
A Realistic Simulation Framework for Learning with Label Noise

Keren Gu*

Xander Masotto*

Vandana Bachani*

Balaji Lakshminarayanan[†]

Jack Nikodem*

Dong Yin*

* DeepMind, [†] Google Research, Brain Team

{kerengu, xanderm, vbachani, balajiln, nikodem, dongyin}@google.com

Abstract

1 We propose a simulation framework for generating realistic instance-dependent
2 noisy labels via a pseudo-labeling paradigm. We show that this framework gener-
3 ates synthetic noisy labels that exhibit important characteristics of the label noise
4 in practical settings via comparison with the CIFAR10-H dataset. Equipped with
5 controllable label noise, we study the negative impact of noisy labels across a few
6 realistic settings to understand when label noise is more problematic. Additionally,
7 with the availability of annotator information from our simulation framework, we
8 propose a new technique, Label Quality Model (LQM), that leverages annotator
9 features to predict and correct against noisy labels. We show that by adding LQM
10 as a label correction step before applying existing noisy label techniques, we can
11 further improve the models' performance.

12 1 Introduction

13 In many applications, training machine learning models requires labeled data. In practice, the training
14 data labeled by human raters are often noisy, leading to inferior model performance. The study of
15 learning in the presence of label noise dates back to the eighties [Angluin and Laird, 1988], and still
16 receives significant attention in recent years [Natarajan et al., 2013, Reed et al., 2014, Malach and
17 Shalev-Shwartz, 2017, Han et al., 2018, Li et al., 2020a].

18 In the research community, some datasets with real noisy human ratings are available, such as
19 Clothing 1M [Xiao et al., 2015], Food 101-N [Lee et al., 2018] (only a small subset has clean
20 labels), WebVision [Li et al., 2017a], and CivilComments [Borkan et al., 2019], which allow testing
21 approaches that address label noise. However, since the level and type of label noise in these datasets
22 cannot be controlled, it becomes hard to conduct ablation study to understand the impact of noisy
23 labels. As a result, the majority of research work in this area uses benchmark datasets generated by
24 simulations. For example, many prior works simulate noisy labels by flipping the labels according
25 to certain transition matrix [Natarajan et al., 2013, Khetan et al., 2017, Patrini et al., 2017, Han
26 et al., 2018, Hendrycks et al., 2018], independently from the model inputs, e.g., the raw images.
27 However, this type of random label noise may not be an ideal way to simulate realistic noisy labels,
28 since the errors in human ratings are often instance-dependent, i.e., harder examples are easier to
29 get wrong labels, whereas the noisy labels generated by random flipping do not have this type of
30 dependency, even if the transition matrix is asymmetric, i.e., class-conditional. In addition, in many
31 applications, we often have additional features of the raters, such as tenure, historical biases, and
32 expertise level [Cabitza et al., 2020]. Leveraging these features properly can potentially lead to better

33 model performance. However, neither the commonly used public datasets with human ratings nor the
34 synthetic datasets created by random label noise have such rater features available.

35 In this work, we focus on creating realistic benchmarks for the research on label noise. We propose
36 a simulation method that is easy to implement, more realistic than random label flipping, and can
37 convert any commonly used public dataset with clean labels into a noisy label dataset with additional
38 rater features. More specifically, we propose a simulation method based on a pseudo-labeling
39 paradigm: given a dataset with clean labels, we use a subset of it to train a set of models (rater
40 models), and use them to label the rest of the data. In this way, we obtain a dataset whose size is
41 smaller than the original one with clean labels, but with multiple instance-dependent noisy labels.
42 Moreover, some characteristics of the rater models, such as the number of training epochs, the number
43 of samples used, the validation performance metrics, and the number of parameters in the model can
44 be used as a proxy for the rater features.

45 We note that this simulation approach is very similar to self-training in semi-supervised learning
46 [Chapelle et al., 2006]. In the research on label noise, methods inspired by semi-supervised learning
47 have been adopted in several prior works for training robust models [Han et al., 2018; Li et al., 2020a]
48 or generate synthetic noisy label dataset [Lee et al., 2019; Robinson et al., 2020]. We intend to
49 exploit this approach for both providing a comprehensive study of how practical label noise affects
50 the performance of machine learning models, and the research of better training algorithm in the
51 presence of label noise. Our main contributions are summarized as follows:

- 52 • We propose a pseudo-labeling simulation framework for learning with realistic label noise. We
53 provide detailed description, including the generation of rater features (Section 2.2).
- 54 • We propose a systematic protocol for evaluating the synthetic dataset generated by our framework.
55 The evaluation protocol focuses on testing whether the synthetic datasets exhibit some important
56 characteristics of realistic label noise (Section 2.3).
- 57 • We study the negative impact of label noise on deep learning models using our synthetic datasets.
58 We find that noisy labels are more detrimental under class imbalanced settings, when pretraining
59 is not used, and on tasks that are easier to learn with clean labels (Section 3).
- 60 • We propose a label correction approach, named Label Quality Model (LQM), that leverages rater
61 features to significantly improve model performance. We also show that LQM can be combined
62 with other existing noisy label techniques to further improve the performance (Section 4).

63 Moreover, we examine the performance of several existing noisy label algorithm on our synthetic
64 datasets without LQM label correction. We find that the behavior of these techniques on our synthetic
65 datasets is different from the datasets generated by independent random label flipping. We present
66 these results in Appendix A.

67 2 Generating synthetic datasets with realistic label noise

68 In this section, we discuss the formulation of generating synthetic noisy labels, and provide details for
69 the dataset generation procedure and the methods we use to evaluate whether the synthetic datasets
70 share certain characteristics of real human labeled data.

71 2.1 Formulation

72 We consider a K -class classification problem with input space \mathcal{X} and label space $\mathcal{Y} = \{1, \dots, K\}$.
73 In addition, we assume that there is a rater space \mathcal{R} , with each element being the feature of a rater who
74 can label any element in \mathcal{X} . Suppose that there is an unknown distribution over $\mathcal{X} \times \mathcal{Y} \times \mathcal{R} \times \mathcal{Y}$,
75 and each tuple (x, y^*, r, y) in this space corresponds to the input feature of an example x , clean label
76 of the example y^* , a rater r , and the label y provided by the rater.

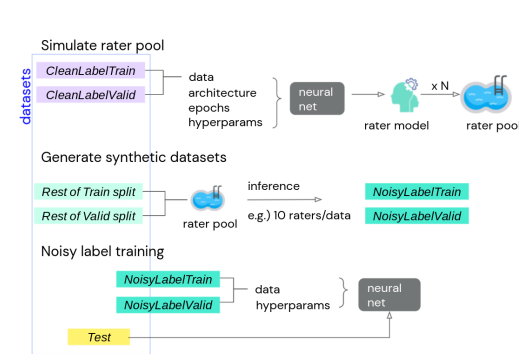
77 The problem of generating synthetic noisy labels can be modeled as generating a noisy label y given
78 a pair of input feature and clean label y^* . Ideally, the probability distribution of the noisy label y
79 should depend on all of x , y^* , and r , i.e., we should generate y according to $p(y | x, y^*, r)$. This
80 means that the label noise should depend on the input—harder and more nuanced examples such
81 as blurred images are more likely to have incorrect labels, as well as the rater—raters with higher
82 expertise level are less likely to make mistakes.

83 However, many prior studies on generating synthetic noisy labels ignore such dependency on x and r
84 and only generate y according to y^* . Here, we specify three approaches for generating noisy labels.

- 85 • *Independent random flipping*. In this method, with probability δ , the label of each example is
86 flipped to an incorrect one, uniformly chosen from all the other $K - 1$ labels [Zhang et al., 2021,
87 Rolnick et al., 2017, Han et al., 2018]. The method is sometimes called symmetric label noise.
88 More specifically, we have $p(y = k | y^*) = (1 - \delta)\mathbf{1}(k = y^*) + \frac{\delta}{K-1}\mathbf{1}(k \neq y^*)$.
- 89 • *Class-conditional random flipping*. In this method, we assume that there is a stochastic matrix
90 $T \in \mathbb{R}^{K \times K}$. The i -th row of T corresponds to the probability distribution of the noisy label
91 y given that the clean label $y^* = i$, i.e., $p(y = j | y^* = i) = T_{i,j}$. This method is sometimes
92 called the asymmetric label noise, and is usually considered more realistic than symmetric noise,
93 since classes that are semantically close are more likely to be confused than classes that have
94 clearer decision boundaries. As mentioned in Section 1, this method has been used in many prior
95 works [Angluin and Laird, 1988, Han et al., 2018, Zhang et al., 2017, Wang et al., 2019, Jiang
96 et al., 2018]; the matrix can be designed with human knowledge or estimated from a small subset
97 of clean data [Patrini et al., 2017, Hendrycks et al., 2018]. Here, we emphasize that the noisy
98 labels in class-conditional label flipping still do not depend on the input feature x and the rater r .
- 99 • *Instance-dependent*, i.e., generating noisy labels according to $p(y | x, y^*, r)$. The method that we
100 propose in this paper satisfies this criterion.

101 2.2 Dataset generation

102 In our framework, we first identify a public dataset that we would like to generate noisy labels for,
103 e.g., CIFAR10 [Krizhevsky and Hinton, 2009] for image classification. We observe that many public
104 datasets already have default training, validation, and test splits. For those without a validation
105 split, we can randomly partition the training data into training and validation splits. We note that in
106 our paper we assume that public datasets have “clean” labels. We acknowledge that many widely
107 used public datasets such as CIFAR10 or ImageNet [Deng et al., 2009] may have mislabeled data
108 points [Northcutt et al., 2021b]; however, the amount of label noise in these public datasets is
109 significantly smaller than what the noisy label research community usually consider [Han et al., 2018,
110 Lee et al., 2019], including our work. Therefore, we believe it is reasonable to consider the labels in
111 public datasets as clean, i.e., less noisy, labels, and we do not expect the label noise in public datasets
112 changes the conclusions in our paper.



125 Figure 1: Pseudo-labeling paradigm for simulating
126 realistic noisy labels.

129 *features*. Then we use all or a subset of models from the rater pool to run inference on the data in
130 the *NoisyLabelTrain* and *NoisyLabelValid* splits. In this way we obtain multiple noisy labels for
131 every data in these two splits, and we replace the clean labels with these noisy labels. We find that in
132 order to control the amount of label noise in these two splits, it is important to train a diverse set of
133 rater models using different combinations of architectures, training steps, learning rate, and batch
134 size. The details for the rater models that we use throughout this paper are provided in Appendix B.
135 To perform label noise research, we can use the *NoisyLabelTrain* split to train models and use the
136 *NoisyLabelValid* split for hyperparameter tuning¹. For the *Test* split, we use the original clean labels.
137 We illustrate our framework in Figure 1.

¹The *NoisyLabelValid* split also contains noisy labels, which may affect the hyperparameters that we select. Understanding the impact of label noise in the validation set is beyond the scope of this paper and will be a future direction.

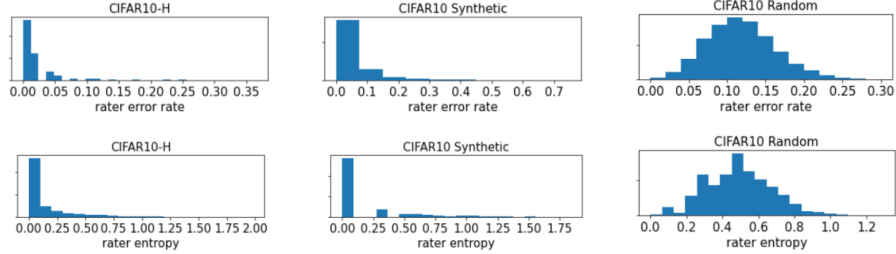


Figure 2: Qualitative comparison among real human labels, our pseudo-labeling framework, and random label flipping. First row: histogram of rater error rate; second row: histogram of rater entropy. Left column: real human labels in CIFAR10-H [Peterson et al., 2019]; middle column: synthetic datasets generated by our pseudo-labeling framework (*NoisyLabelTrain* split); right column: dataset generated by random label flipping. Both CIFAR10-H and our synthetic dataset show right-skewed distribution; whereas the dataset with random label flipping does not.

138 In most of our experiments, the sizes of the *CleanLabelTrain* and *NoisyLabelTrain* splits are around
 139 50% of the original training and validation splits. However, this ratio can be adjusted depending on the
 140 problem of interest. For a synthetic dataset with multiple noisy labels, i.e., the *NoisyLabelTrain* and
 141 *NoisyLabelValid* splits, we use the following two metrics to measure the amount of noise in the
 142 dataset: (1) overall rater error rate, which is defined as the fraction of the incorrect labels among all the
 143 labels given by all the raters, and (2) Krippendorff’s alpha (k-alpha) [Hayes and Krippendorff, 2007],
 144 which measures the agreement between the raters. We note that the computation of Krippendorff’s
 145 alpha does not require the clean labels. Usually, datasets with higher k-alpha are less noisy. All
 146 model training in this and the following sections are performed on TPUs in our internal cluster.

147 2.3 Dataset evaluation

148 Once we have the synthetic datasets, the next step is to evaluate how realistic the synthetic datasets
 149 that we generate are. Here, we qualitatively show that the distribution of the label noise in our dataset
 150 is closer to real human labels, at least when compared to independent label flipping. We evaluate the
 151 following two evaluation criteria for noisy label synthetic datasets:

- 152 (a) The distribution of rater errors should be qualitatively similar to real human labels.
- 153 (b) Label noise should be class-conditional.

154 **Remark 1** For criterion (a), we note that since the label noise is instance-dependent, when we have
 155 multiple labels for every data, the distribution of rater error rate, defined as the ratio of the number of
 156 raters that gave wrong labels for a particular data to the total number of raters that labeled this data,
 157 should have certain dataset-specific characteristics. For example, in a dataset where the majority of
 158 the data are easy to be labeled correctly, while the fraction of the hard examples is relatively small,
 159 we should observe a right-skewed distribution, i.e., the majority of the data have small rater error
 160 rate, and data with large rater error rates account for the tail of the distribution. The similar trend
 161 should also be observed for the entropy (measuring consistency) of the noisy labels. We note that this
 162 right-skewed distribution should not be observed in datasets where labels are flipped with a fixed
 163 probability, *even if the label flipping is class-conditional*. This is because with random flipping, the
 164 rater error rate and rate entropy of the data should concentrate around certain values when we have
 165 multiple labels per data.

166 We emphasize that our evaluation focuses on the overall trend of the distribution of label noise, rather
 167 than each individual instance. This means that we do not expect that our synthetic noisy labels closely
 168 match a certain human labeled dataset for every single data point. Instead, we qualitatively compare
 169 the distribution of the rater error rate (or entropy) and show that the datasets have similar trends.
 170 Making the predictions of the rater models (neural networks) more close to human behavior is an
 171 interesting problem to study, but is beyond the scope of this paper.

172 **Remark 2** For criterion (b), we aim to show that the class confusion matrix has non-uniform off-
 173 diagonal entries. We acknowledge that this can only show that our method is class-conditional rather
 174 than instance-dependent. However, one important difference from prior work is that the confusion
 175 among classes in our framework occurs after the rater models are trained; whereas in prior works it
 176 needs to be designed by human or estimated from a subset of data.

177 **Evaluation results** We now proceed
 178 to present our evaluation results. We
 179 create a synthetic dataset based on the
 180 CIFAR10/100 dataset. For **criterion**
 181 **(a)**, in order to obtain real human labels,
 182 we use the CIFAR10-H dataset,
 183 recently published by [Peterson et al.](#)
 184 [\[2019\]](#). This dataset contains the 10K
 185 data points from the CIFAR10 test
 186 split, and a total of 500K human labels
 187 were crowdsourced, i.e., on average
 188 around 50 labels for each data. In our
 189 proposed simulation framework, we
 190 train 10 rater models using the *Clean-*
 191 *LabelTrain* split and use them to label the *NoisyLabelTrain* split. Details of our synthetic dataset
 192 can be found in Appendix [B](#). For independent random label flipping, we generate 50 labels for each
 193 data, and flipping probability $\delta = 10.8\%$. matching our synthetic dataset. As we can see in Figure [2](#),
 194 the synthetic datasets generated by our pseudo-labeling framework produce similar right-skewed
 195 distribution to the datasets with real human labels in CIFAR10-H, whereas the datasets generated
 196 by random label flipping have different trends. Thus, we conclude that our proposed simulation
 197 framework is more realistic on the first evaluation criterion. As for **criterion (b)**, we compute the
 198 class confusion matrix two synthetic datasets generated by our method for CIFAR10 and CIFAR100,
 199 respectively. As shown in Figure [3](#), the label noise in our dataset is class-conditional. We note that
 200 one can verify that the class confusion matrix of CIFAR10-H is also class-conditional; however,
 201 ensuring that the confusion matrix of our synthetic dataset closely matches certain human labeled
 202 dataset is non-trivial and is beyond the scope of this paper.

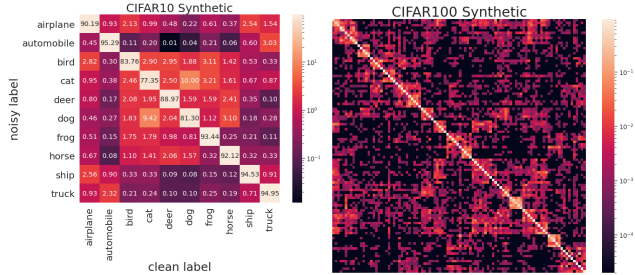


Figure 3: Class confusion matrix for our synthetic dataset. Columns and rows in the CIFAR100 figure (right) are grouped by the 20 coarse labels.

191 *LabelTrain* split and use them to label the *NoisyLabelTrain* split. Details of our synthetic dataset
 192 can be found in Appendix [B](#). For independent random label flipping, we generate 50 labels for each
 193 data, and flipping probability $\delta = 10.8\%$. matching our synthetic dataset. As we can see in Figure [2](#),
 194 the synthetic datasets generated by our pseudo-labeling framework produce similar right-skewed
 195 distribution to the datasets with real human labels in CIFAR10-H, whereas the datasets generated
 196 by random label flipping have different trends. Thus, we conclude that our proposed simulation
 197 framework is more realistic on the first evaluation criterion. As for **criterion (b)**, we compute the
 198 class confusion matrix two synthetic datasets generated by our method for CIFAR10 and CIFAR100,
 199 respectively. As shown in Figure [3](#), the label noise in our dataset is class-conditional. We note that
 200 one can verify that the class confusion matrix of CIFAR10-H is also class-conditional; however,
 201 ensuring that the confusion matrix of our synthetic dataset closely matches certain human labeled
 202 dataset is non-trivial and is beyond the scope of this paper.

203 3 Impact of realistic label noise on deep learning models

204 With the realistic synthetic datasets with noisy labels, our next step is to study the impact of noisy
 205 labels on deep learning models. Interestingly, there exist different views for the impact of noisy
 206 labels to deep neural networks. While most of the recent research works on noisy labels try to design
 207 algorithms that can tackle the negative impact of label noise, some other works claim that deep
 208 learning models are robust to independent random label noise [\[Rolnick et al., 2017, Li et al., 2020b\]](#)
 209 without using sophisticated algorithms. A prominent example is the weak supervision paradigm
 210 [\[Ratner et al., 2016, 2017\]](#), where massive training datasets are generated by weak raters and labeling
 211 functions. Other lines research indicate that large neural network can easily fit all the noisy labels in
 212 the training data [\[Zhang et al., 2021\]](#), while smaller models may be more robust against label noise
 213 due to the regularization effect [\[Advani et al., 2020, Belkin et al., 2019, Northcutt et al., 2021b\]](#).

214 We hypothesize that the negative impact of noisy labels is problem-dependent. While in most cases
 215 the incorrect labels can impair models’ performance, the impact may depend on factors related to the
 216 data distribution and the model. In this section, we choose the following factors to measure the impact
 217 of label noise: the class imbalance, the inductive bias of the model (in particular, pretraining vs
 218 random initialization), and the difficulty of the task (test accuracy that models can achieve when clean
 219 labels are accessible). Note that for better understanding, we decouple these factors with algorithm
 220 design: In this section, we choose simple SGD-style training algorithms with cross-entropy loss and
 221 focus on analyzing the impact of label noise; the discussion on more sophisticated algorithms to
 222 tackle label noise is presented in Sections [A](#) and [4](#). We do not aim to study label aggregation methods
 223 either. Instead, in this and the following sections, given a synthetic dataset with multiple noisy labels,
 224 we generate a dataset with a *single noisy label* by independently and uniformly selecting a random
 225 noisy label for every data point. This is a simulation of the realistic setting where we have a pool of
 226 raters and for each data, we choose a random rater from the pool and request a label.

227 3.1 Label noise has higher impact on more imbalanced datasets

228 One of the important characteristics of many real world datasets is that the classes are usually
 229 imbalanced. When the classes are more imbalanced, the impact of noisy labels may become more
 230 pronounced since the number of data in the minority classes is already small, and noisy labels can
 231 further corrupt these data, making the learning procedure more difficult. We validate this hypothesis
 232 in this section. We use two binary classification tasks, PatchCamelyon (PCam) [\[Veeling et al., 2018\]](#),

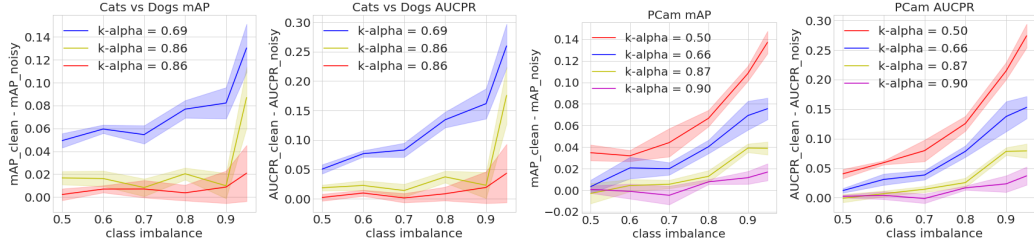


Figure 4: Impact of label noise for tasks with different class imbalance. The x-axis represents class imbalance, measured by the fraction of the majority class. For PCam, we use MobileNet-v1 [Howard et al., 2017] model, and for Cats vs Dogs, we use ResNet50 [He et al., 2016]. The k-alfas correspond to the synthetic datasets before subsampling.

233 [Bejnordi et al., 2017] and Cats vs Dogs (CvD) [Elson et al., 2007]. We generate synthetic noisy
 234 label datasets with different k-alfas, and for each of these datasets, we subsample the two classes
 235 to create several smaller datasets with different class imbalance but the same total number of data.
 236 We note that here we control the class imbalance to be the same for all of the *NoisyLabelTrain*,
 237 *NoisyLabelValid*, and *Test* splits. We train models with clean and noisy labels and use the difference
 238 in mean average precision (mAP) [Zhang and Zhang, 2009] and area under the precision-recall curve
 239 (AUCPR) [Raghavan et al., 1989] as the indicators for the impact of label noise. The results are
 240 shown in Figure 4. As we can see, the impact of label noise becomes more significant as the classes
 241 become more imbalanced.

242 3.2 Pretraining improves robustness to label noise

243 One model training technique that is often used in practice, especially for computer vision and natural
 244 language tasks, is to pretrain the models on some large benchmark datasets and then fine-tune them
 245 using the data for specific tasks. It has been observed that model pretraining can improve robustness
 246 to independent random label noise [Hendrycks et al., 2019] and the web label noise considered by
 247 [Jiang et al., 2020]. Here we show that this can still be observed in our synthetic framework. A simple
 248 explanation is that model pretraining adds strong inductive bias to the models and thus they are less
 249 sensitive to a fraction of noisy labels during fine-tuning.

250 We validate this hypothesis using two datasets, Cats vs Dogs (CvD) and CIFAR10. For both datasets,
 251 we generate three synthetic noisy label datasets using our framework with different rater error rates.
 252 We compare the test accuracy on the *Test* split (with clean labels) between the models that are trained
 253 from random initialization and those that are fine-tuned from models pretrained on ImageNet [Deng
 254 et al., 2009]. We experiment with three different architectures, including Inception-v4 [Szegedy
 255 et al., 2017], ResNet152, and ResNet50 [He et al., 2016]. As we can see in Figure 5, models that
 256 are pretrained on ImageNet achieve better test accuracy. In addition, for pretrained models, the test
 257 accuracy tends to drop more slowly compared to models that are trained from random initialization
 258 as we increase the amount of noise (rater error rate).

259 Meanwhile, we also observe that ImageNet pretraining does not improve the test accuracy under
 260 noisy labels for the PatchCamelyon dataset. This can be explained by the fact that the PatchCamelyon
 261 dataset consists of histopathologic scans of lymph node sections, and these medical images have very
 262 different distribution from the data in ImageNet. Therefore, the inductive bias that the model learned
 263 from ImageNet pretraining may not be helpful on PCam.

264 3.3 Easier tasks are more sensitive to label noise

265 We also study the impact of label noise on tasks with different difficulty levels (the test accuracy
 266 models can achieve when clean labels are accessible). Our hypothesis here is that when a task is
 267 already hard to learn even given clean labels, then the impact of label noise is smaller. The reason
 268 can be that when a classification task is hard, the data distributions of different classes are relatively
 269 close such that even if some data are mislabeled, the final performance may not be heavily impacted.
 270 On the contrary, label noise may be more detrimental to easier tasks as the data distribution can
 271 significantly change when well-separated data points get mislabeled. We validate this hypothesis
 272 with two experiments.

273 **Setup.** Our first experiment involves two binary classification tasks, i.e., PatchCamelyon (PCam)
 274 with MobileNet-v1 [Howard et al., 2017] and Cats vs Dogs with ResNet50 [He et al., 2016]. We

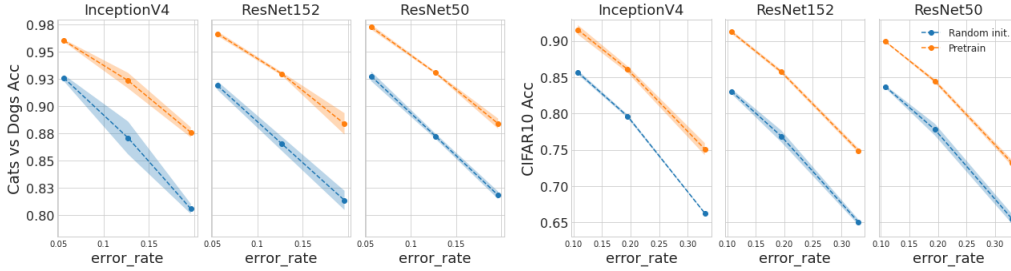


Figure 5: Pretrained models achieve better test accuracy on CvD and CIFAR10. Moreover, as the amount of label noise increases, the amount of test accuracy drop is smaller for pretraining (e.g. the slope is smaller) than training from random initialization.

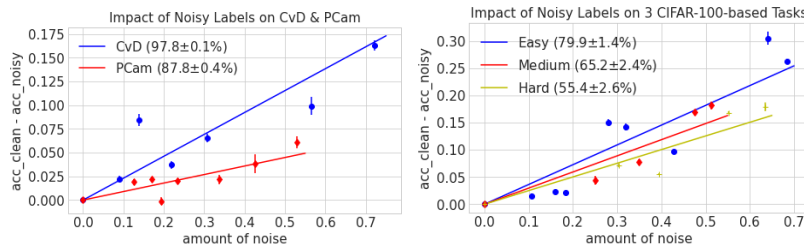


Figure 6: Impact of label noise on tasks with different difficulty levels. The numbers in the legend correspond to test accuracies when training with clean labels for every task. The x-axis represents the amount of noise, measured by $1.0-k$ -alpha. The y-axis represents the negative impact of label noise, measured by the difference in test accuracy when training with clean and noisy labels. Scattered points represent the pairs of noise level and the accuracy drop, and solid lines show the linear fit of the scattered points in the same color.

275 generate synthetic noisy label datasets with different k -alphas using our framework, and compare
 276 the accuracies when the models are trained with clean and noisy labels. We observe that CvD is
 277 easier than PCam (clean label accuracy $97.8 \pm 0.1\%$ vs $87.7 \pm 0.4\%$). In our second experiment,
 278 we design three 20-way classification tasks with the same number of data but different difficulty
 279 levels by subsampling different classes from the CIFAR100 dataset. We call the three tasks the *easy*,
 280 *medium*, and *hard* tasks. Details for the design of the three tasks are provided in Appendix C, and we
 281 observe that with clean labels, we can obtain test accuracies of $79.9 \pm 1.4\%$, $65.2 \pm 2.4\%$, and 55.4
 282 $\pm 2.6\%$ for the three tasks, respectively. We generate synthetic datasets with different amounts of
 283 noise, measured by k -alpha and use the MobileNet-v2 model [Sandler et al., 2018].

284 **Results.** We study the impact of label by measuring the absolute difference in test accuracy when
 285 training with clean and noisy labels. The results are shown in Figure 6. As we can see, the impact of
 286 noisy labels is higher on the easy task: On CvD, the drop in test accuracy grows faster as we increase
 287 the amount of label noise (indicated by $1.0-k$ -alpha) compared to PCam, and similar phenomenon
 288 can be observed on the three CIFAR100-based tasks.

289 4 Leveraging rater features: Label Quality Model

290 Existing work in the noisy label literature commonly assumes that training labels are the only output
 291 of the data curation process. In practice however, the data curation process often produces a myriad
 292 of additional features that can be leveraged in downstream training, e.g., which rater is responsible for
 293 a given label, as well as that rater’s tenure, historical errors, and time spent on a given task. With our
 294 proposed method of simulating instance-dependent noisy labels via rater models, we can additionally
 295 simulate these rater features by extracting metadata from the rater models, e.g., the number of epochs
 296 used to train the rater models is a proxy for rater tenure. Another common practice in label curation
 297 is assigning multiple raters for a single example. This is commonly used to reduce the label noise
 298 via aggregation, or to evaluate the performance of individual raters against the pool. This practice
 299 assumes that agreement among multiple raters are more accurate than individual responses.

300 With understanding of practical data collection setup, we introduce a technique for training with
 301 noisy labels, which we coin *Label Quality Model* (LQM). LQM is an intermediate supervised task
 302 aimed at predicting the clean labels from noisy labels by leveraging rater features and a paired subset
 303 for supervision. The LQM technique assumes the existence of rater features and a subset of training
 304 data with both noisy and clean labels, which we call *paired-subset*. We expect that in real world
 305 scenarios some level of label noise may be unavoidable. LQM approach still works as long as the
 306 clean(er) label is less noisy than a label from a rater that is randomly selected from the pool, e.g.,
 307 clean labels can be from either expert raters or aggregation of multiple raters. LQM is trained on the
 308 paired-subset using rater features and noisy label as input, and inferred on the entire training corpus.
 309 The output of LQM is used during model training as a more accurate alternative to the noisy labels.

310 The intuition for LQM is to correct the labels in a rater-dependent manner. This means that by
 311 learning the patterns from the paired-subset, we can conduct rater-dependent label correction. For
 312 example, LQM can potentially learn that raters with a certain feature often mislabel two breeds of
 313 dogs, then it can possibly correct these two labels from similar raters for the rest of the data. Below
 314 we formally present the details of LQM.

315 **Algorithm design.** Formally, let $D := \{(x_i, y_i, r_i)\}_{i=1}^N$ be a noisy label dataset, e.g., the *Noisy-*
 316 *LabelTrain* split² where x_i is the input, y_i is the one-hot encoded noisy label, and r_i is the rater
 317 feature corresponding to y_i . Let $D_{ps} = \{(x_j, y_j, r_j, y_j^*)\}_{j=1}^M$ be the paired-subset, and y_j^* be a more
 318 accurate label than y_j . We usually have $M \ll N$. We propose to optimize a parameterized model
 319 $LQM(\theta; x, r, y)$ to approximate the conditional probability $P(y^*|x, r, y)$ using D_{ps} . We note that
 320 LQM leverages all the information from the input x , rater features r , and the noisy label y .

321 Once we have the LQM, we proceed to tackle the main task using the noisy label dataset D . Instead of
 322 trying to predict $P(y_i|x_i)$, we replace the noisy labels y_i with the outputs of LQM and train a model
 323 to predict $P(LQM(\theta; x_i, r_i, y_i)|x_i)$. From experimentation, we find that by interpolating between
 324 noisy label y_i and the output of LQM produces even stronger results. Therefore, we recommend
 325 training with target $\tilde{y}_i = \gamma LQM(\theta; x_i, r_i, y_i) + (1 - \gamma)y_i$, where γ is a hyperparameter between 0
 326 and 1 and can be selected using the validation set. This is particularly helpful for datasets with a large
 327 number of classes such as CIFAR100, since it prevents the training target from getting too far from
 328 the original labels y_i . Moreover, since \tilde{y}_i specifies a distribution over the labels, we can also sample a
 329 single one-hot label according to the distribution \tilde{y}_i as the target.

330 We use a small set of rater features in the simulated framework, such as the accuracy of the rater
 331 model on *CleanLabelValid*, the number of epochs trained, and the type of architecture. In addition,
 332 we also use the paired-subset to empirically calculate the confusion matrix for each rater and use it as
 333 a feature for the rater. Instead of training LQM with raw input x , we first train an auxiliary image
 334 classifier $f(x)$ and train LQM using the output logits of $f(x)$. The auxiliary classifier can be trained
 335 over either the full noisy dataset D or the paired-subset D_{ps} . We find that the better option depends
 336 on the task and the amount of noise present. In our experimentation, we train $f(x)$ on both dataset
 337 options and select the better one. Given that LQM has fewer training examples, using an auxiliary
 338 image classifier significantly simplifies training.

339 **Experiment setup and results.** For uniformity, we assume $M = 0.1N$, i.e. 10% of training data
 340 has access to a clean label in all of our following experiments. For the main prediction model, i.e.,
 341 $P(LQM(\theta; x_i, r_i, y_i)|x_i)$, we use the ResNet50 architecture. For the auxiliary model $f(x)$, we use
 342 MobileNet-v2. The LQM itself is trained using a one-hidden-layer MLP architecture with cross-
 343 entropy loss. The number of hidden units in the MLP is chosen in $\{8, 16, 32\}$ as a hyperparameter.
 344 We conduct the following two experiments (see the exact numbers for the results in Appendix D).

- 345 • **LQM vs baseline.** First, we compare the performance of models trained with LQM and the
 346 baseline models that are trained using vanilla cross-entropy loss without leveraging rater features.
 347 Since LQM has access to clean labels of 10% of the data, for fair comparison, we ensure that
 348 the baseline models also have access to the same number of clean labels. The comparison is
 349 presented in Figure 7. As we can see, with rater features and the label correction step, in many
 350 cases, especially in the medium and high noise settings, LQM outperforms the baseline.
- 351 • **Combining LQM with other techniques.** In the second experiment, we investigate the con-
 352 junction of LQM with other noisy labels techniques. We hypothesize that, depending on the

²As mentioned in Section 2.2, the size of the *NoisyLabelTrain* split is around 50% of the training split of the original dataset.

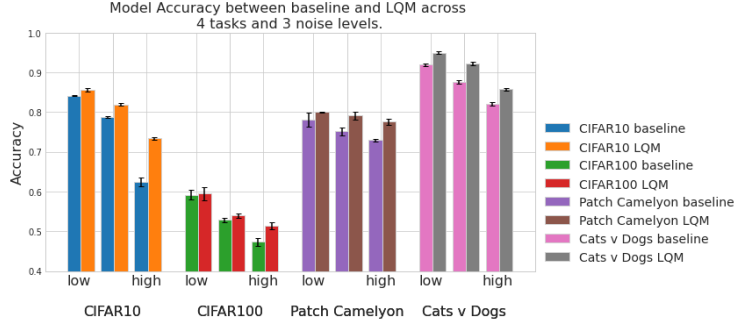


Figure 7: Training with LQM adjusted labels outperforms baselines across all datasets. The improvement from LQM is more prominent in medium and high noise settings. LQM models are trained by randomly sampling a one-hot label from LQM output distribution, and tuning the γ parameter.

353 technique, LQM may be correcting a different kind of noise from existing techniques, and thus
 354 can potentially lead to further performance improvement. To combine LQM with another tech-
 355 nique, we sample a hard label from the soft distribution specified by \tilde{y}_i , and apply other noisy
 356 labels techniques on top of the sampled hard label. We consider the following 4 techniques:
 357 Bootstrap [Reed et al., 2014], Co-Teaching [Han et al., 2018], cross-entropy loss with Monte
 358 Carlo sampling (MCSoftMax) [Collier et al., 2020], and MentorMix [Jiang et al., 2020]. We find
 359 that on CIFAR10 and CIFAR100, combining LQM with other techniques usually lead to further
 360 performance improvement. The improvement can also be observed in the high noise setting for
 361 CvD. The results are illustrated in Figure 8

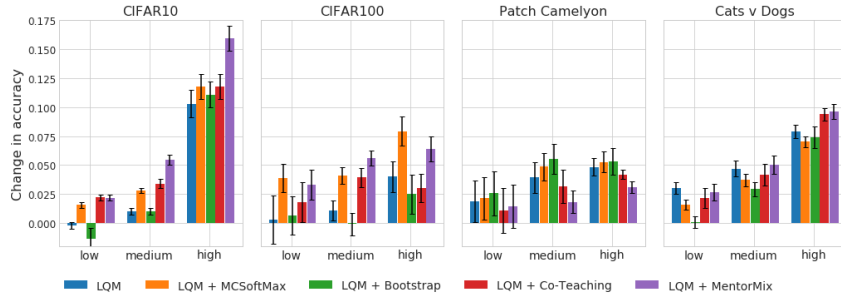


Figure 8: Accuracy improvement of LQM and the combinations of LQM and other techniques compared to the baseline. In many settings (CIFAR10, CIFAR100 and high noise setting in CvD), combining other techniques with LQM further improves the test accuracy. In other cases (PCam and low/medium noise for CvD), the performance gain is less significant.

362 As a final note, since LQM assumes access to a subset of data with clean labels, and also uses an
 363 auxiliary classifier $f(x)$, it has some similarity with semi-supervised learning (SSL). We notice that
 364 several state-of-the-art SSL techniques such as FixMatch [Sohn et al., 2020], UDA [Xie et al., 2019],
 365 self-training with noisy student [Xie et al., 2020] use specifically designed data augmentations that
 366 are only suitable for specific types of data, whereas LQM can be applied to any type of data as
 367 long as we have rater features. We also expect that combining certain SSL techniques (e.g., data
 368 augmentation and consistency training) can further improve the results; however, these extensions are
 369 beyond the scope of this paper.

370 5 Additional related work

371 There is a large body of literature on learning with noisy labels. We mentioned several related work
 372 in the previous sections and it is certainly not an exhaustive list. Since our focus is a more realistic
 373 simulation framework for noisy label research, we first review prior works that try to simulate noisy
 374 labels using methods beyond random label flipping or permutation. As mentioned in Section 1, we
 375 are aware that two prior works by Lee et al. [2019] and Robinson et al. [2020] that also use similar
 376 pseudo-labeling paradigm to generate synthetic datasets with noisy labels. Seo et al. [2019] use a
 377 similar idea of nearest neighbor search in the feature space of a pretrained model with clean labels

378 to generate noisy labels. Compared with these works, our study is much more comprehensive with
379 a diverse set of tasks and model architectures. We conduct a series of analysis on the impact of
380 noisy labels, and propose a method to generate synthetic rater features and use them for improving
381 robustness. These points were not considered in the two prior works. Other approaches to simulating
382 realistic label noise have also been studied in the literature. Jiang et al. [2020] proposes a framework
383 to generate controlled *web label noise*, in which case new images with noisy labels are crawled from
384 the web and then inserted to an existing dataset with clean labels. The framework differs from our
385 approach and the two frameworks should be considered complementary for generating realistic noisy
386 label datasets. In particular, the method by Jiang et al. [2020] is more suitable for web-based data
387 collection, e.g., WebVision [Li et al., 2017a] whereas ours is more suitable for simulating human
388 raters. Moreover, Wang et al. [2018] and Seo et al. [2019] consider open-set noisy labels, where the
389 mislabeled data may not belong to any class of the dataset, similar to Jiang et al. [2020]. Another
390 approach to generating datasets with controllable amount of label noise is to first identify a dataset with
391 noisy labels (potentially some public datasets [Northcutt et al., 2021a,b]) and then use the confident
392 learning (CL) method [Northcutt et al., 2021a] to increase or decrease the amount of label noise
393 proportionally to the distribution of real-world label noise in the dataset. The idea is to model the joint
394 distribution of noisy and true labels and then generate the noisy labels based on the noise-increased
395 or noise-decreased joint distribution of noisy and true labels. This differs from our method since we
396 use rater models, which are trained neural networks to generate noisy labels for each instance.

397 Tackling noisy labels using a small subset of data with clean labels has been considered in a few prior
398 works. Common approaches include pretraining or fine-tuning the network using clean labels [Xiao
399 et al., 2015, Krause et al., 2016], and distillation [Li et al., 2017b]. In loss correction approaches,
400 a subset of clean labels are often used for estimating the confusion matrix [Patrini et al., 2017,
401 Hendrycks et al., 2018]. Veit et al. [2017] propose a method that estimates the residuals between
402 the noisy and clean labels. Ren et al. [2018] use the clean label dataset to learn to reweight the
403 examples. Tsai et al. [2019] combine clean data with self-supervision to learn robust representations.
404 Our approach differs from these prior works since we leverage the additional rater features to learn
405 an auxiliary model that corrects noisy labels in a rater-dependent manner, and can be combined with
406 other techniques to further improve the performance as shown in Section 4.

407 Learning from multiple annotators has been a longstanding research topic. The seminal work
408 by Dawid and Skene [1979] uses the EM algorithm to estimate rater reliability, and much progress
409 has been made since then [Raykar et al., 2010, Zhang et al., 2014, Lakshminarayanan and Teh, 2013].
410 Rater features is commonly available in many human annotation process. In crowdsourcing, several
411 prior work focus on estimating the reliability of raters [Raykar et al., 2010, Tarasov et al., 2014,
412 Moayedikia et al., 2019], and rater aggregation [Vargo et al., 2003, Chen et al., 2013]. Item response
413 theory [Embretson and Reise, 2013] from the psychometrics literature uses a latent-trait model to
414 estimate the proficiency of raters and the difficulty of examples, and has a similar underlying principle
415 to our work.

416 Our method is also broadly related to several other lines of research. Training a pool of rater models
417 is similar to ensemble method [Dietterich, 2000], which is usually used to boost test accuracy [Freund
418 and Schapire, 1997] or improve uncertainty estimation [Lakshminarayanan et al., 2017]. Training new
419 models using noisy labels provided by the rater models is similar to knowledge distillation [Buciluă
420 et al., 2006, Hinton et al., 2015]. Designing instance-dependent noisy label generation framework
421 can be considered as reducing underspecification [D’Amour et al., 2020] in generating label noise.

422 6 Conclusions

423 We propose framework for simulating realistic label noise based on the pseudo-labeling paradigm.
424 We show that the synthetic datasets that we generated are more realistic than independent random
425 label noise. With our synthetic datasets, we evaluate the negative impact of label noise on deep
426 learning models, and demonstrate scenarios where label noise is more detrimental. Using the rater
427 features from our simulation framework, we propose a new technique, Label Quality Model, that
428 leverages annotator features to predict and correct against noisy labels. Our work demonstrates the
429 importance of using realistic datasets for noisy label research, and we hope our framework can serve
430 as an option for the noisy label research community to develop more efficient methods for practical
431 challenges. One limitation of our framework is that, as discussed in Section 5, it focuses more on
432 simulating human errors, whereas there might be other types of label noise in practical settings that
433 need other simulation methods. Another limitation is that LQM requires the paired-subset containing
434 both clean and noisy labels. This requirement may not be satisfied in some applications.

435 **References**

- 436 Madhu S Advani, Andrew M Saxe, and Haim Sompolinsky. High-dimensional dynamics of general-
437 ization error in neural networks. *Neural Networks*, 132:428–446, 2020.
- 438 Dana Angluin and Philip Laird. Learning from noisy examples. *Machine Learning*, 2(4):343–370,
439 1988.
- 440 Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes Van Diest, Bram Van Ginneken, Nico
441 Karssemeijer, Geert Litjens, Jeroen AWM Van Der Laak, Meyke Hermsen, Quirine F Manson,
442 Maschenka Balkenhol, et al. Diagnostic assessment of deep learning algorithms for detection of
443 lymph node metastases in women with breast cancer. *Jama*, 318(22):2199–2210, 2017.
- 444 Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning
445 practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*,
446 116(32):15849–15854, 2019.
- 447 Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced metrics
448 for measuring unintended bias with real data for text classification. In *Companion Proceedings of*
449 *the 2019 World Wide Web Conference*, pages 491–500, 2019.
- 450 Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings*
451 *of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*,
452 pages 535–541, 2006.
- 453 Federico Cabitza, Andrea Campagner, Domenico Albano, Alberto Aliprandi, Alberto Bruno, Vito
454 Chianca, Angelo Corazza, Francesco Di Pietto, Angelo Gambino, Salvatore Gitto, et al. The
455 elephant in the machine: Proposing a new metric of data reliability and its application to a medical
456 case to assess classification reliability. *Applied Sciences*, 10(11):4014, 2020.
- 457 Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. *Semi-Supervised Learning*. The MIT
458 Press, 2006.
- 459 Xi Chen, Paul N Bennett, Kevyn Collins-Thompson, and Eric Horvitz. Pairwise ranking aggregation
460 in a crowdsourced setting. In *Proceedings of the sixth ACM International Conference on Web*
461 *Search and Data Mining*, pages 193–202, 2013.
- 462 Mark Collier, Basil Mustafa, Efi Kokiopoulou, Rodolphe Jenatton, and Jesse Berent. A simple
463 probabilistic method for deep classification under input-dependent label noise. *arXiv preprint*
464 *arXiv:2003.06778*, 2020.
- 465 Alexander D’Amour, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel,
466 Christina Chen, Jonathan Deaton, Jacob Eisenstein, Matthew D Hoffman, et al. Underspecification
467 presents challenges for credibility in modern machine learning. *arXiv preprint arXiv:2011.03395*,
468 2020.
- 469 Alexander Philip Dawid and Allan M Skene. Maximum likelihood estimation of observer error-rates
470 using the EM algorithm. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 28
471 (1):20–28, 1979.
- 472 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-
473 scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern*
474 *Recognition*, pages 248–255. IEEE, 2009.
- 475 Thomas G Dietterich. Ensemble methods in machine learning. In *International Workshop on Multiple*
476 *Classifier Systems*, pages 1–15. Springer, 2000.
- 477 Jeremy Elson, John R Douceur, Jon Howell, and Jared Saul. Asirra: a CAPTCHA that exploits interest-
478 aligned manual image categorization. In *ACM Conference on Computer and Communications*
479 *Security*, volume 7, pages 366–374, 2007.
- 480 Susan E Embretson and Steven P Reise. *Item response theory*. Psychology Press, 2013.
- 481 Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an
482 application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.

- 483 Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi
484 Sugiyama. Co-Teaching: Robust training of deep neural networks with extremely noisy labels.
485 *arXiv preprint arXiv:1804.06872*, 2018.
- 486 Andrew F Hayes and Klaus Krippendorff. Answering the call for a standard reliability measure for
487 coding data. *Communication Methods and Measures*, 1(1):77–89, 2007.
- 488 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
489 recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
490 pages 770–778, 2016.
- 491 Dan Hendrycks, Mantas Mazeika, Duncan Wilson, and Kevin Gimpel. Using trusted data to train
492 deep networks on labels corrupted by severe noise. *arXiv preprint arXiv:1802.05300*, 2018.
- 493 Dan Hendrycks, Kimin Lee, and Mantas Mazeika. Using pre-training can improve model robustness
494 and uncertainty. In *International Conference on Machine Learning*, pages 2712–2721. PMLR,
495 2019.
- 496 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv*
497 *preprint arXiv:1503.02531*, 2015.
- 498 Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand,
499 Marco Andreetto, and Hartwig Adam. MobileNets: Efficient convolutional neural networks for
500 mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- 501 Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. MentorNet: Learning data-
502 driven curriculum for very deep neural networks on corrupted labels. In *International Conference*
503 *on Machine Learning*, pages 2304–2313. PMLR, 2018.
- 504 Lu Jiang, Di Huang, Mason Liu, and Weilong Yang. Beyond synthetic noise: Deep learning on
505 controlled noisy labels. In *International Conference on Machine Learning*, pages 4804–4815.
506 PMLR, 2020.
- 507 Ashish Khetan, Zachary C Lipton, and Anima Anandkumar. Learning from noisy singly-labeled data.
508 *arXiv preprint arXiv:1712.04577*, 2017.
- 509 Jonathan Krause, Benjamin Sapp, Andrew Howard, Howard Zhou, Alexander Toshev, Tom Duerig,
510 James Philbin, and Li Fei-Fei. The unreasonable effectiveness of noisy data for fine-grained
511 recognition. In *European Conference on Computer Vision*, pages 301–320. Springer, 2016.
- 512 Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.
513 *Technical Report*, 2009.
- 514 Balaji Lakshminarayanan and Yee Whye Teh. Inferring ground truth from multi-annotator ordinal
515 data: a probabilistic approach. *arXiv preprint arXiv:1305.0015*, 2013.
- 516 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
517 uncertainty estimation using deep ensembles. *Advances in Neural Information Processing Systems*,
518 30, 2017.
- 519 Kimin Lee, Sukmin Yun, Kibok Lee, Honglak Lee, Bo Li, and Jinwoo Shin. Robust inference via
520 generative classifiers for handling noisy labels. In *International Conference on Machine Learning*,
521 pages 3763–3772. PMLR, 2019.
- 522 Kuang-Huei Lee, Xiaodong He, Lei Zhang, and Linjun Yang. CleanNet: Transfer learning for
523 scalable image classifier training with label noise. In *Proceedings of the IEEE Conference on*
524 *Computer Vision and Pattern Recognition*, pages 5447–5456, 2018.
- 525 Junnan Li, Richard Socher, and Steven CH Hoi. DivideMix: Learning with noisy labels as semi-
526 supervised learning. *arXiv preprint arXiv:2002.07394*, 2020a.
- 527 Mingchen Li, Mahdi Soltanolkotabi, and Samet Oymak. Gradient descent with early stopping is
528 provably robust to label noise for overparameterized neural networks. In *International Conference*
529 *on Artificial Intelligence and Statistics*, pages 4313–4324. PMLR, 2020b.

- 530 Wen Li, Limin Wang, Wei Li, Eirikur Agustsson, and Luc Van Gool. WebVision database: Visual
531 learning and understanding from web data. *arXiv preprint arXiv:1708.02862*, 2017a.
- 532 Yuncheng Li, Jianchao Yang, Yale Song, Liangliang Cao, Jiebo Luo, and Li-Jia Li. Learning from
533 noisy labels with distillation. In *Proceedings of the IEEE International Conference on Computer
534 Vision*, pages 1910–1918, 2017b.
- 535 Eran Malach and Shai Shalev-Shwartz. Decoupling “when to update” from “how to update”. *arXiv
536 preprint arXiv:1706.02613*, 2017.
- 537 Alireza Moayedikia, William Yeoh, Kok-Leong Ong, and Yee Ling Boo. Improving accuracy and
538 lowering cost in crowdsourcing through an unsupervised expertise estimation approach. *Decision
539 Support Systems*, 122:113065, 2019.
- 540 Nagarajan Natarajan, Inderjit S Dhillon, Pradeep Ravikumar, and Ambuj Tewari. Learning with noisy
541 labels. In *Neural Information Processing Systems*, volume 26, pages 1196–1204, 2013.
- 542 Curtis Northcutt, Lu Jiang, and Isaac Chuang. Confident learning: Estimating uncertainty in dataset
543 labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021a.
- 544 Curtis G Northcutt, Anish Athalye, and Jonas Mueller. Pervasive label errors in test sets destabilize
545 machine learning benchmarks. *arXiv preprint arXiv:2103.14749*, 2021b.
- 546 Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making
547 deep neural networks robust to label noise: A loss correction approach. In *Proceedings of the
548 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1944–1952, 2017.
- 549 Joshua C Peterson, Ruairidh M Battleday, Thomas L Griffiths, and Olga Russakovsky. Human
550 uncertainty makes classification more robust. In *Proceedings of the IEEE/CVF International
551 Conference on Computer Vision*, pages 9617–9626, 2019.
- 552 Vijay Raghavan, Peter Bollmann, and Gwang S Jung. A critical investigation of recall and precision
553 as measures of retrieval system performance. *ACM Transactions on Information Systems (TOIS)*, 7
554 (3):205–229, 1989.
- 555 Alexander Ratner, Christopher De Sa, Sen Wu, Daniel Selsam, and Christopher Ré. Data program-
556 ming: Creating large training sets, quickly. *Advances in Neural Information Processing Systems*,
557 29:3567, 2016.
- 558 Alexander Ratner, Stephen H Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré.
559 Snorkel: Rapid training data creation with weak supervision. In *Proceedings of the VLDB
560 Endowment. International Conference on Very Large Data Bases*, volume 11, page 269. NIH
561 Public Access, 2017.
- 562 Vikas C Raykar, Shipeng Yu, Linda H Zhao, Gerardo Hermosillo Valadez, Charles Florin, Luca
563 Bogoni, and Linda Moy. Learning from crowds. *Journal of Machine Learning Research*, 11(4),
564 2010.
- 565 Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew
566 Rabinovich. Training deep neural networks on noisy labels with bootstrapping. *arXiv preprint
567 arXiv:1412.6596*, 2014.
- 568 Mengye Ren, Wenyan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for
569 robust deep learning. In *International Conference on Machine Learning*, pages 4334–4343. PMLR,
570 2018.
- 571 Joshua Robinson, Stefanie Jegelka, and Suvrit Sra. Strength from weakness: Fast learning using
572 weak supervision. In *International Conference on Machine Learning*, pages 8127–8136. PMLR,
573 2020.
- 574 David Rolnick, Andreas Veit, Serge Belongie, and Nir Shavit. Deep learning is robust to massive
575 label noise. *arXiv preprint arXiv:1705.10694*, 2017.

- 576 Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mo-
577 bileNetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on*
578 *Computer Vision and Pattern Recognition*, pages 4510–4520, 2018.
- 579 Paul Hongsuck Seo, Geeho Kim, and Bohyung Han. Combinatorial inference against label noise. In
580 *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 581 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image
582 recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 583 Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D Cubuk, Alex
584 Kurakin, Han Zhang, and Colin Raffel. FixMatch: Simplifying semi-supervised learning with
585 consistency and confidence. *arXiv preprint arXiv:2001.07685*, 2020.
- 586 Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Du-
587 mitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In
588 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9,
589 2015.
- 590 Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the
591 inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer*
592 *Vision and Pattern Recognition*, pages 2818–2826, 2016.
- 593 Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander Alemi. Inception-v4, Inception-
594 Resnet and the impact of residual connections on learning. In *Proceedings of the AAAI Conference*
595 *on Artificial Intelligence*, volume 31, 2017.
- 596 Alexey Tarasov, Sarah Jane Delany, and Brian Mac Namee. Dynamic estimation of worker reliability
597 in crowdsourcing for regression tasks: Making it work. *Expert Systems with Applications*, 41(14):
598 6190–6210, 2014.
- 599 Tsung Wei Tsai, Chongxuan Li, and Jun Zhu. Countering noisy labels by learning from auxiliary
600 clean labels. *arXiv preprint arXiv:1905.13305*, 2019.
- 601 John Vargo, John C Nesbit, Karen Belfer, and Anne Archambault. Learning object evaluation:
602 computer-mediated collaboration and inter-rater reliability. *International Journal of Computers*
603 *and Applications*, 25(3):198–205, 2003.
- 604 Bastiaan S Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation equivari-
605 ant CNNs for digital pathology. In *International Conference on Medical image computing and*
606 *computer-assisted intervention*, pages 210–218. Springer, 2018.
- 607 Andreas Veit, Neil Alldrin, Gal Chechik, Ivan Krasin, Abhinav Gupta, and Serge Belongie. Learning
608 from noisy large-scale datasets with minimal supervision. In *Proceedings of the IEEE Conference*
609 *on Computer Vision and Pattern Recognition*, pages 839–847, 2017.
- 610 Yisen Wang, Weiyang Liu, Xingjun Ma, James Bailey, Hongyuan Zha, Le Song, and Shu-Tao Xia.
611 Iterative learning with open-set noisy labels. In *Proceedings of the IEEE Conference on Computer*
612 *Vision and Pattern Recognition*, pages 8688–8696, 2018.
- 613 Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross
614 entropy for robust learning with noisy labels. In *Proceedings of the IEEE/CVF International*
615 *Conference on Computer Vision*, pages 322–330, 2019.
- 616 Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. Learning from massive noisy
617 labeled data for image classification. In *Proceedings of the IEEE Conference on Computer Vision*
618 *and Pattern Recognition*, pages 2691–2699, 2015.
- 619 Qizhe Xie, Zihang Dai, Eduard Hovy, Minh-Thang Luong, and Quoc V Le. Unsupervised data
620 augmentation for consistency training. *arXiv preprint arXiv:1904.12848*, 2019.
- 621 Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student
622 improves ImageNet classification. In *Proceedings of the IEEE/CVF Conference on Computer*
623 *Vision and Pattern Recognition*, pages 10687–10698, 2020.

- 624 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep
625 learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115,
626 2021.
- 627 Ethan Zhang and Yi Zhang. Average precision. *Encyclopedia of Database Systems*, pages 192–193,
628 2009.
- 629 Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical
630 risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- 631 Yuchen Zhang, Xi Chen, Dengyong Zhou, and Michael I Jordan. Spectral methods meet EM:
632 A provably optimal algorithm for crowdsourcing. *Advances in Neural Information Processing*
633 *Systems*, 27:1260–1268, 2014.
- 634 Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures
635 for scalable image recognition. In *Proceedings of the IEEE Conference on Computer Vision and*
636 *Pattern Recognition*, pages 8697–8710, 2018.

637 Checklist

- 638 1. For all authors...
- 639 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
640 contributions and scope? [Yes]
- 641 (b) Did you describe the limitations of your work? [Yes] See Sections [5](#) and [6](#).
- 642 (c) Did you discuss any potential negative societal impacts of your work? [No] We don’t
643 think our work has negative societal impacts.
- 644 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
645 them? [Yes]
- 646 2. If you are including theoretical results...
- 647 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 648 (b) Did you include complete proofs of all theoretical results? [N/A]
- 649 3. If you ran experiments (e.g. for benchmarks)...
- 650 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
651 mental results (either in the supplemental material or as a URL)? [Yes] We made the
652 data that we generated using our simulation framework accessible for the reviewers
653 through the following link: [link will be shared with reviewers privately](#). This folder
654 also includes instructions for how to load the data. In our paper, we provide sufficient
655 details, such as hyperparameters in Appendix [B](#) for training rater models, and sufficient
656 details in Section [4](#) for LQM. We are working on making the datasets publicly available.
- 657 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
658 were chosen)? [Yes] For LQM, we believe we provide enough experiment details
659 in Section [4](#). For training rater models, we provide significant amount of details in
660 Appendix [B](#).
- 661 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
662 ments multiple times)? [Yes]
- 663 (d) Did you include the total amount of compute and the type of resources used (e.g., type
664 of GPUs, internal cluster, or cloud provider)? [Yes] See Section [2.2](#).
- 665 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 666 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 667 (b) Did you mention the license of the assets? [Yes] See Appendix [E](#).
- 668 (c) Did you include any new assets either in the supplemental material or as a URL?
669 [Yes] The noisy label datasets generated by our framework at [link will be shared with](#)
670 [reviewers privately](#).
- 671 (d) Did you discuss whether and how consent was obtained from people whose data you’re
672 using/curating? [Yes] See Appendix [E](#) for discussion on data license.

- 673 (e) Did you discuss whether the data you are using/curating contains personally identifiable
674 information or offensive content? [No] We don't think our data has such information
675 or content.
- 676 5. If you used crowdsourcing or conducted research with human subjects...
- 677 (a) Did you include the full text of instructions given to participants and screenshots, if
678 applicable? [N/A]
- 679 (b) Did you describe any potential participant risks, with links to Institutional Review
680 Board (IRB) approvals, if applicable? [N/A]
- 681 (c) Did you include the estimated hourly wage paid to participants and the total amount
682 spent on participant compensation? [N/A]