

Learning from Crowds by Modeling Common Confusions

Zhendong Chu, Jing Ma, Hongning Wang

Department of Computer Science, University of Virginia
{zc9uy, jm3mr, hw5x}@virginia.edu

Abstract

Crowdsourcing provides a practical way to obtain large amounts of labeled data at a low cost. However, the annotation quality of annotators varies considerably, which imposes new challenges in learning a high-quality model from the crowdsourced annotations. In this work, we provide a new perspective to decompose annotation noise into *common noise* and *individual noise* and differentiate the source of confusion based on instance difficulty and annotator expertise on a per-instance-annotator basis. We realize this new crowdsourcing model by an end-to-end learning solution with two types of noise adaptation layers: one is shared across annotators to capture their commonly shared confusions, and the other one is pertaining to each annotator to realize individual confusion. To recognize the source of noise in each annotation, we use an auxiliary network to choose from the two noise adaptation layers with respect to both instances and annotators. Extensive experiments on both synthesized and real-world benchmarks demonstrate the effectiveness of our proposed common noise adaptation solution.

Introduction

The availability of large amounts of labeled data is often a prerequisite for applying supervised learning solutions in practice. Crowdsourcing makes it possible to collect massive labeled data in both time- and cost-efficient manner (Buecheler et al. 2010). However, because of varying and unknown expertise of annotators, crowdsourced labels are usually noisy, which naturally lead to an important research problem: *how to train an accurate learning model with only crowdsourced annotations?*

The first step to estimate an accurate learning model from crowdsourced annotations is to properly model the generation of such data. In this work, we focus on the crowdsourced classification problem. The seminal work from Dawid and Skene (1979) (known as the DS model) assumes that each annotator has his/her own class-dependent confusion when providing annotations to instances. This is modeled by an annotator-specific confusion matrix, whose entries are the probability of flipping one class into another. The DS model has become the cornerstone of most learning from crowds solutions; and mainstream solutions perform label aggregation prior to classifier training: their key difference lies on

different label aggregation methods based on the DS model (Venzani et al. 2014; Zhang et al. 2014; Whitehill et al. 2009). Recent developments focus more on unified solutions, where variants of the Expectation-Maximization (EM) algorithm are proposed to integrate label aggregation and classifier training (Albarqouni et al. 2016; Cao et al. 2019; Raykar et al. 2010). Typically, such solutions treat the classifier’s predictions as latent variables, which are then mapped to the observed crowdsourced labels using individual confusion matrices of annotators. Rodrigues and Pereira (2018) further fuse label inference and classifier training in an end-to-end approach using neural networks, where the gradient from label aggregation is directly propagated to estimate the annotators’ confusion matrices. Tanno et al. (2019) propose a similar solution but encourage the annotator confusion matrix to be close to an identity matrix by trace regularization.

All existing DS-model-based solutions assume noise in crowdsourced labels is only caused by individual annotators’ expertise. However, it is not uncommon that different annotators would share common confusion about the labels. For example, when a *bird* in an image is too small, every annotator has a chance to confuse it with an *airplane* because of the background sky. We hypothesize that on an instance the annotator is confident about, he/she is more likely to use his/her expertise to provide a label (i.e., introducing individualized noise), while he/she would use common sense to label those unconfident ones. We empirically evaluate this hypothesis on two public crowdsourcing datasets, one for image labeling and one for music genre classification (more details of the datasets can be found in the Experiment Section), and visualize the results in Figure 1. On both datasets, there are quite some commonly made mistakes across annotators. For example, on the image labeling dataset LabelMe, 61.0% annotators mistakenly labeled *street* as *inside city* and 44.1% of them mislabeled *open country* as *forest*; on the music classification dataset, 63.6% annotators mislabeled *metal* as *rock* and 38.6% of them mislabeled *disco* as *pop*. The existence of such shared confusions across annotators directly affects label aggregation: the majority of annotators are not necessarily correct, as their mistakes are no longer independent (e.g., those large off-diagonal entries in Figure 1). This is against the fundamental assumption in the DS model, and strongly urges new noise modeling to better handle real-world crowdsourced data.

Moving beyond the independent noise assumption in the

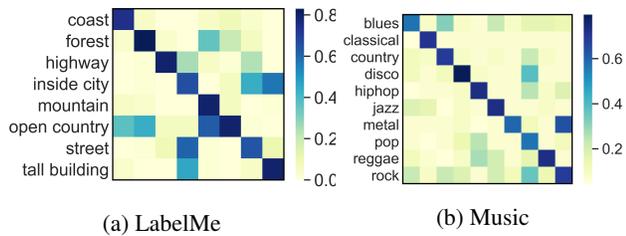


Figure 1: Analysis of commonly made mistakes across annotators on two real-world crowdsourcing datasets. The value of each entry in the heatmap denotes the percentage of annotators with this confusion pair (e.g., mistakenly label *street* as *inside city* on LabelMe dataset).

family of DS models (Dawid and Skene 1979; Rodrigues and Pereira 2018), we decompose annotation noise into two sources, *common noise* and *individual noise*, and differentiate the source of noise based on both annotators and instances. We refer to the annotation confusions shared across annotators as common noise, and model it by a global confusion matrix shared by all annotators. In the meanwhile, we also maintain annotator-specific confusion matrices for individual noise modeling. We still treat ground-truth labels of instances as latent variables, but map them to noisy annotations by two parallel confusion matrices, to capture these different sources of noise. We determine the choice of confusion matrices on a per-instance-annotator basis, by explicitly modeling of annotator expertise and instance difficulty (Whitehill et al. 2009; Yin et al. 2017). To leverage the power of representation learning to model annotator expertise and instance difficulty, we realize all our model components using neural networks. In particular, we model the two types of confusion matrices as two parallel noise adaptation layers (Goldberger and Ben-Reuven 2016). For each annotator-instance pair, the classifier first maps the instance to a latent class label, then an auxiliary network decides which noise adaptation layer to map the latent class label to the observed annotation. Cross-entropy loss is counted on the predicted annotations for end-to-end training of these components. We name this approach *CoNAL* - learning from crowds with **common noise adaptation layers**. Extensive experiments show considerable improvement of our new noise modeling approach against a rich set of baselines on two synthesized datasets, including a fully synthesized dataset and one based on CIFAR-10 dataset with various settings of noise generation, as well as two real-world datasets, e.g., LabelMe for image classification, and Music for music genre classification.

Related Works

Several existing studies focused on modeling the different roles of instance and annotator in crowdsourced data. Whitehill et al. (2009) model the accuracy of each annotation, which depends on instance difficulty and annotator expertise, to weigh each instance in final majority vote. Welinder et al. (2010) model each annotator as a multi-dimensional classifier and consider instance difficulty as single dimension latent variable. Zhou et al. (2012) propose a minimax

entropy principle on a probability distribution over annotators, instances and annotations, in which by minimizing entropy instance confusability and annotator expertise are naturally inferred. Khetan and Oh (2016) and Shah, Balakrishnan, and Wainwright (2016) consider generalized DS models which model the instance difficulty. Instead of simply using a single scalar to model instance difficulty and annotator expertise as in previous works, we model them by learning their corresponding representations via an auxiliary network, which can better capture the shared statistical pattern across observed annotations.

Our method is closely related to several existing DS-based models considering relations among annotators; but it is also clearly distinct from them. Kamar, Kapoor, and Horvitz (2015) use a global confusion matrix to capture the identical mistakes by all annotators, and it is designed to replace the individual matrix when observations of an annotator are rare. Moreover, the choice of confusion matrix in this solution only depends on the number of annotations an annotator provided. This unnecessarily reflects the annotator expertise, as the task assignment is typically out of their control in crowdsourcing. Venanzi et al. (2014) and Imamura, Sato, and Sugiyama (2018) cluster annotators to generate their own confusion matrices from a shared community-wide confusion matrix. However, the above approaches still assume a single underlying noise source, and thus they do not consider the difference between global (or community-level) and individual confusions. Li, Rubinstein, and Cohn (2019) explore the correlation of annotation across annotators by classifying them into auxiliary subtypes under different ground-truth classes. However, the characteristics of each annotator are missing since they are only represented by a specific subtype. In our work, we still characterize individual annotators by modeling their own confusions.

Common Confusion Modeling in Crowdsourced Data

In this section, we formulate our problem-solving framework for training classifiers directly from crowdsourced labels, based on the insight of common confusion modeling across annotators. We first describe the notations and our probabilistic modeling of the noisy annotation process, considering both common and individual confusions. This probabilistic model of noisy annotations is the basis of the end-to-end neural solution we develop in this paper.

Notations and Probabilistic Modeling

Assume we have N instances labeled by R annotators out of C possible classes. We define \mathbf{x}_i as the feature vector of the i -th instance and y_i^r as its label provided by the r -th annotator. Denote z_i as the unobservable ground-truth label for the i -th instance, which is considered as a latent variable sampled from a multinomial distribution parameterized by $\{p(z_i = c | \mathbf{x}_i)\}_{c=1}^C$. For simplicity, we collectively define $X = \{\mathbf{x}_i\}_{i=1}^N$, $Y = \{y_i^r\}_{i=1, r=1}^{N, R}$ and $Z = \{z_i\}_{i=1}^N$. The final goal of learning from crowds is to obtain the classifier $P(Z|X)$ only with crowdsourced annotations Y .

Similar to the DS-based models (see Figure 2a for reference), the confusion of the r -th annotator is measured

by an annotator-specific confusion matrix π^r , in which the (z, z') -element $\pi_{z, z'}^r$ denotes the probability that annotator r will label the true label z as z' . Aside from individual confusion, the key assumption of our solution is that annotation mistakes can also be introduced by common confusion, which is modeled by a globally shared confusion matrix π^g across all annotators. We define the confusion matrices set as $\Pi = \{\pi^{1:R}, \pi^g\}$. We associate a Bernoulli random variable $s_i^r \sim B(\omega_i^r)$ with each annotation y_i^r to differentiate the source of noise on it: $s_i^r=1$ if the confusion is caused by the common noise, where ω_i^r is the probability of the global confusion matrix being chosen by annotator r on instance i (see Figure 2b). Denote the set of parameters governing the generation of s_i^r across all annotations as Ω .

Suggested by the successful practice in modeling crowd-sourced data, we also impose the following two commonly made assumptions: 1) each annotator provides their annotations independently (Dawid and Skene 1979); and 2) each annotation is independent from the instance’s features given the ground-truth labels (Yan et al. 2014; Rodrigues and Pereira 2018). We should note the first assumption is *not* contradicting to our common confusion modeling: as the annotators can independently choose the shared common noise model to generate their annotations, the resulting observed annotations are no longer independent across annotators. As a result, the complete data likelihood of observed annotations under our model can be defined as,

$$p(Y, Z|X, \Pi, \Omega) = \prod_{i=1}^N \prod_{r=1}^R \sum_{z=1}^C p(y_i^r|z_i; \Pi, \omega_i^r) p(z_i|\mathbf{x}_i), \quad (1)$$

$$p(y_i^r|z_i; \Pi, \omega_i^r) = \omega_i^r p(y_i^r|z_i, \pi^g) + (1 - \omega_i^r) p(y_i^r|z_i, \pi^r).$$

Based on the above imposed problem structure, we derive an information-theoretical lower bound about the resulting noise modeling quality. Let \hat{Z} be the estimated true labels of all instances. Noise modeling quality is measured by the error rate given by $\mathcal{L}(\hat{Z}, Z) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(\hat{z}_i \neq z_i)$, where $\mathbb{I}(\cdot)$ is an indicator function. Given the ground-truth instance-specific class distribution $\rho_i = \{\rho_{ic}\}_{c=1}^C$ and confusion matrices Π , we have the following theorem about the lower bound of minimax error rate of our model.

Theorem 1. *The minimax error rate of our model is lower bounded by*

$$\begin{aligned} & \inf_{\hat{Z}} \sup_{Z \in [C]^N} \mathbb{E} \left[\mathcal{L}(\hat{Z}, Z) \right] \\ & \geq \frac{1}{N^2 \log C} \sum_{i=1}^N F(\rho_i, \Pi, \Omega) - \frac{\log 2}{N^2 \log C}, \\ & F(\rho_i, \Pi, \Omega) = H(\rho_i) - \sum_{r=1}^R \sum_{c=1}^C \sum_{c'=1}^C \rho_{ic} \rho_{ic'} \left(\omega_i^r \text{KL}(\pi_{c*}^g \parallel \pi_{c'*}^g) \right. \\ & \quad \left. + (1 - \omega_i^r) \text{KL}(\pi_{c*}^r \parallel \pi_{c'*}^r) \right). \end{aligned} \quad (2)$$

where $H(\rho_i) = -\sum_{c=1}^C \rho_{ic} \log \rho_{ic}$ is the entropy of ground-truth class distribution and π_{c*} is the c -th row in confusion matrix π . The proof and further discussion of Theorem 1 is provided in Appendix A.

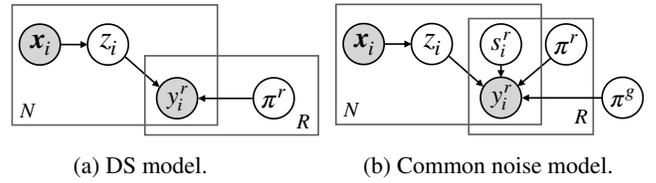


Figure 2: Graphical model presentations of DS model and our common noise model.

Remarks. This result extends the known lower bound result of DS models (Imamura, Sato, and Sugiyama 2018). Lower bound on the error rate measures the difficulty of a crowdsourcing problem. Theorem 1 suggests the proposed decomposition has the potential to further reduce the lower bound, i.e., to obtain better inferred true labels. To understand this result, we should first note that the lower bound mainly depends on the KL distance between the class distributions conditioned on different ground-truth classes, as defined in $F(\rho_i, \Pi, \Omega)$, i.e., how two different classes will be confused with other classes. The more different they are (i.e., a larger KL distance), the easier one can differentiate the two from the observed noisy labels. For example, consider a crowdsourced dataset where an annotator labels a set of instances as *airplane*; but among them, 50% cases should be *bird*, and the other 50% should be *spacecraft*. Intuitively, without any additional knowledge, it is hard to determine the true label when he/she labels an instance as *airplane*. And this is asserted by Theorem 1: If we only used a single confusion matrix for this annotator, the conditional class distributions for *bird* and *spacecraft* will be pushed closer, because their entries on *airplane* are close. This causes a smaller KL term in $F(\rho_i, \Pi, \Omega)$ between *bird* and *spacecraft* (e.g., setting $\omega_i^r=0$ for all instances in annotator r). But if we knew that the confusion between *bird* and *airplane* is caused by common noise, and the confusion between *spacecraft* and *airplane* is caused by individual noise, these mistakes could be attributed to two confusion matrices separately, which eliminates the misleading similarity between the conditional probabilities for *bird* and *spacecraft* caused by *airplane*.

End-to-end Learning Framework

To apply our noise modeling in crowdsourced data, we need to estimate the confusion matrices Π together with the classifier. Instead of building a vanilla tabular model for them, we realize them using neural models, to take advantage of the power of representation learning. In particular, we map the output of the classifier to noisy annotations by two types of confusion layers, which we refer to as noise adaptation layers (Goldberger and Ben-Reuven 2016). We also introduce an auxiliary network that takes both annotator and instance as input to predict the choice of these two noise adaptation layers. Since we treat the ground-truth label of an instance as a latent variable, the Expectation Maximization (EM) algorithm becomes a natural choice for model learning, as typically done in literature (Albarqouni et al. 2016; Rodrigues and Pereira 2018; Bertsekas 2014). For the integrity of work, we provide the derived EM algorithm in Appendix B for interested readers. However, the EM-based

algorithm has several clear drawbacks in our solution: 1) In crowdsourced data, because the annotators typically only label a small proportion of instances, EM-based algorithm becomes very sensitive to the initialization of model parameters. It can easily cause instability issues in training a neural network model. 2) In every EM iteration, we need to retrain the neural network, which causes a huge overhead when handling large networks. Instead, we take an end-to-end approach to jointly perform latent variable inference and model parameter estimation. We define cross-entropy loss on the observed annotations and use error back-propagation to update the classifier’s output and the network parameters simultaneously.

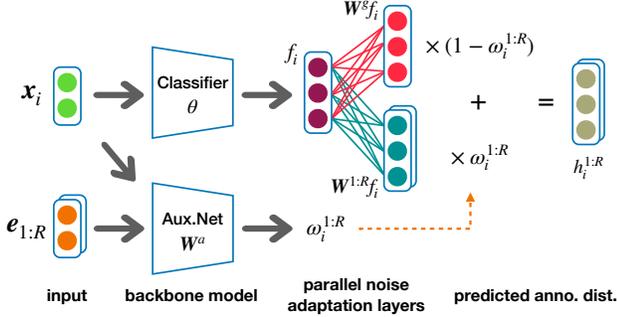


Figure 3: Overview of our framework for classification with 3 classes and R annotators.

We construct a neural network classifier with non-linear intermediate layers and a softmax output layer. The probability distribution of the predicted true label z_i given the instance feature vector \mathbf{x}_i is thus specified as $p_\theta(z_i|\mathbf{x}_i)$, where θ is the network parameter set including the softmax layer. We denote the immediate output of the classifier as $f_i = f(\mathbf{x}_i) \in \mathbb{R}^C$. We then use noise adaptation layers to map the classifier’s output into noisy annotations, which are implemented by introducing additional softmax output layers on top of the output layer of the classifier (see overview in Figure 3). The weight matrices of the noise adaptation layers resemble confusion matrices Π in a probabilistic sense. The output of the noise adaptation layer is thus the probability distribution of predicted annotation $p_{\mathbf{W}}(\hat{y}_i^r|f(\mathbf{x}_i))$, where \mathbf{W} is the parameter set of the noise adaptation layer.

We consider two types of noise adaptation layers: one individual noise adaptation layer for every annotator parameterized by \mathbf{W}^r , and a common noise adaptation layer shared across all annotators parameterized by \mathbf{W}^g . The final probability distribution of annotations is obtained as,

$$p(\hat{y}_i^r|\mathbf{x}_i) = \omega_i^r p_{\mathbf{W}^g}(\hat{y}_i^r|f(\mathbf{x}_i)) + (1 - \omega_i^r) p_{\mathbf{W}^r}(\hat{y}_i^r|f(\mathbf{x}_i)).$$

where ω_i^r governs the distribution that the mistake of annotator r on instance i is caused by common confusion π^g , denoted by the noise source indicator s_i^r .

As s_i^r is unobservable, we introduce an auxiliary network to model $s_i^r \sim B(\omega_i^r)$ by parameterizing it over annotator expertise and instance difficulty, both of which are modeled via learnt representations by the auxiliary network. Specifically, as in our problem setup, every instance is associated

with raw features, the auxiliary network takes instance feature \mathbf{x}_i as input for learning instance i ’s embedding \mathbf{v}_i . The same can be applied to annotator r , if any raw feature \mathbf{e}^r is available about the annotator, otherwise we use its one-hot encoding as input for learning annotator embedding \mathbf{u}_r . Then ω_i^r can be obtained as follows,

$$\begin{aligned} \mathbf{v}_i &= \mathbf{W}_v \mathbf{x}_i + b_v, \mathbf{u}_r = \mathbf{W}_u \mathbf{e}^r + b_u, \\ \omega_i^r &= \sigma(\mathbf{u}_r^\top \mathbf{v}_i). \end{aligned} \quad (3)$$

where (\mathbf{W}_v, b_v) and (\mathbf{W}_u, b_u) are weight matrices and bias terms for annotator and instance embeddings, and σ is a sigmoid function. To simplify our notations, we collectively refer the parameters in this auxiliary network as \mathbf{W}^a . To avoid the magnitude of learnt \mathbf{u} and \mathbf{v} becoming extremely large or small, which causes numerical issues in estimating ω_i^r , we normalize the learnt annotator and instance embeddings before computing their inner product.

Based on the above full specifications of our probabilistic modeling using neural networks, we are ready to estimate the network parameters. We can easily verify that, maximizing the likelihood of observed annotations given the input feature vectors as defined in Eq (1) is equivalent to minimizing the cross-entropy loss between the observed annotations and predicted annotation distributions,

$$\mathcal{L}(\theta, \mathbf{W}^g, \mathbf{W}^{1:R}, \mathbf{W}^a) = -\frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \sum_{j=1}^C y_{ij}^r \log p_j(\hat{y}_i^r|\mathbf{x}_i).$$

where $y_{ij}^r = 1$ if $y_i^r = j$; otherwise $y_{ij}^r = 0$; and $p_j(\hat{y}_i^r|\mathbf{x}_i)$ refers to the j -th entry of the predicted annotation distribution. All parameters can be trained by back-propagation using gradient descent techniques, such as Adam (Kingma and Ba 2014) and SGD (Goodfellow, Bengio, and Courville 2016). Once trained, in the testing phase, we can directly use the classifier to make predictions on new instances.

The gradient flow in back-propagation reveals how our common confusion modeling handles crowdsourced data. In the context of classification, we can simply view the introduced noise adaptation layer as performing a projection of gradients; and with a slight abuse of notations, we denote the output of our noise adaptation layers as $h_i^r = \omega_i^r \mathbf{W}^g f_i + (1 - \omega_i^r) \mathbf{W}^r f_i$. Under the chain rule, the gradients are naturally decoupled with respect to different sources of noise,

$$\frac{\partial \mathcal{L}}{\partial f_i} = \sum_{r=1}^R \frac{\partial \mathcal{L}}{\partial h_i^r} \frac{\partial h_i^r}{\partial f_i} = \sum_{r=1}^R \omega_i^r \frac{\partial \mathcal{L}}{\partial h_i^r} \mathbf{W}^g + (1 - \omega_i^r) \frac{\partial \mathcal{L}}{\partial h_i^r} \mathbf{W}^r. \quad (4)$$

It clearly shows confusion matrices reshape the gradients, which informs the classifier layer what the true label should be on an instance given its noisy annotations. The importance of each confusion matrix in shaping the classifier is determined by ω_i^r , which infers the source of noise based on annotator expertise and instance difficulty.

The gradients in Eq (4) also suggest a potential bottleneck of our proposed solution: if the common and individual noise adaptation layers are unidentifiable, we cannot correctly attribute the noise, which is the key for our solution to

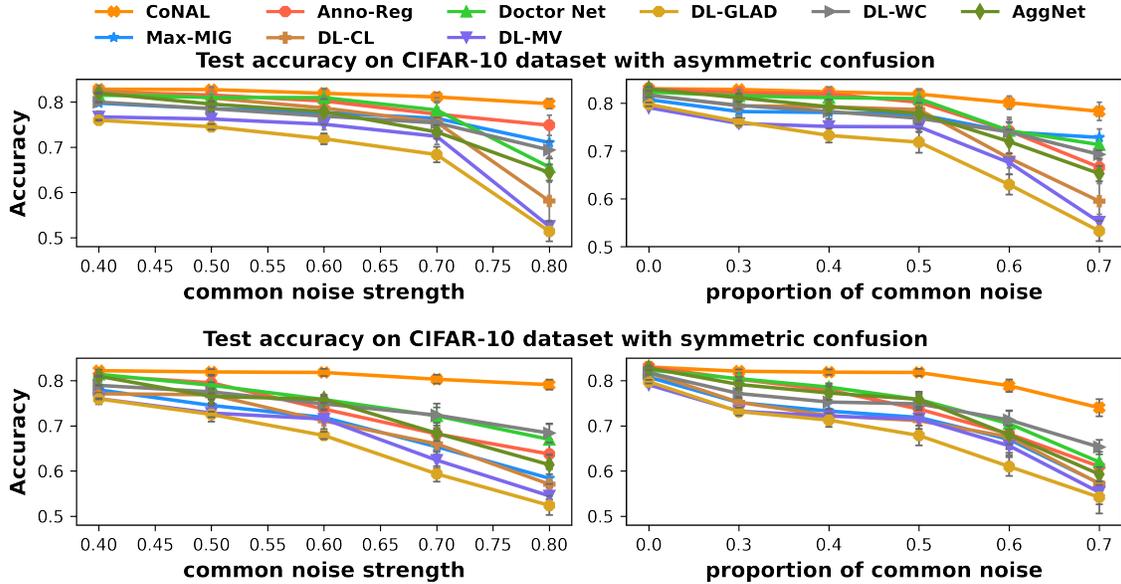


Figure 4: Results on CIFAR-10 dataset.

perform according to Theorem 1. To avoid this, we add ℓ_2 -norm on the difference between the common and individual noise adaptation layers as a regularization term, to enforce them to be different. This presents our final loss function,

$$\mathcal{L}(\theta, \mathbf{W}^g, \mathbf{W}^{1:R}, \mathbf{W}^a) = -\frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \sum_{j=1}^C y_{ij}^r \log p_j(\hat{y}_i^r | \mathbf{x}_i) - \lambda \sum_{r=1}^R \|\mathbf{W}^g - \mathbf{W}^r\|_2$$

where λ is a hyper-parameter to control regularization.

Experiments

We evaluate our method on both synthesized and real-world datasets. We consider a rich set of related solutions as our baselines, which can be divided into two categories:

1) Methods with simple noise models. **DL-MV**: it learns a neural network classifier with labels aggregated by majority voting. **DL-CL** (Rodrigues and Pereira 2018): it learns a neural classifier with designated layers to fit individual annotator confusions (so-called crowd layer). **Anno-Reg** (Tanno et al. 2019): it improves DL-CL by imposing additional trace regularization on individual confusion matrices. **Doctor Net** (Guan et al. 2018): it learns a neural network for every annotator’s annotations and aggregates the networks’ output by weighted majority voting. **Max-MIG** (Cao et al. 2019): it jointly estimates a neural classifier and a label aggregation network using an information-theoretical loss function.

2) Methods with complex noise models. **DL-GLAD**: it learns a neural classifier with labels aggregated by GLAD (Whitehill et al. 2009), where annotator ability and instance difficulty are modeled. **DL-WC**: it learns a neural classifier with labels aggregated by WC (Imamura, Sato, and

Sugiyama 2018), where similar annotators are clustered to share the same confusion matrix. **AggNet** (Albarqouni et al. 2016): an EM-based deep model considering annotator sensitivity and specificity.

Experiments on Synthesized Datasets

We evaluate the proposed method under various settings of synthesized data. Particularly, we demonstrate the effectiveness of our model with different (1) *common confusion types*; (2) *common noise strength*, which is defined as the sum of off-diagonal entries in the common confusion matrix; and (3) *proportion of common noise*, which reflects the percentage of annotations introduced by common confusion.

Datasets description. We generate synthesized crowd-sourced data on two datasets, where we directly manipulate the number of annotators and annotation generation under a variety of settings. On the **Synthetic** dataset, we completely synthesized everything. We first sample a mean vector for every class and then sample instance features from a multi-variate Gaussian distribution parameterized by this mean vector. In particular, we randomly generate 10,000 instances with 6 classes, which are split into a 8,000-instance training set, a 1,000-instance validation set and a 1,000-instance testing set. The **CIFAR-10** dataset is generated based on the CIFAR-10 image classification dataset (Krizhevsky, Hinton et al. 2009). It consists of 60,000 32×32 color images from 10 classes, which are split into a 40,000-instance training set, a 10,000-instance validation set and a 10,000-instance testing set. Image features are used to train the neural classifier on this dataset. In both datasets, each instance in the training set is labeled by averaging 3 randomly selected annotators out of 30 in total.

Synthesizing annotations. We consider two representative noise patterns in common noise: (1) *Asymmetric confusion*. Every class is mapped to another uniformly chosen class

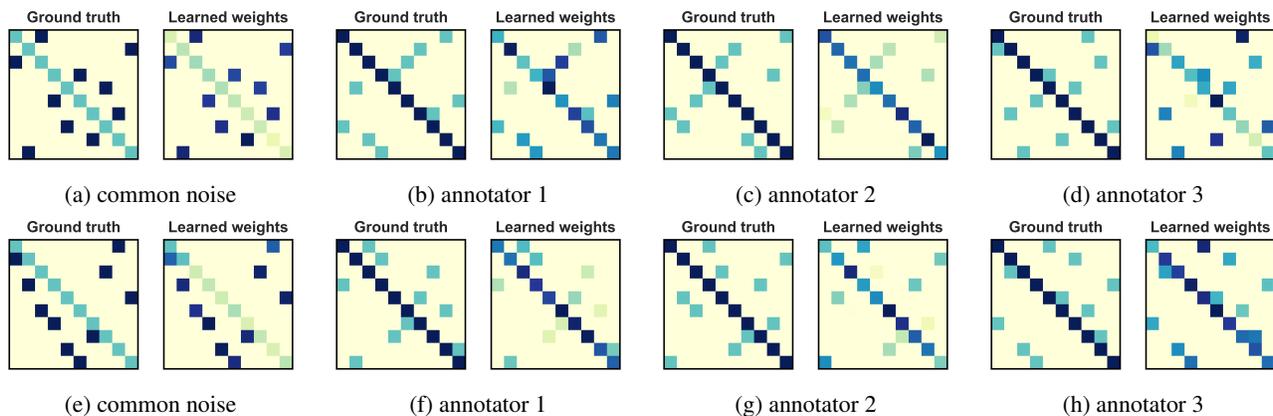


Figure 5: Comparison between ground truth confusion matrices and learned ones on CIFAR-10 dataset. The top row is the result of asymmetric common noise. The bottom row is the result of symmetric common noise.

on both datasets. (2) *Symmetric confusion*. On Synthetic dataset, two random classes are paired and flipped into each other. And on CIFAR-10 dataset, we manually paired similar classes (e.g., *bird* and *airplane*) to be flipped with each other. For individual confusion matrices, we use asymmetric confusion. We generate one global confusion matrix, and one individual confusion matrix for every annotator. In our experiments, the range of common noise strength is set to $[0.4, 0.8]$, while the individual noise strength of annotators is fixed to 0.7. In both noise generation patterns, the noise strength is evenly distributed among the chosen off-diagonal entries.

To control the source of noise in each annotation, i.e., s_i^r , we randomly generate a set of annotator features \mathbf{u} , which are not disclosed to the learners. Given instance feature vector \mathbf{v}_i and annotator feature vector \mathbf{u}_r , we compute ω_i^r by Eq (3) with the ground-truth weight matrices (\mathbf{W}_u, b_u) and (\mathbf{W}_v, b_v) . These weight matrices are not disclosed to the learner. The bias terms are used to control the average proportion of common noise across annotations into a range of $[0.3, 0.7]$. When we generate annotation y_i^r for instance i by annotator r , we first sample $s_i^r \sim B(\omega_i^r)$. If $s_i^r = 1$, the common confusion matrix π^g will be used; otherwise, individual confusion matrix π^r will be used. Then we sample y_i^r from the chosen confusion matrix based on the true label z_i of this instance. We also include a special case that the proportion is 0, where there is no common confusion.

In our experiments, when studying the influence of common noise strength on the learnt classifier, the average proportion of common noise is controlled to be around 0.5. When studying the influence of the proportion of common noise in each annotation, the common and individual noise strength is controlled to 0.4 and 0.7 respectively.

Backbone networks & training details. On the Synthetic dataset, we apply a simple network with only one fully connected (FC) layer (with 128 units and ReLU activations), along with a softmax output layer, using 50% dropout. On the CIFAR-10 dataset, we follow the setting of Cao et al. (2019) to use VGG-16 as the backbone network. We trained the network using the Adam optimizer (Kingma and Ba 2014) with default parameters and learning rate searched

from $\{0.02, 0.01, 0.005\}$. The dimension of annotator and instance embedding is chosen from $\{20, 40, 60, 80\}$. The regularization term λ is searched from $\{10^{-4}, 10^{-5}, 10^{-6}\}$. All experiments are repeated 5 times with different random seeds. Model selection is achieved by choosing the model with the highest accuracy on the validation set. We report mean and standard deviation of test accuracy on the five runs. To make the comparisons fair, all the evaluated methods used the same backbone networks. We implement our framework with PyTorch, and run it on a CentOS system with one NVIDIA 2080Ti GPU with 10 GB memory.

Results. We report the results on the CIFAR-10 dataset in Figure 4, where our solution demonstrated consistent improvement against all baselines across all settings. The observation on the Synthetic dataset is similar, and we present the results in Appendix C due to space limit. All the baselines assumed single source of noise, i.e., annotator-specific noise; as a result, they are heavily influenced when noise become complicated, e.g., a large proportion of mistakes from common confusion and the strength of common noise is strong. Our solution is less sensitive to the environment by decomposing and separately modeling the confusion. When there is no common confusion, the empirical result shows no significant difference between our solution and baselines in this extreme setting, which should also be expected. But we argue that this extreme setting rarely holds in reality, as annotators always share some commonsense about the world.

All models are influenced by symmetric common noise, which directly makes the swapped classes similar. Based on the lower bound provided in Theorem 1, similar conditional class distributions in the confusion matrices will make the problem more difficult, so that the degeneration of all methods are expected under symmetric confusion. In the most extreme case where the proportion of common noise is set to 0.7 and the common noise strength is set to 0.6, nearly 42% annotations are pairwise flipped. However, our method can still outperform baselines with a large margin. Mix-MIG is believed to be robust to correlated mistakes if high-quality annotator exists. However, our experiments show that common confusion poisoned the classifier obtained in Max-MIG even though every annotator is of high quality (individual

	DL-MV	DL-CL	Doctor Net	Anno-Reg	Max-MIG	DL-GLAD	DL-WC	AggNet	CoNAL
LabelMe	79.83±0.34	83.27±0.52	82.12±0.43	82.77±0.48	85.33±0.61	83.12±0.34	82.74±0.33	84.75±0.27	87.12±0.55
Music	72.53±0.41	81.46±0.53	76.58±0.47	79.12±0.36	81.37±0.33	77.82±0.37	75.76±0.24	81.92±0.41	84.06±0.42

Table 1: Test accuracy on two real-world crowdsourcing datasets.

noise strength is set to 0.7). DL-CL and Anno-Reg failed because they could not differentiate the source of noise, such that the gradients from the modeled annotations cannot be properly adjusted to update the classifier. Both Doctor Net and DL-MV are based on majority vote, so that they fail when the annotations across annotators are no longer independent, i.e., caused by the common confusion. Compared to methods with complex noise models, DL-GLAD directly models the annotation accuracy, which is not suitable for class-dependent confusion. DL-WC clusters correlated annotators to share confusion matrix, which can reduce the influence of common confusion. But the expertise of each annotator is missing, which leads to its bad performance. AggNet shows the advantage of directly learning from annotations rather than from aggregated labels. But it still assumes the only noise source thus cannot handle common noise well.

To understand how accurate our solution can distinguish common and individual noise, we report the learnt weights of noise adaptation layers against the ground-truth confusion matrices on the CIFAR-10 dataset in Figure 5. In this experiment, we set the common noise strength to 0.7 and the proportion of common noise to 0.5. We can find that in most cases the ground-truth common noise pattern is well recovered, especially under the asymmetric noise pattern.

Experiments on Real-world Datasets

Datasets description. We consider two real-world datasets. **LabelMe** (Rodrigues and Pereira 2018; Russell et al. 2008) is an image classification dataset, consists of 2,688 images from 8 classes, where 1,000 of them are labeled by annotators from Amazon Mechanical Turk (AMT)¹ and the remainings are used for validation and testing. Each image is labeled by an average of 2.5 annotators, with a mean accuracy of 69.2%. Standard data augmentation techniques are used on training data, including horizontal flips, rescaling and shearing, to enrich the training set to 10,000 images. **Music** (Rodrigues, Pereira, and Ribeiro 2014) is a music genre classification dataset, consisting of 1,000 samples of songs with 30 seconds length from 10 music genres, where 700 of them are labeled by AMT annotators and the rest are used for testing. Each sample is labeled by an average of 4.2 annotators, with a mean annotation accuracy of 73.2%.

Backbone networks & training details. For LabelMe dataset, we followed the setting of Rodrigues and Pereira (2018): we apply a pre-trained VGG-16 network followed by a FC layer with 128 units and ReLU activations, and a softmax output layer, using 50% dropout. For Music dataset, we use the same FC layer and softmax layer as LabelMe. Batch normalization (Ioffe and Szegedy 2015) is performed

in each layer. Other hyper-parameters are the same as the synthesized experiments.

Results. As reported in Table 1, CoNAL achieved new state-of-the-art performance on both real-world datasets. In particular, we looked into the accuracy on classes where commonly made mistakes across annotators are observed (see Figure 1). For example, for *open country* on LabelMe, its accuracy in CoNAL is 67.21%, while the best baseline Max-MIG only achieved 54.19%. The good performance aligns with our analysis in Theorem 1, by differentiating common and individual confusions, it is easier to find the true labels. We provide the visualization of the learned confusion matrices and the training and testing accuracy plots on real-world datasets in Appendix C.

Influence of the regularization term λ . We studied the influence of different λ in Table 2. The results show by enforcing the noise adaptation layers to be different, the performance is improved on both datasets. The value of λ also matters, and 10^{-5} achieves best performance empirically.

λ	0	10^{-4}	10^{-5}	10^{-6}
LabelMe	85.68±0.38	86.61±0.41	87.12±0.55	86.26±0.47
Music	82.14±0.31	83.52±0.25	84.06±0.42	82.98±0.37

Table 2: Model performance under different λ .

Conclusion & Future works

In this paper, we study the problem of learning from crowds with noisy annotations. Aside from the widely employed independent noise assumptions across annotators, we decompose annotation noise into common and individual confusions. We used neural networks to realize our probabilistic modeling of crowdsourced data, and estimate each component in our solution in an end-to-end fashion. Extensive empirical evaluations confirm the advantage of our solution in learning from complicated real-world crowdsourced data. Our solution is also flexible: it can be easily applied to any existing neural classifiers by simply connecting with the proposed noise adaptation layers. In our current solution, all annotators share the same global confusion matrix. An interesting extension is to consider group-wise confusion, where we keep a shared confusion matrix for each annotator group, and identify the groups by optimization. It is also worthwhile to extend the solution to a proactive setting, e.g., probe annotators for more annotations so as to improve common confusion modeling.

¹<https://www.mturk.com/>

Acknowledgments

We thank our anonymous reviewers for their helpful comments. This work was supported by NSF 1718216, 1553568, and Department of Energy DE-EE0008227.

Ethics statement

Our study focuses on tackling an urgent problem in this deep learning era: learning from crowds. High-quality labels are needed for real-world deep learning applications; however, they are typically difficult and expensive to collect in practice. Hence, we propose to directly learn from labels given by non-expert annotators, considering both common mistakes and individualized mistakes. On the one hand, industrial applications will benefit from this work since non-expert labels are both cost- and time-effective to enable deployment of deep learning systems. On the other hand, our work also has academic impact. Our method can be applied to new research problems where high-quality labeled data is rare but crowdsourced labels are easy to obtain, such as medical image classification.

The potential issue of common noise modeling is it might open the door for adversarial annotators. When previously modeled independently, they need to provide a large number of annotations to poison a learner. But if an attacker gets access to common noise, he/she only needs to provide a few annotations consistent with the common noise to amplify the influence of common noise. This will also make other ordinary annotators inadvertently contribute to the attack. Another potential issue of learning from crowds is when modeling annotator expertise, we are learning an annotator profile, which has risk in disclosing their privacy, especially in privacy sensitive annotation problems. Data masking or distortion (e.g., differential privacy) is needed to protect annotators' privacy.

References

- Albarqouni, S.; Baur, C.; Achilles, F.; Belagiannis, V.; Demirci, S.; and Navab, N. 2016. Aggnet: deep learning from crowds for mitosis detection in breast cancer histology images. *IEEE transactions on medical imaging* 35(5): 1313–1321.
- Bertsekas, D. P. 2014. *Constrained optimization and Lagrange multiplier methods*. Academic press.
- Buecheler, T.; Sieg, J. H.; Fuchslin, R. M.; and Pfeifer, R. 2010. Crowdsourcing, open innovation and collective intelligence in the scientific method: a research agenda and operational framework. In *The 12th International Conference on the Synthesis and Simulation of Living Systems, Odense, Denmark, 19–23 August 2010*, 679–686. MIT Press.
- Cao, P.; Xu, Y.; Kong, Y.; and Wang, Y. 2019. Max-mig: an information theoretic approach for joint learning from crowds. *arXiv preprint arXiv:1905.13436*.
- Dawid, A. P.; and Skene, A. M. 1979. Maximum likelihood estimation of observer error-rates using the EM algorithm. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 28(1): 20–28.
- Goldberger, J.; and Ben-Reuven, E. 2016. Training deep neural-networks using a noise adaptation layer. In *International Conference on Learning Representations*.
- Goodfellow, I.; Bengio, Y.; and Courville, A. 2016. *Deep learning*. MIT press.
- Guan, M. Y.; Gulshan, V.; Dai, A. M.; and Hinton, G. E. 2018. Who said what: Modeling individual labelers improves classification. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Imamura, H.; Sato, I.; and Sugiyama, M. 2018. Analysis of minimax error rate for crowdsourcing and its application to worker clustering model. *arXiv preprint arXiv:1802.04551*.
- Ioffe, S.; and Szegedy, C. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.
- Kamar, E.; Kapoor, A.; and Horvitz, E. 2015. Identifying and accounting for task-dependent bias in crowdsourcing. In *Third AAAI Conference on Human Computation and Crowdsourcing*. Citeseer.
- Khetan, A.; and Oh, S. 2016. Achieving budget-optimality with adaptive schemes in crowdsourcing. In *Advances in Neural Information Processing Systems*, 4844–4852.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Krizhevsky, A.; Hinton, G.; et al. 2009. *Learning multiple layers of features from tiny images*. Citeseer.
- Li, Y.; Rubinstein, B.; and Cohn, T. 2019. Exploiting worker correlation for label aggregation in crowdsourcing. In *International Conference on Machine Learning*, 3886–3895.
- Raykar, V. C.; Yu, S.; Zhao, L. H.; Valadez, G. H.; Florin, C.; Bogoni, L.; and Moy, L. 2010. Learning from crowds. *Journal of Machine Learning Research* 11(Apr): 1297–1322.
- Rodrigues, F.; Pereira, F.; and Ribeiro, B. 2014. Gaussian process classification and active learning with multiple annotators. In *International conference on machine learning*, 433–441.
- Rodrigues, F.; and Pereira, F. C. 2018. Deep learning from crowds. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Russell, B. C.; Torralba, A.; Murphy, K. P.; and Freeman, W. T. 2008. LabelMe: a database and web-based tool for image annotation. *International journal of computer vision* 77(1-3): 157–173.
- Shah, N. B.; Balakrishnan, S.; and Wainwright, M. J. 2016. A permutation-based model for crowd labeling: Optimal estimation and robustness. *arXiv preprint arXiv:1606.09632*.
- Tanno, R.; Saeedi, A.; Sankaranarayanan, S.; Alexander, D. C.; and Silberman, N. 2019. Learning from noisy labels by regularized estimation of annotator confusion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 11244–11253.

- Venanzi, M.; Guiver, J.; Kazai, G.; Kohli, P.; and Shokouhi, M. 2014. Community-based bayesian aggregation models for crowdsourcing. In *Proceedings of the 23rd international conference on World wide web*, 155–164.
- Welinder, P.; Branson, S.; Perona, P.; and Belongie, S. J. 2010. The multidimensional wisdom of crowds. In *Advances in neural information processing systems*, 2424–2432.
- Whitehill, J.; Wu, T.-f.; Bergsma, J.; Movellan, J. R.; and Ruvolo, P. L. 2009. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. In *Advances in neural information processing systems*, 2035–2043.
- Yan, Y.; Rosales, R.; Fung, G.; Subramanian, R.; and Dy, J. 2014. Learning from multiple annotators with varying expertise. *Machine learning* 95(3): 291–327.
- Yin, L.; Han, J.; Zhang, W.; and Yu, Y. 2017. Aggregating crowd wisdoms with label-aware autoencoders. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, 1325–1331.
- Zhang, Y.; Chen, X.; Zhou, D.; and Jordan, M. I. 2014. Spectral methods meet EM: A provably optimal algorithm for crowdsourcing. In *Advances in neural information processing systems*, 1260–1268.
- Zhou, D.; Basu, S.; Mao, Y.; and Platt, J. C. 2012. Learning from the wisdom of crowds by minimax entropy. In *Advances in neural information processing systems*, 2195–2203.

A. Proof of Theorem 1

Proof. In our setting, the ground-truth class distribution ρ_i depends on the instance features. Then the minimax error rate of the crowdsourcing problem can be lower bounded by the following,

$$\inf_{\hat{Z}} \sup_{Z \in [C]^N} \mathbb{E} \left[\mathcal{L}(\hat{Z}, Z) \right] \geq \frac{1}{N^2 \log C} \sum_{i=1}^N R(\rho_i, \Pi') - \frac{\log 2}{N^2 \log C} \quad (5)$$

where

$$R(\rho_i, \Pi') = H(\rho_i) - \sum_{r=1}^R \sum_{c=1}^C \sum_{c'=1}^C \rho_{ic} \rho_{ic'} \text{KL}(\pi_{c*}^{r'} \parallel \pi_{c'*}^{r'}) \quad (6)$$

and $\Pi' = \{\pi^{r'}\}_{r=1}^R$ denotes the set of annotator-level confusion matrices. We use π' to differentiate with our defined individual confusion matrix in the main paper. The proof of Eq (5) is similar to (Imamura, Sato, and Sugiyama 2018). Based on our new noise generation assumption, the annotation noise can be decomposed by common noise and individual noise. Thus we can further bound the minimax error rate under this noise assumption.

Under our new noise assumption, we can evaluate the confusion matrix on a per-instance-annotator basis. Specifically, in each annotation, the effective confusion matrix is a weighted combination of the global and individual confusion matrices, where the weight is ω_i^r . In a mixture model, the Kullback–Leibler divergence can be decomposed accordingly by,

$$\begin{aligned} \text{KL}(\pi_{c*}^{r'} \parallel \pi_{c'*}^{r'}) &= \text{KL}(\omega_i^r \pi_{c*}^g + (1 - \omega_i^r) \pi_{c*}^r \parallel \\ &\quad \omega_i^r \pi_{c'*}^g + (1 - \omega_i^r) \pi_{c'*}^r) \\ &\leq \text{KL}(\omega_i^r \parallel \omega_i^r) + \omega_i^r \text{KL}(\pi_{c*}^g \parallel \pi_{c'*}^g) \\ &\quad + (1 - \omega_i^r) \text{KL}(\pi_{c*}^r \parallel \pi_{c'*}^r) \end{aligned} \quad (7)$$

$$\begin{aligned} &= \omega_i^r \text{KL}(\pi_{c*}^g \parallel \pi_{c'*}^g) \\ &\quad + (1 - \omega_i^r) \text{KL}(\pi_{c*}^r \parallel \pi_{c'*}^r) \end{aligned} \quad (8)$$

where $\omega_i^r = (\omega_i^r, 1 - \omega_i^r)$. The inequality can be derived by the log-sum inequality. Substitute Eq (8) back to Eq (6), we can get the new term $F(\rho, \Pi, \Omega)$ in Theorem 1. Plug it back into Eq (5), we can get the refined result in our Theorem 1,

$$\inf_{\hat{Z}} \sup_{Z \in [C]^N} \mathbb{E} \left[\mathcal{L}(\hat{Z}, Z) \right] \geq \frac{1}{N^2 \log C} \sum_{i=1}^N F(\rho_i, \Pi, \Omega) - \frac{\log 2}{N^2 \log C} \quad \square$$

Corollary 1.1. When $N \geq \frac{2 \log 2}{\max F(\rho_i, \Pi, \Omega)}$, increasing the number of instances N will decrease the error rate bound.

Proof.

$$\begin{aligned} &\frac{1}{N^2 \log C} \sum_{i=1}^N F(\rho_i, \Pi, \Omega) - \frac{\log 2}{N^2 \log C} \\ &\leq \frac{\max F(\rho_i, \Pi, \Omega)}{N \log C} - \frac{\log 2}{N^2 \log C}, \end{aligned}$$

When the gradient of the upper bound is less than 0, the upper bound will decrease when N is growing. This can be achieved by setting N by the following,

$$\begin{aligned} &-\frac{1}{N^2} \frac{\max F(\rho_i, \Pi, \Omega)}{\log C} + \frac{2 \log 2}{N^3 \log C} \leq 0 \\ &\Rightarrow N \geq \frac{2 \log 2}{\max F(\rho_i, \Pi, \Omega)} \quad \square \end{aligned}$$

Remarks. The corollary shows when the number of instances is growing, the label aggregation quality gets improved. Also, we need to point out the structure of confusion matrices Π is more important than the number of classes C in this lower bound. With a larger KL distance between every pair of rows in Π , we can expect an improved error lower bound.

B. EM algorithm for learning from crowds by modeling common confusions

The EM algorithm is a generic solution for aggregating crowdsourced labels in classic crowdsourcing problems (Dawid and Skene 1979; Imamura, Sato, and Sugiyama 2018), and it can also be used under our common confusion assumption. Though we have pointed out the main drawbacks of EM-based algorithms in our solution framework, we still list the procedures of using EM algorithm in our problem for interested readers. In particular, we demonstrate a two-step solution, where the latent indicator s_i^r is drawn from a Bernoulli distribution directly parameterized by ω_i^r and the ground-truth label z_i is drawn from a multinomial distribution $p_{\theta}(z_i | \mathbf{x}_i)$ parameterized by θ , which is essentially the soft-classifier we are estimating from the crowdsourced data. Once the ground-truth labels $\{z_i\}_{i=1}^N$ on instances are inferred, we estimate the parameters θ in $p_{\theta}(z_i | \mathbf{x}_i)$ by treating the inferred labels as ground-truth. When the instance features are unavailable, we can use another multinomial distribution $p(z_i | \rho)$ to replace $p_{\theta}(z_i | \mathbf{x}_i)$, where $\rho = \{\rho_c\}_{c=1}^C$ is the corresponding Dirichlet prior, to perform answer aggregation by EM as well.

Under our common confusion assumption, the conditional probability $p(y_i^r | z_i)$ can be written as

$$p(y_i^r | z_i; \Pi, \omega_i^r) = \sum_{s_i^r \in \{0,1\}} p(s_i^r | \omega_i^r) p(y_i^r | z_i, s_i^r, \pi^r)$$

Based on this conditional probability, we derive the EM procedure to infer the ground-truth labels as follows. In the E-step, we estimate hidden ground-truth label z_i and latent indicator s_i^r . The posterior $q(z_i)$ and $q(s_i^r)$ are obtained using Bayes' rule,

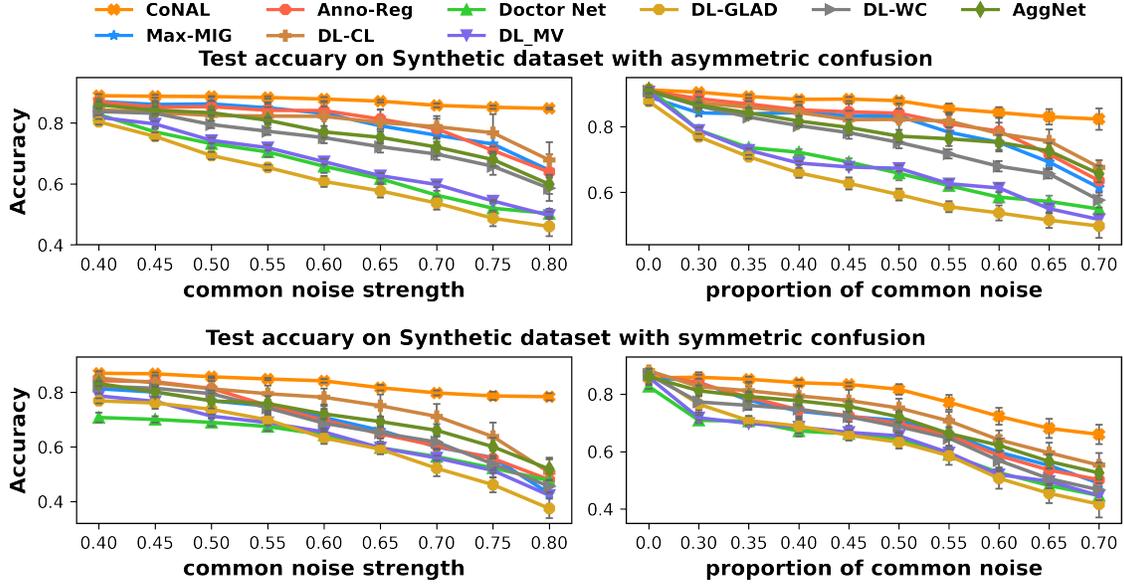


Figure 6: Results on Synthetic dataset.

$$\begin{aligned}
 q(z_i = c) &\propto p_{\theta_0}(z_i = c | \mathbf{x}_i) \prod_{r=1}^R p(y_i^r | z_i = c; \Pi_0, \omega_{i0}^r), \\
 q(s_i^r = 1) &\propto \omega_{i0}^r p(y_i^r | z_i, \pi_0^g), \\
 q(s_i^r = 0) &\propto (1 - \omega_{i0}^r) p(y_i^r | z_i, \pi_0^r).
 \end{aligned}$$

where θ_0, ω_0 and Π_0 are the current estimated parameters. In the M-step, we update the parameters of neural network θ , proportion of common noise ω and confusion matrices Π . The proportion of common noise and confusion matrices have closed-form solutions by using the Lagrange multiplier method (Bertsekas 2014),

$$\begin{aligned}
 \omega_i^r &= q(s_i^r = 1) \\
 \pi_{c,l}^g &= \frac{\sum_{i=1}^N \sum_{r=1}^R q(z_i = c) q(s_i^r = 1) \mathbb{I}(y_i^r = l)}{\sum_{i=1}^N \sum_{r=1}^R q(z_i = c) q(s_i^r = 1)}, \\
 \pi_{c,l}^r &= \frac{\sum_{i=1}^N q(z_i = c) q(s_i^r = 0) \mathbb{I}(y_i^r = l)}{\sum_{i=1}^N q(z_i = c) q(s_i^r = 0)}
 \end{aligned}$$

To update the neural network parameter θ , we follow the approach in (Goldberger and Ben-Reuven 2016; Albarqouni et al. 2016) and use the inferred posterior of ground-truth $q(z_i)$ as the target. Specifically, we compute the cross-entropy loss and backpropagate the error using stochastic gradient optimization techniques such as Adam (Kingma and Ba 2014). For the generic setting where instance features are unavailable, we can update the class distribution ρ using its closed-form solution,

$$\rho_c = \frac{1}{N} \sum_{i=1}^N q(z_i = c)$$

C. Additional experiment results

Results on Synthetic dataset. Figure 6 presents the test accuracy on Synthetic dataset under the same settings as we described in Section 3. We report mean and standard deviation of test accuracy on five runs. The results align with our analysis in Section 3. Under the asymmetric confusion, our proposed approach is robust to the settings of common noise strength and the proportion of common noise. Under the symmetric confusion, all methods' performance is influenced (becomes worse); however, CoNAL still outperforms baselines with a large margin by differentiating the source of noise.

Figure 7 shows the learnt weights of noise adaptation layers against the ground-truth confusion matrices on the Synthetic dataset. We set the common noise strength to 0.7 and the average proportion of common noise around 0.5. We reconstruct the confusion matrix from the learnt weights by normalizing them using softmax on each row. From the results, we can clearly observe most confusion matrices (especially the confusion matrix for common noise) are well recovered under both confusion settings.

Visualization of learnt confusion matrices on real-world datasets. We provide visualization of learnt confusion matrices on both real-world datasets in Figure 8 and 9. We can clearly observe these two types of learnt confusion matrices, i.e., for common confusion and individual confusion, capture different mistake patterns across annotators. For example, on LabelMe dataset, the commonly made mistake from *inside city* to *street* was covered by the learnt common confusion. The same observation is also obtained on the Music dataset, such as the common mistake from *jazz* to *blues* is reflected in our learnt common confusion matrix. On the other hand, the individual noise on annotators captures their own specific mistakes. For example, on LabelMe dataset, both annotator 1 and 2 confused about *tall building* and *inside city*; but this mistake does not appear in annotator 3, 4 and

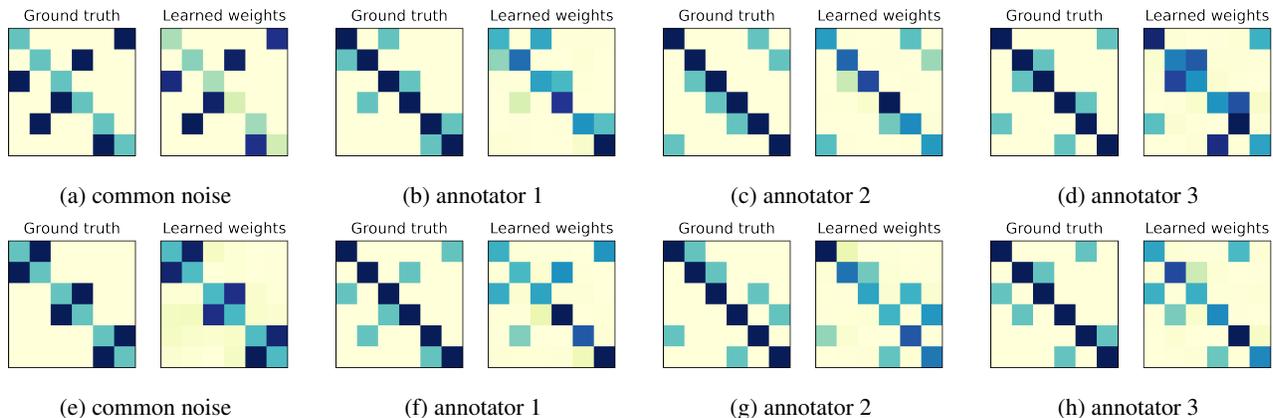


Figure 7: Comparison between ground truth confusion matrices and learned ones on Synthetic dataset. The top row is the result of asymmetric common noise. The bottom row is the result of symmetric common noise.

5, nor the global confusion matrix.

We also notice some low-quality annotators in the Music dataset, such as annotator 1 (with ground-truth annotation accuracy of 0.182) and annotator 5 (with ground-truth annotation accuracy of 0.108), whose annotations are almost random. By separately modeling the annotation noise at a per-annotation basis, our solution reduces the influence from such low-quality annotators in learning the common noise model and maintains the quality of inferred true labels overall.

We visualized the distribution of inferred proportion of common noise (i.e., ω_i^r) across annotations to better understand how CoNAL differentiates the source of noise in individual annotations. We rank the instances by their average ω_i^r over all annotators who have labeled this instance in a descending order. Then we count the frequency of ground-truth labels in the top 50% and bottom 50% instances respectively and report the results in Figure 10. The larger the average ω_i^r in an instance is, the more likely the annotators made similar mistakes on it (i.e., the common confusion matrix can better explain the observed annotations on this instance). On LabelMe dataset, we can observe that annotators tend to make similar annotations on *forest* (with ground-truth annotation entropy 0.660) and *mountain* (with ground-truth annotation entropy 0.192), but make their own mistakes on *open country* (with ground-truth annotation entropy 1.287) and *inside city* (with ground-truth annotation entropy 1.116). In other words, the annotations on *forest* and *mountain* are much more consistent than those on *open country* and *inside city*. This observation can also be explained by the learnt confusion matrices. From the learnt global confusion matrix, we can observe confusion patterns in *open country* and *inside city* are quite scattered, and different annotators (e.g., all those five visualized annotators) have distinct confusions. While for *forest* and *mountain*, the global confusion matrix correctly maps them to the correct annotation, and individual annotators might occasionally make their own mistakes, e.g., annotator 3. Similar observations are also obtained on the Music dataset, where annotators tend to make similar mistakes on *hiphop* and *reggae*, and make their own distinct mistakes on *jazz* and *rock*.

Discussion about overparameterized models. To prove that the improved performance of our solution comes from its unique modeling of crowdsourced data other than simply an increased number of parameters to fit, we compare our model with the overparameterized DL-CL (Rodrigues and Pereira 2018), which has a similar structure as ours to capture individual confusions, but without the notion of modeling common confusion. Rodrigues and Pereira (2018) discussed that simply adding more parameters can make the output of the learnt classifier lose its interpretability as a shared ground-truth estimate across annotators, so that they only used one softmax layer for each annotator upon the classifier’s output layer. We add another softmax layer for each annotator, to introduce more parameters but avoid losing the interpretability of the bottleneck layer, we name it as DL-CL_Over.

Model	DL-CL	DL-CL_Over	CoNAL
#Params in NAL	$R \times C^2$	$2 \times R \times C^2$	$(R + 1) \times C^2$
LabelMe	83.27±0.52	82.34±0.34	87.12±0.55
Music	81.46±0.53	80.47±0.27	84.06±0.42

Table 3: Comparison with the overparameterized model.

We present the results on real-world datasets, along with the number of parameters in the noise adaptation layers (NAL). Even though DL-CL_Over has the most number of parameters to fit, its performance did not increase but decreased, which indicates that blindly adding more parameters will not help model crowdsourced data. Our model adds a global noise adaptation layer, which has fewer parameters than DL-CL_Over. The results prove the advantage of our model comes from its unique design to annotation confusions, but not simply more parameters to fit.

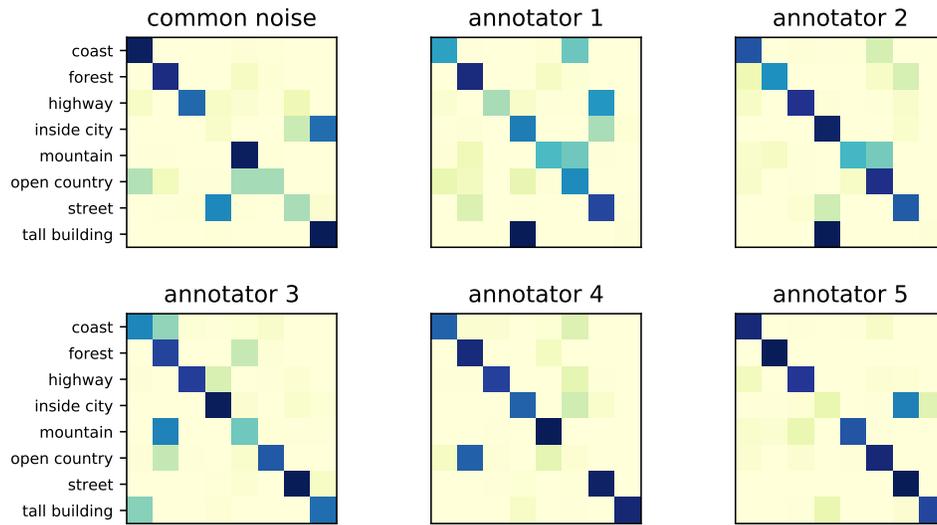


Figure 8: Learnt global confusion matrix and individual confusion matrices of 5 annotators on LabelMe dataset.

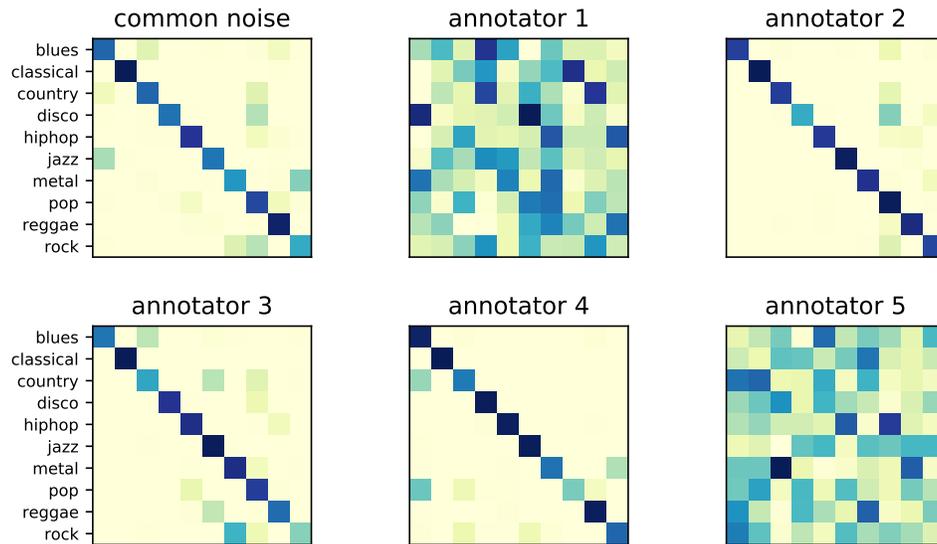


Figure 9: Learnt global confusion matrix and individual confusion matrices of 5 annotators on Music dataset.

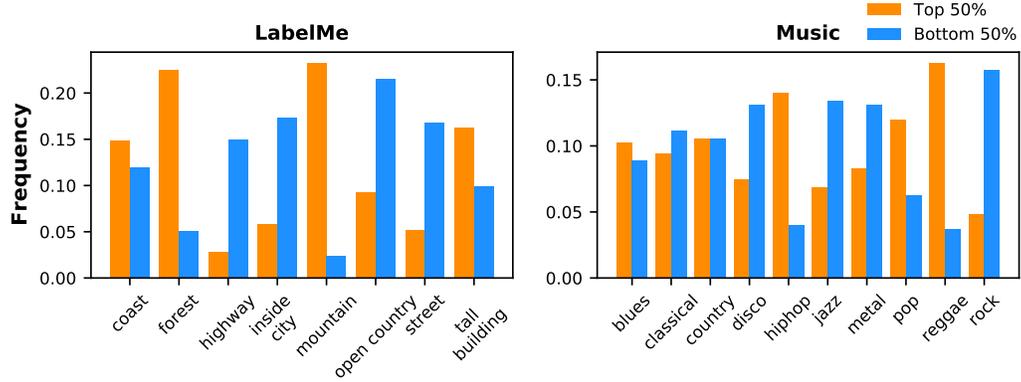


Figure 10: ω -label distribution on real-world datasets. We rank the instances by average ω over all annotators and visualize the ground-truth label distribution of top 50% and bottom 50% instances.