization of Eq. (3) constructs our final algorithm to learn with noisy supervision. Specifically, when P(z|x) approach the true distribution of clean labels, the classifier training is similar to that on the clean dataset. This yields the asymptotically unbiased estimation as on the clean datasets.

Autoencoded Gibbs sampling. Firstly, according to the aforementioned generative process, we can easily deduce the posterior of z conditioned on the observed training pair $\{(x_n,y_n)\}_{n=1}^N$ and the Dirichlet parameter α . This is implemented by factorizing the target conditional probability based on Fig. 2 and applying the Bayes theorem as follows,

$$P(Z|X,Y;\alpha) = \int_{\phi} \prod_{k=1}^{K} P(\phi_{k};\alpha) \prod_{n=1}^{N} P(z_{n}|x_{n},y_{n},\phi) d\phi$$

$$= \int_{\phi} \prod_{k=1}^{K} P(\phi_{k};\alpha) \prod_{n=1}^{N} \frac{P(z_{n}|x_{n})P(y_{n}|z_{n},\phi)}{P(y_{n}|x_{n})} d\phi$$

$$= S * \int_{\phi} \prod_{k=1}^{K} \frac{\Gamma(\sum_{k'}^{K} \alpha_{k'})}{\prod_{k'}^{K} \Gamma(\alpha_{k'})} \prod_{k'}^{K} \phi_{kk'}^{\alpha_{k'}-1} \prod_{n=1}^{N} \phi_{z_{n}y_{n}} d\phi,$$
(4)

where S represents $\prod_{n=1}^N \frac{P(z_n|x_n)}{P(y_n|x_n)}$ to simplify above equation. If we use the notation $N_{(\cdot)(\cdot)}$ to represent the confusion matrix

¹Note that ¬ means removing the current object statistic from the whole collection of all object statistics.

 $^{^2} The probabilistic prediction from the model is clipped between <math display="inline">\xi$ and $1-\xi$ for the computational stability, where ξ is set to 10^{-20} in this paper.

of the noisy dataset, then we have $\sum_{k}^{K}\sum_{k'}^{K}N_{kk'}=N$ and $\prod_{n=1}^{N}\phi_{z_ny_n}=\prod_{k}^{K}\prod_{k'}^{K}\phi_{kk'}^{N_{kk'}}$. Putting the later equation into Eq. (4) and then using the conjugation characteristic between the Dirichlet distribution and the Multinomial distribution, the following form can be further deduced,

$$P(Z|X,Y;\alpha) = S * \int_{\phi} \prod_{k=1}^{K} \frac{\Gamma(\sum_{k'}^{K} \alpha_{k'})}{\prod_{k'}^{K} \Gamma(\alpha_{k'})} \prod_{k'}^{K} \phi_{kk'}^{N_{kk'} + \alpha_{k'} - 1} d\phi$$

$$= S * \prod_{k=1}^{K} \frac{\Gamma(\sum_{k'}^{K} \alpha_{k'})}{\prod_{k'}^{K} \Gamma(\alpha_{k'})} \prod_{k=1}^{K} \frac{\prod_{k'}^{K} \Gamma(\alpha_{k'} + N_{kk'})}{\Gamma(\sum_{k'}^{K} (\alpha_{k'} + N_{kk'}))}.$$
(5)

Unfortunately, Eq. (5) is non-analytical and cannot be used to generate the samples of z directly, which can be solved by Gibbs sampling. According to the Gibbs sampling, we need to compute $P(z_n|Z^{-n})$ first. And then based on $P(z_n|Z^{-n})$, a sequence of observations can be sampled, which are approximately from $P(z_n|x_n,y_n,\phi)$. The following deduction facilitates Eq. (5) and $\Gamma(x+1)=x\Gamma(x)$ to acquire the final conditional probability for our autoencoded Gibbs sampling.

$$P(z_{n}|Z^{\neg n}, X, Y; \alpha) = \frac{P(Z|X, Y; \alpha)}{P(Z^{\neg n}|X, Y; \alpha)}$$

$$= \frac{P(z_{n}|x_{n})}{P(y_{n}|x_{n})} \frac{\alpha_{y_{n}} + N_{z_{n}y_{n}}^{\neg n}}{\sum_{k'}^{K} (\alpha_{k'} + N_{z_{n}y_{n}}^{\neg n})}$$

$$\propto \underbrace{P(z_{n}|x_{n})}_{\text{Classifier encoder}} \underbrace{\frac{\alpha_{y_{n}} + N_{z_{n}y_{n}}^{\neg n}}{\sum_{k'}^{K} (\alpha_{k'} + N_{z_{n}y_{n}}^{\neg n})}}_{\text{Conditional transition}}.$$
(6)

With Eq. (6), we can sample a collection of latent labels $\{z_n\}$. Such samples are then used to solve the optimization problem in Eq. (3). Iterating the procedure of Eq. (6) and Eq. (3), we gradually approach the latent label, and at the same time train the classifier and estimate the noise transition.

Lemma 1. For a reversible, irreducible and aperiodic Markov chain with state space Ω , let λ^* be the maximal absolute eigenvalue of the state transition matrix and π be the underlying stationary probability measure where $\pi_{min} = \min_{Z \in \Omega} \pi(Z)$. Then, the ϵ -mixing time from the initial arbitrary state to the equilibrium is characterized by the following bounds,

$$\frac{\lambda^*}{1-\lambda^*} \ln\left(\frac{1}{2\epsilon}\right) \le \tau_{mix}(\epsilon) \le \frac{1}{1-\lambda^*} \ln\left(\frac{1}{\pi_{min}\epsilon}\right), \quad (7)$$

where $\tau_{mix}(\epsilon) = \min\{t : ||P_t(Z) - \pi||_{TV} \le \epsilon\}$ and $||\cdot||_{TV}$ is the total variation distance between two probability measures.

The above lemma indicates [31] the mixing time of LCCN is at most constantly linear to the inverse of $1-\lambda^*$. Although it is hard for Gibbs sampling to accurately quantify λ^* due to the evolving state transition matrix, the recent work [32] shows Gibbs sampling is efficient enough and almost proportional to the logarithm of the dataset size N. In experiments, we will show LCCN is well converged after same epochs as baselines.

In statistical learning theory, the excess risk ³ and the error bound w.r.t. the expected risk and Bayes risk, are two impor-

tant quantities to measure model generalization performance. In the setting of noisy labels, such two quantities are bounded by the following generalization bound (see the Appendix A)

$$\Delta_{\mathcal{F}} = \sup_{f_{\theta} \in \mathcal{F}} \left| \mathbf{E} \left[\ell_1(z, f_{\theta}(x)) \right] - \mathbf{E}^{(D_N)} \left[\ell_1(z, f_{\theta}(x)) \right] \right|,$$

where $\mathbf{E}\left[\cdot\right]$ and $\mathbf{E}^{(D_N)}\left[\cdot\right]$ respectively represents the expectation on the clean data distribution and the empirical estimation with the data whose labels are from the Gibbs sampling. Thus, by analyzing the upper bound of $\Delta_{\mathcal{F}}$, we can then understand which factors affect the generalization performance of LCCN. Specifically, we deduce the following theorem to interpret this.

Theorem 2. Assume f_{θ}^* and f_{θ}^{\dagger} respectively are the underlying groundtruth labeling functions $\mathcal{X} \to \mathcal{Y}$ of clean test data and data from the Gibbs sampling. Define the composite function class $\mathcal{G} = \{x \mapsto \ell_1(f_{\theta}'(x), f_{\theta}(x)) : f_{\theta}', f_{\theta} \in \mathcal{F}\}$. Then, for any probability $\delta > 0$, with probability at least $1 - \delta$,

$$\Delta_{\mathcal{F}} \le \Delta + \widehat{\mathcal{R}}(\mathcal{G}) + 3\rho \sqrt{\frac{\ln(\frac{2}{\delta})}{2N}}$$
 (8)

where $\Delta = \sup_{f_{\theta} \in \mathcal{F}} \left| \mathbf{E} \left[\ell_1(f_{\theta}^*(x) - f_{\theta}^{\dagger}(x), f_{\theta}(x)) \right] \right|, \widehat{\mathcal{R}}(\mathcal{G})$ is the Rademacher complexity [33] of \mathcal{G} and ρ is the maximum of the ξ -clipped cross entropy loss, i.e., $-\ln \xi$.

The above theorem indicates the generalization performance of the classifier learned by LCCN depends upon three factors, i.e., the inherent gap Δ between the noisy training domain and the clean test domain, the function complexity $\widehat{\mathcal{R}}(\mathcal{G})$ and the sample number N. In particular, if LCCN can exactly infer all the latent labels of noisy data and eliminate the domain bias, we will have $f_{\theta}^* = f_{\theta}^{\dagger}$ and $\Delta = 0$. In this case, Eq. (8) will degenerate to the Rademacher bound [34] after scaling the loss to [0,1], and equal to the training on the clean data. However, it is usually hard to completely remove the domain bias, since the distribution of the corrected samples could still be different from that of the clean test data. For example, the web data may contain many outlier classes. Thus, Δ is an important factor to the generalization performance of LCCN.

D. Safeguarded Transition Update

In this section, we will show that our method safeguards the bounded update of the noise transition by a batch of samples, avoiding the arbitrarily tuning via a Softmax layer in [13].

Theorem 3. Suppose α_i is a positive smoothing scalar, N_i is the current sample number of the ith category $(i=1,\ldots,K)$, M_i is the sum of the sample numbers newly allocated into (positive) and removed from (negative) the ith category after a batch of training samples, and \widehat{M}_i is its absolute sum of such two cases. Then, for the transition vector ϕ_i of the ith category, its variation via a training batch is characterized by the following equation,

$$\left|\phi_i^{new} - \phi_i^{old}\right| \le \frac{|r_i| + \widehat{r}_i}{1 + r_i} \tag{9}$$

where $r_i = \frac{M_i}{N_i + \sum_{j=1}^K \alpha_j}$ and $\widehat{r}_i = \frac{\widehat{M}_i}{N_i + \sum_{j=1}^K \alpha_j}$. According to the definition, we have $r_i > -1$, $\widehat{r}_i \geq 0$ and $\widehat{r}_i \geq |r_i|$.

³https://en.wikipedia.org/wiki/Risk_difference

Proof. The variation of ϕ_i after a training batch is,

$$\begin{aligned} &\left|\phi_{i}^{new} - \phi_{i}^{old}\right| \\ &= \sum_{j=1}^{K} \left|\phi_{ij}^{new} - \phi_{ij}^{old}\right| \\ &= \sum_{j=1}^{K} \left|\frac{N_{ij} + \alpha_{j} + M_{ij}}{N_{i} + \sum_{j'=1}^{K} \alpha_{j'} + M_{i}} - \frac{N_{ij} + \alpha_{j}}{N_{i} + \sum_{j'=1}^{K} \alpha_{j'}}\right| \\ &\leq \sum_{j=1}^{K} \frac{\left|(N_{i} + \sum_{j'=1}^{K} \alpha_{j'})M_{ij}\right| + \left|(N_{ij} + \alpha_{j})M_{i}\right|}{(N_{i} + \sum_{j'=1}^{K} \alpha_{j'})(N_{i} + \sum_{j'=1}^{K} \alpha_{j'} + M_{i})} \\ &= \frac{(N_{i} + \sum_{j'=1}^{K} \alpha_{j'})\widehat{M}_{i} + (N_{i} + \sum_{j'=1}^{K} \alpha_{j'})\left|M_{i}\right|}{(N_{i} + \sum_{j'=1}^{K} \alpha_{j'})(N_{i} + \sum_{j'=1}^{K} \alpha_{j'} + M_{i})} \\ &= \frac{\left|r_{i}\right| + \widehat{r}_{i}}{1 + r_{i}} \end{aligned}$$

Corollary 3.1. Suppose M is the batch size in the training. If it satisfies the condition $M \ll N_i$, we have $\widehat{r}_i < \frac{M}{N_i}$ in a small scale. Then the variation of ϕ_i after a training batch will be bounded by $\frac{|r_i|+\widehat{r}_i}{1+r_i} \leq \frac{2\widehat{r}_i}{1-\widehat{r}_i} \approx 2\widehat{r}_i$ in a small scale.

The core drawback in [13] is the noise transition modeled by a Softmax layer can be arbitrarily updated via a batch of samples. This is because the gradients of the parameters estimated by a "local" batch can be arbitrarily large in the backpropagation. Then, the noise transition decided by the "global" dataset might be pushed into a bad local minimum by a batch of some extremely noisy training samples, yielding a serious harm on the classifier. The later experimental analysis in Fig. 7 will confirm this point. Instead, our dynamic label regression theoretically safeguards the bounded update of the noise transition via a batch of samples. Specifically, with the bounded update, the conditional transition in Equation (6) is gradually changing towards at a true distribution when the classifier is well trained. Similarly, with more reliable sampled labels, the classifier is better trained and the noise modeling is refined. Finally, we acquire a virtuous cycle for optimization.

E. Complexity Analysis

The learning procedure is summarized in Algorithm 1. Note that, we give the complete implementation including details, like pretraining and warming-up used in the experiments.

As we know, the stochastic optimization of a DNN model involves two steps, the forward and backward computations. In each mini-batch update, its time complexity is $\mathcal{O}(M\Lambda)$, where M is the mini-batch size and Λ is the parameter size. Here, in Algorithm 1, we additionally add a sampling operation via Eq. (6) whose complexity is $\mathcal{O}(M+K^2)$ (K is the class size). Note that, the first term in the RHS of Eq. (6) has been computed in the forward procedure. Since M and K is usually significant smaller than Λ , the extra cost for the sampling is negligible compared to $\mathcal{O}(M\Lambda)$. Besides, the optimization for noise modeling in Eq. (3) can be ignored, as this only involves the normalization of a confusion matrix whose complexity is $\mathcal{O}(K^2)$. In total, since the big-O complexity of each minibatch remains the same, our method is scalable to big data.

Algorithm 1 Dynamic Label Regression for LCCN

Require: A noisy dataset $\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N$, a classifier $P(\cdot|x)$ modeled by DNN f_{θ} , warming-up steps δ , the running epoch number L and the batch-size M.

- 1: Directly pretrain the classifier f_{θ} on the noisy dataset \mathcal{D} .
- 2: Compute the warming-up noise transition matrix ϕ' .
- 3: **for** epoch i = 1 to L **do**
- 4: **for** batch j = 1 to $\lceil N/M \rceil$ **do**
 - Let step= $i \times \lceil N/M \rceil + j$ and hook a batch of samples.
 - if step $< \delta$ then
- 7: Substitute the transition in Equation (6) with ϕ' , and then sample z_n for each x_n in the batch.
- 8: else

5:

6:

- 9: Sample z_n with Equation (6) for the batch.
- 10: **end if**
- Update the confusion matrix $N_{(\cdot)(\cdot)}$ based on the existing sampling observations $\{(z_n, y_n)\}$.
- 12: Optimize Equation (3) to learn the classifier f_{θ} and estimate the noise transition matrix ϕ .
- 13: end for
- 14: **end for**
- 15: Output the classifier f_{θ} and the noise transition ϕ .

F. Extensions on Outlier and Semi-supervised Learning

In this section, we will extend the original model in Fig. 2 to the generalized version in Fig. 4, which shares the optimization procedure but is more useful in the real-world applications.

Extension on Outlier Learning: In practise, datasets collected from online websites or real-world scenarios, usually contains the open-set label noise [20]. That means data from other distributions might be disturbed as the given class samples involving in the training. Previous class-conditional noise model [12-14, 18, 35], mainly focus on the closed-set label perturbation and thus can not handle the outlier classes. For example, it is impossible to estimate the transition matrix via the mentioned two-step formulation in [12] since this requires the selection of the representative outlier samples. Similarly, learning the transition matrix by a noise adaptation layer [13] still suffers from the instability. Therefore, it is useful to extend LCCN to deal with this open-set noisy label setting. Actually, the modification for LCCN only requires to add an outlier choice for the latent variable z, i.e., $z \in \{1, ..., K, K + 1\}$, where K+1 indexes the collapsed outlier classes and then change the noise transition from \mathbf{R}^{KxK} to $\mathbf{R}^{(K+1)xK}$. The model modification is shown in Fig. 4 and the the network modification is illustrated in Fig. 3. Importantly, the above modifications do not alter the aforementioned deduction.

Extension on Semi-supervised Learning: It is common to improve the model performance by augmenting the large scale noisy dataset with a small set of clean samples. Many works [10–12, 35, 36] have leveraged such a semi-supervised setting to calibrate the classifier and achieve a better result. In our model, it is naturally compatible with this case, where we can directly utilize the clean labels instead of labels from sampling when they are available. As illustrated in Fig. 3,

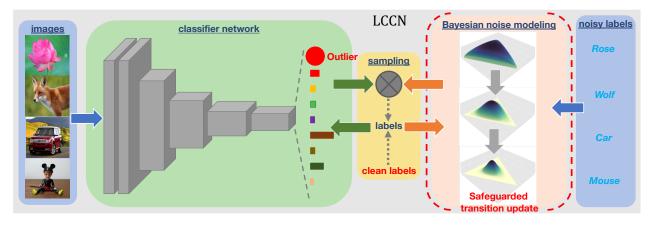


Fig. 3. Extensions based on the original safeguarded dynamic label regression for LCCN. The images and noisy labels are respectively input to the classifier and the safeguarded Bayesian noise modeling to compute the prediction (including the outlier component) and the conditional transition. When clean labels are not available as the latent labels, they are sampled based on the product of previous two quantities. Then, the latent labels composite by both the clean parts and those inferred from sampling are used to train the classifier, and only the latent labels inferred from the sampling are used to refine the noise model.

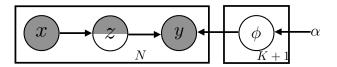


Fig. 4. The Generalized Latent Class-Conditional Noise model. x and y is the observed training pair. z is the partially observed latent label (the observed samples are for semi-supervised learning). $\phi \in \mathbb{R}^{(K+1)xK}$ is the unknown noise transition (the extra dimension is for outlier learning). α is a Dirichlet parameter. N is the sample number and K is the class number.

this configuration is specifically marked in red in parallel to sampling. The corresponding model modification is indicated in Fig. 4. In a broad sense, clean labels can be as accurate as a given category or as weak as a coarse hint that tells outlier or not. In the former case, a standard cross entropy loss can be applied to the classifier; while in the latter case, a collapsed cross entropy loss that defines on outliers *vs.* non-outliers could be applied. Since this is tightly related to the work on domain adaptation [37], we leave the latter case in the future and only validate the former semi-supervised setting in this paper.

IV. EXPERIMENTS

The experiments involve both the simulated noisy datasets and the real-world noisy datasets. We verify the performance of our model by comparing with state-of-the-art methods.

A. Datasets and Baselines

1) Datasets: We conduct the toy experiments on CIFAR-10, CIFAR-100 and the real-world experiments on Clothing1M and WebVision. CIFAR-10 and CIFAR-100 [38] respectively consist of 60,000 32x32 color images from 10 and 100 classes. Both of them contain 50,000 training samples and 10,000 test samples. For the toy experiments without outliers, we inject the asymmetric noise to disturb their labels to form the noisy datasets. Concretely, on CIFAR-10, we set a probability r to disturb the label to its similar class, i.e., truck \rightarrow automobile, bird \rightarrow airplane, deer \rightarrow horse, cat \rightarrow dog. For CIFAR-100, a similar r is set but the label flip only happens in each

super-class. The label is randomly disturbed into the next class circularly within the super-classes. For the toy experiments that consider the open-set noisy labels, we randomly select 10,000 samples from the original datasets and shuffle the order of the pixel values as the outliers. In the semi-supervised learning, we utilize the clean labels of the first 5,000 clean samples and the first 500 outlier samples for the training.

Clothing1M [35] dataset has 1 million noisy clothes samples collected from the shopping websites. The authors in [35] predefined 14 categories and assigned the clothes images with the labels extracted from the surrounding text provided by sellers, which thus might be very noisy. According to [35], only about 61.54% labels are reliable. Besides, this dataset contains 50k, 14k and 10k clean samples respectively for auxiliary training, validation and test. WebVision⁴ [39] is a more challenging noisy dataset, which contains more than 2.4 million images. It is crawled from the Internet by using the 1,000 concepts of ILSVRC [40] as queries. In addition, a clean validation set which contains 50,000 annotated images, are provided to boost and validate the proposed models in diverse applications. We use the validation set of ImageNet [40] as its test set.

2) Baselines: For the toy experiments, we compare LCCN with the classifier that is directly trained on the dataset (termed as CE), the method Bootstrapping proposed in [18], the transition based method Forward [12] and the method that fine-tunes the transition S-adaptation [13]. Note that, we choose the hard mode for Bootstrapping, since it is empirically better than the soft mode. In the outlier corrupted datasets, we denote the extension of our model that considers the outlier as LCCN*. In the semi-supervised learning setting, we denote our model as LCCN+, meaning the clean samples are used in the training. For the experiments on real-world datasets, we also report the result of Joint Optimization [26] that leverages the auxiliary noisy label distribution and the state-of-the-art result Forward+ [12] that finetunes on clean samples.

⁴Due to the reason of time and the computational resource, in this paper, we only use the original WebVision 1.0 dataset. The newest version, namely WebVision 2.0 dataset, contains more images and more classes.

	Dataset	CIFAR-10				CIFAR-100					
# Method \ Noise Ratio		0.1	0.3	0.5	0.7	0.9	0.1	0.2	0.3	0.4	0.5
$\frac{n}{1}$	CE	90.10	88.12	76.93	59.01	56.85	66.15	64.31	60.11	51.68	33.37
2	Bootstrapping	90.73	88.12	76.29	57.04	56.79	66.48	64.61	63.01	55.27	34.52
3	Forward	90.86	89.03	82.47	67.11	57.29	65.43	62.72	61.28	52.64	33.82
4	S-adaptation	91.02	88.83	86.79	72.74	60.92	65.52	64.11	62.39	52.74	30.07
5	LCCN	91.35	89.33	88.41	79.48	64.82	67.83	67.63	66.86	65.52	33.71
6	CE with the clean data	91.63				69.41					

TABLE I THE AVERAGE ACCURACY (%) OVER 5 TRIALS ON CIFAR-10 AND CIFAR-100 WITH DIFFERENT NOISE LEVELS.

TABLE II

The average accuracy (%) over 5 trials on outlier-corrupted CIFAR-10 and CIFAR-100 with different noise levels.

	Dataset	CIFAR-10				CIFAR-100					
#	Method \ Noise Ratio	0.1	0.3	0.5	0.7	0.9	0.1	0.2	0.3	0.4	0.5
1	CE	89.13	87.06	74.63	62.29	57.07	62.94	59.73	54.71	45.57	31.74
2	Bootstrapping	90.13	84.58	74.76	54.87	55.56	63.73	60.88	59.77	40.23	31.86
3	Forward	88.63	84.97	78.47	58.23	56.52	63.69	62.63	61.86	51.47	35.71
4	S-adaptation	88.58	87.28	61.17	57.12	56.73	63.51	61.50	60.59	53.22	32.19
5	LCCN	88.63	88.06	82.15	69.48	55.12	63.97	62.84	61.79	60.34	33.52
6	LCCN*	89.59	88.43	84.34	72.33	56.28	64.71	63.05	62.48	62.02	32.37
7	LCCN+	90.30	88.93	88.21	87.42	86.33	65.67	64.24	63.52	63.19	62.39

B. Implementation Details

For CIFAR-10 and CIFAR-100, the PreAct ResNet-32 [41] is adopted as the classifier. The image data is augmented by horizontal random flip and 32×32 random crops after padding with 4 pixels. Then, the per-image standardization is used to normalize pixel values. For the optimizer, we utilize SGD with a momentum of 0.9 and a weight decay of 0.0005. The batch size is set to 128. The training runs totally 120 epochs and is divided into three phases in 40 and 80 epochs. In these three phases, we respectively set the learning rate as 0.5, 0.1 and 0.01. Note that, the reason that we adopt a large learning rate (others may set the learning rate smaller than 0.001), is that the small learning rate will lead to overfitting on the noisy dataset as claimed in [8]. Following the benchmark in [12], we use CE to initialize the classifier in other baselines and LCCN. For S-adaptation, the following transition is computed to warm-up the transition parameters in the first 80 epochs.

$$\phi'_{ij} = \frac{\sum_{t} 1_{y_t = j} p(z_t = i | x_t)}{\sum_{t} p(z_t = i | x_t)}$$
(11)

Similarly on CIFAR-10, we use above transition to warm up the sampling procedure in LCCN for the first 20,000 steps. However, on CIFAR-100, we set $\phi'_{ij}=1[i=j]$ in the warming-up since Equation (11) will induce the high sampling variance and need long time to converge.

For Clothing1M and WebVision, the ResNet-50 is leveraged as the classifier. We resize the short side of their images to 224 and do the random crop of 224×224 . The training images are augmented with the random flip, whiteness and saturation. For the optimizer, we deploy SGD with a momentum of 0.9 with a weight decay of 10^{-3} . The batch size for Clothing1M is set to 32 and we fix the learning rate as 0.001 to run 5 epochs. For the warming-up transition, we both validate the one [35] from

manual annotation and the one estimated by Equation (11) for 40,000 steps. Note that, on the large real-world datasets, due to the strong capacity of ResNet-50, it is easy for LCCN to occur the sampling collapsed problem, i.e., the sampled latent label is identical to the noisy label. Thus, we norm Equation (6) with a power annealed coefficient $\max\{\exp\left(-\frac{\text{step}}{\max_s \text{step}} * 0.8\right), 0.5\}$ to introduce the sufficient perturbation in avoid of this issue. On WebVision, the batch size is set to 128, and the learning rate is initialized with 0.1 is divided by 10 every 30 epochs until 90 epochs. We use the diagonal transition for 10,000 steps of warming-up and then update the confusion matrix to the end, since it contains 1,000 categories. The similar power-annealed strategy for sampling is leveraged. Finally, to fairly compare LCCN+ and Forward+ in semi-supervised learning, we use the similar fine-tuning in [12] to run the experiments.

C. Results on CIFAR10 and CIFAR-100

1) Classification experiments: Table I summarizes the performance of LCCN and baselines on the noisy datasets without outilers. Compared with the baselines, LCCN achieves the best performance at most of noise levels. In particular, even in the large noise rates, our model still acquires the competitive classification accuracy. For example, when r=0.7 on CIFAR-10 and r=0.4 on CIFAR-100, LCCN reaches 79.48% and 65.52%, outperforming the best results of baselines by about 7% and 13% respectively. This demonstrates that our model is significantly better than baselines. Regarding r=0.5 on CIFAR-100, the way to disturb the labels [12] leads that there is one undesired minimum, since two classes are mixed into one class by equal quota after injecting noise. In this case, it is hard to say which model can achieve the best result. We include it here for the complete comparison as in other works [12, 13].

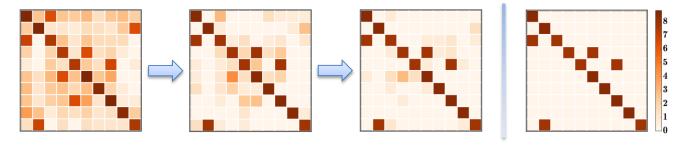


Fig. 5. The colormap of the confusion matrix on CIFAR-10 with r=0.5. We utilize the log-scale for each element in the confusion matrix for the fine-grained visualization. The left three maps are respectively learned by LCCN at the beginning, 30,000 step and the end, and the right one is the groundtruth.

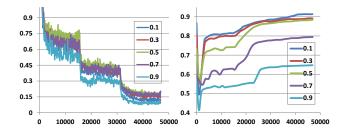


Fig. 6. The training loss (left) and the test accuracy (right) of LCCN on the CIFAR-10 dataset with different noise rates r = 0.1, 0.3, 0.5, 0.7, 0.9.

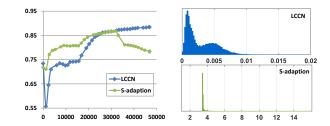


Fig. 7. The test accuracy of LCCN and S-adaptation in the training on CIFAR-10 with r=0.5 (left), and the corresponding histograms for the change of noise transition ϕ via a mini-batch of samples (right).

Table II presents the results of LCCN and baselines on the outlier-corrupted datasets. Compared to those in Table I, all the methods have a slight performance drop. Nevertheless, LCCN achieved the best performance in such an setting at r=0.3, 0.5,0.7 on CIFAR-10 and r=0.1, 0.2, 0.3 and 0.4 on CIFAR-100. The baselines without the outlier detection mechanism usually have a significant degeneration. Specifically, on CIFAR-10, all the other baselines are even not better than CE, i.e., directly training. Instead, LCCN* that considers the outlier achieves a further improvement based on LCCN. Besides, as expected, by adding clean data, LCCN+ performs better than LCCN*, since the classifier is calibrated by the information of the clean data domain. In the scenario of the extreme noise, as marked by the grey color in Table II, this is quite useful and even necessary to guarantee an acceptable performance. In total, according to the quantitative analysis of Table I and Table II, we demonstrate the superiority of LCCN compared to baselines on toy datasets.

2) On convergence visualization: In Fig. 6, we trace the training of LCCN on CIFAR-10 to visualize its convergence. As can been in the left panel of Fig. 6, LCCN has a stable convergence on loss after the given epochs. Besides, we find the loss converges to irregular scales in different noise rates. Concretely, in most cases, i.e., r=0.1, 0.3, 0.5, the final training loss increases as r increases, while the loss shows attenuation as r >0.5. It is because in the low-level noise, the model can easily correct the labels via the sampling in LCCN, yielding a small loss. While in the case of the extreme noise, it is more challenging to prevent the model from fitting on noise, which incurs difficulty for optimization and thus achieves a big loss. Furthermore, according to the right penal of Fig. 6, the test accuracy also approximately converges and persists to the end of the training without performance drop. Actually, it is not a

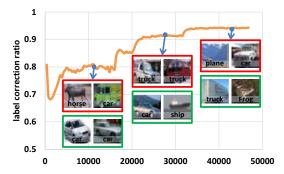


Fig. 8. The label correction ratio in the training of LCCN on CIFAR-10 with r=0.5 as well as some negatively corrected samples (the red box) and some positively corrected samples (the green box) with the high probabilities.

common phenomenon for previous methods to own this merit, since all baselines tends to overfitting on noise more or less in the final few epochs. This demonstrates the advantages of LCCN in the robust training with the noisy datasets.

3) Safeguarded transition update: To show LCCN safeguards the noise transition update compared to S-adaptation, we compute the statistics about their update of noise transition on CIFAR-10 at r=0.5, and illustrate the histogram of changes in Fig. 7. Firstly, from the left panel of Fig. 7, we can see that there is a significant performance drop in the training of S-adaptation. The clue to this phenomenon can be found by inspecting the update of noise transition. As shown in the right panel of Fig. 7, the change magnitude of ϕ in S-adaptation is higher than that of LCCN. One is in a large scale ranging from 0 to 16, while the other one is in a very small scale ranging from 0 to 0.02. This leads to S-adaptation suffering from a high risk of over-tuning to undesired local minimums in the



Fig. 9. Some representative samples in the training set that are considered as the outliers by LCCN*. We intuitively summarize these photos into four categories based on their contents, multiple different objects (RED), implicit categories (GREEN), uncertain types (BLACK) and confusing appearance (BLUE), which are respectively marked by the color of the surrounded boxes. Outliers are relative to LCCN and may contain hard examples of the pre-defined 14 categories.

 $TABLE \; III \\ The average accuracy over 5 trials on Clothing 1 M. \\$

#	Method	Accuracy
1	CE	68.94
2	Bootstrapping	69.12
3	Forward	69.84
4	S-adaptation	70.36
5	Joint Optimization	72.16
	LCCN	71.63
6	LCCN warmed-up by ϕ in [35]	73.07
	LCCN*	72.80
7	CE on the clean data	75.28
8	Forward+	80.38
9	LCCN+	81.25

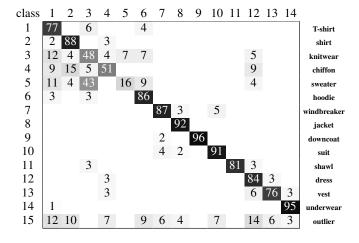
presence of noise. Instead, according to the histogram, LCCN updates ϕ in a safeguarded small scale when approaching to the minimum. In summary, this quantitative analysis confirms the claim of our Theorem 3 in the perspective of experiments.

4) The latent label and noise analysis: Fig. 5 and Fig. 8 respectively depict the colormap of the confusion matrix and the label correction ratio when training LCCN on CIFAR-10 with r=0.5. First, as can be seen in Fig. 5, the initial confusion matrix does not approach the true matrix and there are many incorrect entries. However, with the training progressing, the matrix is gradually corrected and at the end of training, it is approximately similar to the provided groundtruth. Besides, As shown in Fig. 8, the ratio of the image with the correct label increases along with the training progress. This reflects LCCN successfully models the class-conditional noise and gradually infer the latent labels. Specially, by visualizing the mis-corrected examples in the training process, we can find that the classifier at first make mistakes in even some simple samples, while finally has the wrong classification in only the hard examples. These two figures visualize how the dynamic label regression optimizes LCCN to infer the latent label and model the noise.

D. Results on Clothing 1M and WebVision

Table III lists the performance of LCCN and baselines on the large-scale Clothing1M. According to the results, we can see that Forward does not show the significant improvement

TABLE IV THE LEARNED NOISE TRANSITION ON CLOTHING $1\,\mathrm{M}$ BY LCCN.



in this dataset, even though they use the annotated noise transition matrix [35]. And S-adaptation only improves Forward by 0.5%. Joint Optimization that trains the classifier with label correction [26] achieves better results than the other baselines. Nevertheless, this method requires the provided noisy label distribution to prevent degeneration and is not scalable to the large number of classes [26]. Instead, LCCN that contains both the label correction and Bayesian noise modeling, gets the competitive performance 71.63%. With the warming-up of the auxiliary noise transition [35], it further achieves the best 73.07%. Even although there is no auxiliary information available, our extension LCCN* still outperforms the current stateof-the-art result. This demonstrates the potential of LCCN in handling the real-world noisy dataset. In addition, the results of LCCN+ indicates our model also has the advantages in the semi-supervised learning setting.

In Table IV, we present the noise transition matrix learned by LCCN*, where only significant transition probabilities are marked. From this noise transition, we can find the training samples in some classes are very noisy. For example, knitwear and sweater are two classes which transit most of labels to other classes. This is because such two classes usually occur with other categories like jacket or shawl in the common dress collocation, which may incur the label transition according to

the visual appearance. Besides, in Table IV, we can observe the outlier transition to find which class contains a lot of outliers. Furthermore, to better understand the outliers, we give some represenative samples in Fig. 9. According to the image contents, we intuitively summarize the outlier into four sub-classes, multiple different objects, implicit categories, uncertain types and confusing appearance. As can be seen, it is usually improper to asign an unique label to these outliers, since some may contain multiple kinds of clothes. And in some challenging cases, e.g., the images in the blue box in Fig. 9, the hard example is also considered as the outliers by LCCN*, even if the label is correct. Actually, this can be seen as the imperfect sample for training or the potential drawback of our model that requires more explore in the future research.

 $\label{table v} TABLE\ V$ The average accuracy over 5 trials on WebVision.

#	Method	Accuracy@1	Accuracy@5
1	CE	58.61	80.94
2	Bootstrapping	58.48	80.81
3	Forward	58.93	81.06
4	S-adaptation	58.00	80.16
5	LCCN	58.73	81.25
	LCCN*	59.09	81.33
6	CE on the clean data	53.52	77.84
7	LCCN+	59.72	80.34

Table V gives the performance of LCCN and baselines on a more challenging noisy dataset WebVision. Both Top-1 and Top-5 accuracies are reported in the experiment. According to the results either in the perspective of Top-1 accuracy or Top-5 accuracy, LCCN achieves the best performance. Similarly, LCCN* and LCCN+ respectively have the further refinement based on LCCN. Nevertheless, as can be seen, the results of all methods do not present the significant gap. One possible explanation is that this task couples two challenging sub-tasks, perfectly decoupling the clean samples from the noisy dataset in the 1000 classes and well fitting the clean samples in the 1000 classes. From the current limiting performance of image recognition in ImageNet [40], we know either sub-task mentioned above needs a long way to go for a satisfying result in such a large-scale challenging scenario.

V. CONCLUSION AND FUTURE WORK

In this paper, we present a Latent Class-Conditional Noise model to learn with the noisy supervision. Besides, a dynamic label regression method is deployed for LCCN to iteratively infer the latent labels and jointly train the classifier and model the noise. The theoretical analysis on the model convergence and the essential gap between training and test are also provided. Most importantly, we demonstrate that our method safeguards the bounded update of the noise transition to avoid previous arbitrarily tuning via a mini-batch of samples. Finally, we generalize our model to the open-set noisy labels setting and the semi-supervised learning setting. However, although we have shown the advantages of LCCN in a range of experiments with the generalized noisy supervision, other specific settings that considers more complex noise, e.g., image content

based noise, could be explored. Besides, it is important to explore more effective models on the large scale noisy dataset. To the end, more works based on LCCN can be extended to train with noisy datasets.

APPENDIX A

Proof of $\Delta_{\mathcal{F}}$ regarding to the two quantities

By definition, the excess risk of the estimated model \hat{f}_{θ} is directly bounded by the absolute supremum $\Delta_{\mathcal{F}}$ as follows,

$$\mathbf{E}\left[\ell_{1}(z,\hat{f}_{\theta}(x))\right] - \mathbf{E}^{(D_{N})}\left[\ell_{1}(z,\hat{f}_{\theta}(x))\right]$$

$$\leq \sup_{f_{\theta}\in\mathcal{F}}\left|\mathbf{E}\left[\ell_{1}(z,f_{\theta}(x))\right] - \mathbf{E}^{(D_{N})}\left[\ell_{1}(z,f_{\theta}(x))\right]\right|.$$
(12)

Assume the Bayes optimal classifier is f_{θ}^* . Since \hat{f}_{θ} minimizes the empirical loss on D_N , we have the following inequality,

$$\mathbf{E}^{(D_N)}\left[\ell_1(z,f_\theta^*(x))\right] - \mathbf{E}^{(D_N)}\left[\ell_1(z,\hat{f}_\theta(x))\right] \geq 0.$$

Then, for the error bound w.r.t. the expected risk and the Bayes risk, we can deduce with above inequalities as follows,

$$\mathbf{E}\left[\ell_{1}(z,\hat{f}_{\theta}(x))\right] - \mathbf{E}\left[\ell_{1}(z,f_{\theta}^{*}(x))\right]$$

$$\leq \mathbf{E}\left[\ell_{1}(z,\hat{f}_{\theta}(x))\right] - \mathbf{E}^{(D_{N})}\left[\ell_{1}(z,\hat{f}_{\theta}(x))\right]$$

$$-\left(\mathbf{E}\left[\ell_{1}(z,f_{\theta}^{*}(x))\right] - \mathbf{E}^{(D_{N})}\left[\ell_{1}(z,f_{\theta}^{*}(x))\right]\right)$$

$$\leq 2 * \sup_{f_{\theta} \in \mathcal{F}}\left|\mathbf{E}\left[\ell_{1}(z,f_{\theta}(x))\right] - \mathbf{E}^{(D_{N})}\left[\ell_{1}(z,f_{\theta}(x))\right]\right|.$$
(13)

Eq. (12) and Eq. (13) show the supremum of both the excess risk and the error bound are characterized by $\Delta_{\mathcal{F}}$. Thus, we can universally analyze the corresponding upper bound of $\Delta_{\mathcal{F}}$ to investigate the generalization performance of LCCN.

APPENDIX B PROOF OF THEOREM 2

Assume f_{θ}^* and f_{θ}^{\dagger} are the underlying groundtruth labeling functions $\mathcal{X} \to \mathcal{Y}$ of clean test data and data from the Gibbs sampling respectively. Then, the $\Delta_{\mathcal{F}}$ can be reformulated by applying these notations and further deduced as follows,

$$\sup_{f_{\theta} \in \mathcal{F}} \left| \mathbf{E} \left[\ell_{1}(z, f_{\theta}(x)) \right] - \mathbf{E}^{(D_{N})} \left[\ell_{1}(z, f_{\theta}(x)) \right] \right|$$

$$= \sup_{f_{\theta} \in \mathcal{F}} \left| \mathbf{E} \left[\ell_{1}(f_{\theta}^{*}(x), f_{\theta}(x)) \right] - \mathbf{E}^{(D_{N})} \left[\ell_{1}(f_{\theta}^{\dagger}(x), f_{\theta}(x)) \right] \right|$$

$$\leq \sup_{f_{\theta} \in \mathcal{F}} \left| \mathbf{E} \left[\ell_{1}(f_{\theta}^{*}(x), f_{\theta}(x)) \right] - \mathbf{E} \left[\ell_{1}(f_{\theta}^{\dagger}(x), f_{\theta}(x)) \right] \right|$$

$$+ \sup_{f_{\theta} \in \mathcal{F}} \left| \mathbf{E} \left[\ell_{1}(f_{\theta}^{*}(x), f_{\theta}(x)) \right] - \mathbf{E}^{(D_{N})} \left[\ell_{1}(f_{\theta}^{\dagger}(x), f_{\theta}(x)) \right] \right|$$

$$\leq \sup_{f_{\theta} \in \mathcal{F}} \left| \mathbf{E} \left[\ell_{1}(f_{\theta}^{*}(x), f_{\theta}(x)) \right] - \mathbf{E}^{(D_{N})} \left[\ell_{1}(f_{\theta}^{\prime}(x), f_{\theta}(x)) \right] \right|$$

$$+ \sup_{f_{\theta}, f_{\theta}^{\prime} \in \mathcal{F}} \left| \mathbf{E} \left[\ell_{1}(f_{\theta}^{\prime}(x), f_{\theta}(x)) \right] - \mathbf{E}^{(D_{N})} \left[\ell_{1}(f_{\theta}^{\prime}(x), f_{\theta}(x)) \right] \right|$$

$$+ \underbrace{\sum_{f_{\theta}, f_{\theta}^{\prime} \in \mathcal{F}} \left| \mathbf{E} \left[\ell_{1}(f_{\theta}^{\prime}(x), f_{\theta}(x)) \right] - \mathbf{E}^{(D_{N})} \left[\ell_{1}(f_{\theta}^{\prime}(x), f_{\theta}(x)) \right] \right|}_{\Delta_{disc}}$$

$$(14)$$

The second term Δ_{disc} in the right-hand side of Eq. (14) is the popular discrepancy distance. It has been demonstrated by the following Rademacher bound [34] for any probability $\delta > 0$,

$$\Delta_{disc} \le \widehat{\mathcal{R}}(\mathcal{G}) + 3\rho \sqrt{\frac{\ln(\frac{2}{\delta})}{2N}},$$

where \mathcal{G} is defined by the composite functional class $\{x \mapsto \ell_1(f'_{\theta}(x), f_{\theta}(x)) : f'_{\theta}, f_{\theta} \in \mathcal{F}\}$ and N is the sample number. Then, combined with the given Rademacher bound, we finally proof the Theorem 2. i.e., for any probability δ , the generalization bound of the models for learning with noisy labels is

$$\Delta_{\mathcal{F}} \le \Delta + \widehat{\mathcal{R}}(\mathcal{G}) + 3\rho \sqrt{\frac{\ln(\frac{2}{\delta})}{2N}}.$$
(15)

This bound theoretically points out three important factors to affect our model generalization performance, i.e., domain gap, the function complexity and the support sample number.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012.
- [2] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *NIPS*, 2014.
- [3] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath *et al.*, "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal processing magazine*, vol. 29, no. 6, pp. 82–97, 2012.
- [4] M. Hu, Y. Yang, F. Shen, L. Zhang, H. T. Shen, and X. Li, "Robust web image annotation via exploring multi-facet and structural knowledge," *IEEE Transactions on Image Processing*, vol. 26, no. 10, pp. 4871–4884, Oct 2017.
- [5] Y. Zhang, X. Chen, J. Li, W. Teng, and H. Song, "Exploring weakly labeled images for video object segmentation with submodular proposal selection," *IEEE Transactions on Image Processing*, vol. 27, no. 9, pp. 4245–4259, Sep. 2018.
- [6] J. Yang, X. Sun, Y. Lai, L. Zheng, and M. Cheng, "Recognition from web data: A progressive filtering approach," *IEEE Transactions on Image Processing*, vol. 27, no. 11, pp. 5303–5315, Nov 2018.
- [7] J. Yao, J. Wang, I. W. Tsang, Y. Zhang, J. Sun, C. Zhang, and R. Zhang, "Deep learning from noisy image labels with quality embedding," *IEEE Transactions on Image Processing*, vol. 28, no. 4, pp. 1909–1922, April 2019.
- [8] D. Arpit, S. Jastrzebski, N. Ballas, D. Krueger, E. Bengio, M. S. Kanwal, T. Maharaj, A. Fischer, A. Courville, Y. Bengio *et al.*, "A closer look at memorization in deep networks," in *ICML*, 2017.
- [9] X. Ma, Y. Wang, M. E. Houle, S. Zhou, S. M. Erfani, S.-T. Xia, S. Wijewickrema, and J. Bailey, "Dimensionality-driven learning with noisy labels," in *ICML*, 2018.
- [10] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in *ICML*, 2018.

- [11] S. Guo, W. Huang, H. Zhang, C. Zhuang, D. Dong, M. R. Scott, and D. Huang, "Curriculumnet: Weakly supervised learning from large-scale web images," in *The European Conference on Computer Vision (ECCV)*, September 2018.
- [12] G. Patrini, A. Rozza, A. K. Menon, R. Nock, and L. Qu, "Making deep neural networks robust to label noise: A loss correction approach," in CVPR, 2017.
- [13] J. Goldberger and E. Ben-Reuven, "Training deep neural-networks using a noise adaptation layer," *ICLR*, 2017.
- [14] S. Sukhbaatar, J. Bruna, M. Paluri, L. Bourdev, and R. Fergus, "Training convolutional networks with noisy labels," *ICLR*, 2014.
- [15] I. Misra, C. L. Zitnick, M. Mitchell, and R. Girshick, "Seeing through the Human Reporting Bias: Visual Classifiers from Noisy Human-Centric Labels," in CVPR, 2016.
- [16] B. Han, J. Yao, G. Niu, M. Zhou, I. W. Tsang, Y. Zhang, and M. Sugiyama, "Masking: A new perspective of noisy supervision," in *NIPS*, 2018.
- [17] T. Liu and D. Tao, "Classification with noisy labels by importance reweighting," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 38, no. 3, pp. 447–461, 2016.
- [18] S. Reed, H. Lee, D. Anguelov, C. Szegedy, D. Erhan, and A. Rabinovich, "Training deep neural networks on noisy labels with bootstrapping," *ICLR*, 2014.
- [19] Y. Li, J. Yang, Y. Song, L. Cao, J. Luo, and L.-J. Li, "Learning from noisy labels with distillation." in *ICCV*, 2017.
- [20] Y. Wang, W. Liu, X. Ma, J. Bailey, H. Zha, L. Song, and S.-T. Xia, "Iterative learning with open-set noisy labels," in CVPR, 2018.
- [21] L. Jiang, Z. Zhou, T. Leung, L.-J. Li, and L. Fei-Fei, "Mentornet: Regularizing very deep neural networks on corrupted labels," in CVPR, 2018.
- [22] G. Song and W. Chai, "Collaborative learning for deep neural networks," in *Advances in Neural Information Processing Systems 31*, 2018, pp. 1837–1846.
- [23] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. W. Tsang, and M. Sugiyama, "Co-teaching: Robust training deep neural networks with extremely noisy labels," in *NIPS*, 2018.
- [24] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," *ICLR*, 2016.
- [25] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *ICLR*, 2017.
- [26] D. Tanaka, D. Ikami, T. Yamasaki, and K. Aizawa, "Joint optimization framework for learning with noisy labels," in CVPR, 2018.
- [27] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *ICLR*, 2014.
- [28] A. Mnih and K. Gregor, "Neural variational inference and learning in belief networks," in *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32*, 2014, pp. II–1791–II–

1799.

- [29] E. Jang, S. Gu, and B. Poole, "Categorical reparameterization with gumbel-softmax," *ICLR*, 2017.
- [30] C. J. Maddison, A. Mnih, and Y. W. Teh, "The concrete distribution: A continuous relaxation of discrete random variables," *ICLR*, 2017.
- [31] D. A. Levin and Y. Peres, *Markov chains and mixing times*. American Mathematical Soc., 2017, vol. 107.
- [32] J. Johan, "Fast mixing for latent dirichlet allocation," arXiv preprint arXiv:1701.02960v2, 2017.
- [33] P. L. Bartlett and S. Mendelson, "Rademacher and gaussian complexities: Risk bounds and structural results," *Journal of Machine Learning Research*, vol. 3, no. Nov, pp. 463–482, 2002.
- [34] Y. Mansour, M. Mohri, and A. Rostamizadeh, "Domain adaptation: Learning bounds and algorithms," COLT, 2009.
- [35] T. Xiao, T. Xia, Y. Yang, C. Huang, and X. Wang, "Learning from massive noisy labeled data for image classification," in *CVPR*, 2015.
- [36] A. Veit, N. Alldrin, G. Chechik, I. Krasin, A. Gupta, and S. J. Belongie, "Learning from noisy large-scale datasets with minimal supervision." in *CVPR*, 2017.
- [37] G. Csurka, "Domain adaptation for visual applications: A comprehensive survey," *CoRR*, vol. abs/1702.05374, 2017.
- [38] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," Citeseer, Tech. Rep., 2009.
- [39] W. Li, L. Wang, W. Li, E. Agustsson, and L. Van Gool, "Webvision database: Visual learning and understanding from web data," *arXiv* preprint arXiv:1708.02862, 2017.
- [40] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in CVPR09, 2009.
- [41] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.

Michael Shell Biography text here.

John Doe Biography text here.

PLACE PHOTO HERE