

Estimating Noisy Class Posterior with Part-level Labels for Noisy Label Learning

Supplementary Material

A. Classifier training with PLM

In Eq. (4) of main paper, we discussed the empirical risk for estimating the noisy class posterior and the single-to-multiple transition matrix. In this supplementary material, we will provide a detailed discussion on how to train a consistent classifier using PLM and loss correction techniques [6].

As discussed in the main paper, the classification task aims to learn a classifier $f : \mathcal{X} \rightarrow \mathcal{C}$ that maps each instance \mathbf{x}_i to its corresponding label y_i . Given the network for estimating the clean class posterior as $g : \mathcal{X} \rightarrow \mathcal{R}^c$, the classifier can be represented as follows: $f(\mathbf{x}) = \arg \max_{i \in \mathcal{C}} g_i(\mathbf{x})$. Here, $g_i(\mathbf{x})$ refers to the i -th element of the vector $g(\mathbf{x})$, which represents the estimated probability $\hat{P}(Y = i | X = \mathbf{x})$. Given a noisy dataset $\tilde{\mathcal{D}} = \{(\mathbf{x}_i, \tilde{y}_i)\}_{i=1}^n$, the empirical risk of the classifier is defined as:

$$\tilde{R}(f) = \frac{1}{n} \sum_{i=1}^n \ell_1(f(\mathbf{x}_i), \tilde{y}_i), \quad (1)$$

where ℓ_1 denotes a classification loss. Loss correction methods typically introduce a transition matrix $T(\mathbf{x})$ to establish a connection between the posterior of the noisy and clean classes. This allows training a clean classifier by minimizing the empirical risk with noisy dataset. Based on existing loss correction methods, the noise transition matrix $T(\mathbf{x})$ can be estimated, and we have $P(\tilde{Y} | X = \mathbf{x}) = T(\mathbf{x})^\top P(Y | X = \mathbf{x})$. Let the noisy class posterior estimation network be denoted as $g^e : \mathcal{X} \rightarrow \mathcal{R}^c$ where $g^e_i(\mathbf{x}) = P(\tilde{Y} = i | X = \mathbf{x})$. The noisy class classifier $f^e(\mathbf{x})$ can be represented as:

$$f^e(\mathbf{x}) = \arg \max_{i \in \mathcal{C}} g^e_i(\mathbf{x}) = \arg \max_{i \in \mathcal{C}} (T(\mathbf{x})^\top g)_i(\mathbf{x}). \quad (2)$$

Therefore, the empirical risk in loss correction methods can be expressed as:

$$\tilde{R}(g) = \frac{1}{n} \sum_{i=1}^n \ell_1(f^e(\mathbf{x}_i), \tilde{y}_i). \quad (3)$$

By minimizing this loss it is possible to construct classifier-consistent algorithms.

We denote the part-level labels estimation network as $g^p : \mathcal{X} \rightarrow \mathcal{R}^c$ where $g^p_i(\mathbf{x}) = P(Y'_i = 1 | X = \mathbf{x})$. Given the single-to-multiple transition matrix $U(\mathbf{x})$ where $U_{ij}(\mathbf{x}) = P(Y'_j = 1 | \tilde{Y} = i, X = \mathbf{x})$, the part-level multi-

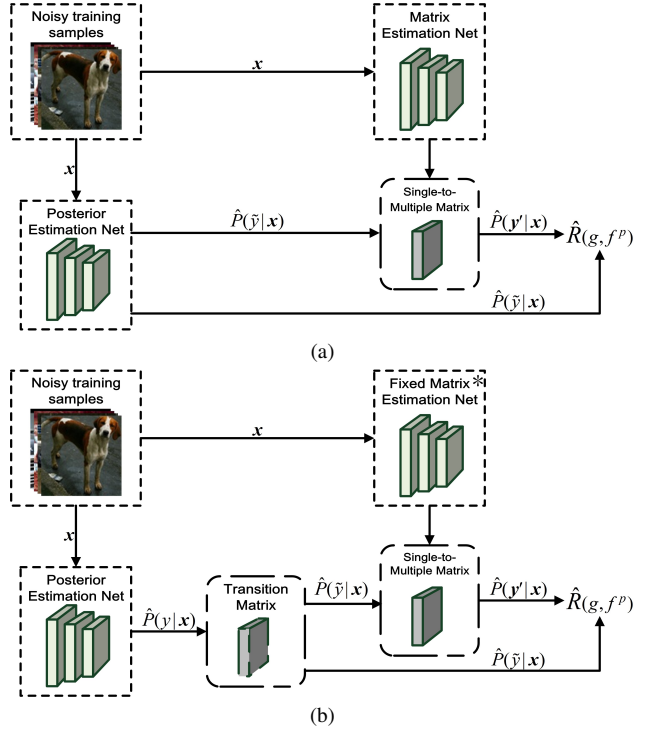


Figure 1. Illustration of neural network training using PLM. (a) Utilizing PLM for estimating the noisy class posterior while simultaneously training the matrix estimation network. (b) Fixing the matrix estimation network (with "**") and integrating loss correction method to facilitate noisy label learning.

label classifier $f^p(\mathbf{x})$ can be represented as: 040

$$\begin{aligned} f^p(\mathbf{x}) &= \{i | g^p_i(\mathbf{x}) > \frac{1}{2}\} = \{i | (U(\mathbf{x})^\top g^e)_i(\mathbf{x}) > \frac{1}{2}\} \\ &= \{i | (U(\mathbf{x})^\top T(\mathbf{x})^\top g)_i(\mathbf{x}) > \frac{1}{2}\}. \end{aligned} \quad (4) \quad 041$$

Similarly, given a dataset $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ with multiple part-level labels, the empirical risk of the training with part-level labels is defined as: 042

$$R'(f^p) = \frac{1}{n} \sum_{i=1}^n \ell_2(f^p(\mathbf{x}_i), \mathbf{y}_i). \quad (5) \quad 046$$

where ℓ_2 denotes a multi-label classification loss. 047

Then, the empirical risk of the joint training framework 048

Algorithm 1 PLM training framework

Input: Noisy training dataset \mathcal{D} , noise transition matrix $T(\mathbf{x})$ derived from existing methods.

Output: Classifier model f .

- 1: Minimize a classification loss to learn a labeling classifier f^l from \mathcal{D} .
- 2: Obtain the set of sub-instances \mathcal{S} by cropping the instances in \mathcal{D} .
- 3: Construct multi-labels by using f^l to label the sub-instances in \mathcal{S} .
- 4: Train a single-to-multiple transition matrix estimation network g^u by minimizing the loss defined in Eq. (4) of the main paper.
- 5: Fix the parameters of g^u , set $g^e(\mathbf{x}) = (T(\mathbf{x})^\top g)(\mathbf{x})$, then minimize the Eq. (4) of the main paper with updated g^u to optimize g .
- 6: Obtain the classifier $f(\mathbf{x}) = \arg \max_{i \in \mathcal{C}} g_i(\mathbf{x})$. Here, $g_i(\mathbf{x})$ represents the i -th element of the network output vector $g(\mathbf{x})$.
- 7: **return** Optimized classifier f .

049 is defined as:

$$\begin{aligned} \hat{R}(g, f^p) &= \frac{1}{2}(\tilde{R}(g) + R(f^p)) \\ &= \frac{1}{2n} \sum_{i=1}^n [\ell_1(f^e(\mathbf{x}_i), \tilde{y}_i) + \ell_2(f^p(\mathbf{x}_i), \mathbf{y}'_i)]. \end{aligned} \quad (6)$$

050
051 We minimize the empirical risk to obtain a robust classi-
052 fier. The training process for the single-to-multiple transi-
053 tion matrix is depicted in Figure 1a, in which we achieve
054 it by minimizing the empirical risk defined in Eq. (4) of
055 the main paper. The training of the classifier is illustrated
056 in Figure 1b, where we keep the trained matrix estimation
057 network fixed and combine it with the noise transition ma-
058 trix obtained through existing loss correction methods. We
059 then optimize the empirical risk discussed before for train-
060 ing. The training procedure is outlined concisely in Algo-
061 rithm 1.

062 In Section 4.3 of the main paper, we have conducted a
063 comparison of performance using various matrix estimation
064 methods. More specifically, for PLM-F, PLM-D, and PLM-
065 V, we adopted the matrix estimation techniques outlined in
066 Forward [6], Dual-T [10], and VolMinNet [5], respectively.
067 Subsequently, we minimize the empirical risk as defined in
068 Eq. (6). For PLM-R, we combined PLM with T-Revision
069 [8] and introduced the slack variable ΔT , then $f^e(\mathbf{x})$ in Eq.
070 (2) can be modified to

$$071 \quad f^{er}(\mathbf{x}) = \arg \max_{i \in \mathcal{C}} ((T(\mathbf{x}) + \Delta T)^\top g)_i(\mathbf{x}). \quad (7)$$

072 Following T-Revision, we also incorporated a importance
073 reweighting strategy. The minimized empirical risk of

PLM-R is defined as follows:

$$\hat{R}(f^{er}, f^p) = \frac{1}{2n} \sum_{i=1}^n [w \ell_1(f^{er}(\mathbf{x}_i), \tilde{y}_i) + \ell_2(f^p(\mathbf{x}_i), \mathbf{y}'_i)], \quad (8)$$

where $w = \frac{g_{y_i}(\mathbf{x}_i)}{((T(\mathbf{x}) + \Delta T)^\top g)_{y_i}(\mathbf{x}_i)}$ denotes the weight.

B. Identifiability of single-to-multiple transition matrix

In the main paper, we introduce a brand-new single-to-multiple transition matrix. In this section, we will discuss the identifiability of this transition matrix. Specifically, regarding $P(\mathbf{Y}'|\mathbf{x}) = U^\top(\mathbf{x})P(\tilde{\mathbf{Y}}|\mathbf{x})$, when the matrix $U(\mathbf{x})$ is unconstrained, the following issue arises: there exists an infinite number of non-singular matrices $Q \in \mathbb{R}^{c \times c}$ such that

$$P(\mathbf{Y}'|\mathbf{x}) = (U^\top(\mathbf{x})Q)(Q^{-1}P(\tilde{\mathbf{Y}}|\mathbf{x})). \quad (9)$$

This situation emerges from the network training process in joint training framework of Section 3.4:

$$g^p(\mathbf{x}) = g^u(\mathbf{x})g^e(\mathbf{x}), \quad (10)$$

where $g^p(\mathbf{x})$, $g^e(\mathbf{x})$ and $g^u(\mathbf{x})$ correspond to the estimates of part-level labels, noisy class posterior, and the matrix respectively. The specific concern appears to center on the scenario where $g^p(\mathbf{x}) = \hat{P}(\mathbf{Y}'|\mathbf{x})$ and $g^e(\mathbf{x}) = Q^{-1}\hat{P}(\tilde{\mathbf{Y}}|\mathbf{x})$, yielding $g^u(\mathbf{x}) = U^\top(\mathbf{x})Q \neq U^\top(\mathbf{x})$.

To address this issue, we employ joint training to simultaneously utilize noisy labels and part-level labels for optimizing both g^e and g^p . More precisely, g^e is directly guided by \tilde{Y} , aligning with a coarse $g^e(\mathbf{x}) = \hat{P}(\tilde{\mathbf{Y}}|\mathbf{x})$, while $g^p(\mathbf{x})$ is supervised by Y' to meet $g^p(\mathbf{x}) = \hat{P}(\mathbf{Y}'|\mathbf{x})$. This dual supervision constrains $g^u(\mathbf{x})$ to comply with $\hat{P}(\mathbf{Y}'|\mathbf{x}) = g^u(\mathbf{x})\hat{P}(\tilde{\mathbf{Y}}|\mathbf{x})$, resulting in $g^u(\mathbf{x}) = \hat{U}(\mathbf{x})^\top$. This means that the potential scenario, where Q leads to $g^e(\mathbf{x}) = Q^{-1}\hat{P}(\tilde{\mathbf{Y}}|\mathbf{x})$ and then $g^u(\mathbf{x}) = \hat{U}(\mathbf{x})^\top Q$, is preemptively negated through supervision from \tilde{Y} . Therefore, during training with the joint framework, the matrix's identifiability is ensured. This approach also echoes the matrix estimation strategy presented in MEIDTM [1].

C. Analysis of time complexity

In the main paper, we introduced additional modules to aid in estimating the noisy class posterior, which to some extent increases the algorithm's time complexity. Therefore, in this section, we will discuss the time complexity and efficiency of the proposed method in comparison to our baseline model (Forward [6]).

Let us assume that the time complexity for training the baseline model for one epoch is denoted as $O(T)$, and the

time complexity for making predictions on the entire training set is $O(P)$. Additionally, the introduction of an extra noise transition matrix layer contributes an additional time complexity of $O(E)$. Considering Forward as the baseline, the time complexity of the proposed method can be expressed as $O(e_1T + cP + e_2(2T + E) + e_3(T + P + 2E))$, where e_1 , e_2 , and e_3 represent the number of epochs for the annotator, transition matrix estimator, and classifier training, respectively, and c indicates the cropping frequency.

For comparison, the time complexity of the Forward method is expressed as $O(e_4T + P + e_5(T + E))$, where e_4 and e_5 represent the number of epochs for the anchor estimation network and classifier training, respectively. For the sake of facilitating comparison, we assume that each section underwent an equal number of training epochs e , i.e., $e = e_1 = e_2 = e_3 = e_4 = e_5$. The complexity of the proposed method is represented as $O(4eT + (c + e)P + (1 + 2e)E)$, whereas the complexity of the Forward method is denoted as $O(2eT + P + eE)$. Then, the transition matrix can be considered as a noise adaptation layer with fixed parameters, which has a total of c^2 parameters where c represents the number of categories. In comparison, deep neural networks have a much larger number of parameters. Taking Resnet-18 as an example, it has a total of 11.7M parameters. Therefore, in this paper and in most cases, we have $T \ll P$ and $T \ll E$. Consequently, we can simplify the two computational complexities to $O(4eT + (c + e)P)$ and $O(2eT + P)$. Since the training process involves backpropagation and gradient computation, it takes more time than the prediction process, leading to $T > P$. Additionally, in this paper, the number of pruning iterations satisfies $e \gg c$. As a result, the computational complexity of the proposed method follows $O(4eT + (c + e)P) < O(6eT)$, and the Forward follows $O(2eT + P) > O(2eT)$. Hence, under the assumption of setting the same number of epochs in each stage, the time overhead of the proposed method should be less than three times that of the Forward. Additionally, for the purpose of evaluating the efficiency of our approach, we conducted a comparison of the code’s runtime based on the CIFAR-10 dataset.

Table 1. The time consumption of PLM and Forward (used as the baseline).

Method	Time Consumption (min)
PLM	35.56
Forward	19.21

Analysis and experiments indicate that our approach significantly enhances the performance of noisy label learning (NLL), with only a linear increase in time consumption.

D. Analysis of cropping strategies

The instance cropping method is related to the multi-labeling of the proposed approach. In the paper, we selected the four corners and the central part of the image data for cropping and determined the cropping size through empirical analysis on the validation set. Table 2 displays the experimental results of different cropping sizes on CIFAR-10 data with sym-50% noise, and we additionally attempted two other cropping strategies. The cropping strategies used in table 2 are as follows: the uniform strategy involves five uniform crops at the four corners and center of the image as used in the paper. The random strategy entails five crops at random positions. The emphasized strategy constructs two sub-instances based on feature emphasis, with one sub-instance masking the top emphasized number of features and the other sub-instance masking the remaining features.

Table 2. The classification accuracy (expressed in percentage) with different cropping sizes and strategies.

Size	Uniform	Random	Emphasized
9	83.58 ± 0.45	83.42 ± 0.86	84.14 ± 0.59
36	83.80 ± 0.31	83.40 ± 0.55	84.24 ± 0.29
81	82.69 ± 2.32	83.81 ± 0.54	84.48 ± 0.47
144	83.17 ± 0.91	83.81 ± 0.39	84.19 ± 0.42
256	83.49 ± 0.95	83.40 ± 0.34	84.28 ± 0.86
361	83.62 ± 0.31	83.65 ± 0.74	84.32 ± 0.46
484	84.99 ± 0.40	84.36 ± 0.32	84.24 ± 0.61
625	85.08 ± 0.16	83.97 ± 0.72	84.26 ± 0.57
784	84.24 ± 0.24	83.91 ± 0.70	84.47 ± 0.33

The emphasized strategy demonstrates superior performance and displays enhanced stability, suggesting the potential for further refinement of cropping strategies in the context of NLL classification, as discussed in Section 5 of the paper. Furthermore, within the established cropping strategy, the method shows robustness to the cropping size.

E. Visualization of focused features

In Figure 2, we employ a visualization approach to provide a visual interpretation of the effectiveness of the PLM method. The STL-10 [2] dataset is used for visualization purposes. In Figure 2b, it is shown that when the labels contain noise, the network emphasizes the background region associated with those labels. Consequently, the model tends to overfit to the noise, hindering the network’s ability to learn features that truly capture the distinctive characteristics of the instances. As a result, the estimation of the posterior for the noisy labels becomes excessively confident. As depicted in Figure 2c, removing the overemphasized features through cropping effectively redirects the model’s attention to other more informative features. By generating



Figure 2. Illustration of class activation maps (CAM) for overemphasized region correction: the highlighted area (with more intense red color) indicates the emphasized area of a model trained from noisy labels. (a) Original images with noisy labels: car, bird, dog, monkey, ship, deer. (b) CAMs for estimating noisy class posterior by the classifier. (c) CAMs when excluding the overemphasized regions after cropping. (d) CAMs for estimating noisy class posterior by the model after PLM training.

196 labels associated with these features and providing addi-
 197 tional supervised information during network training, the
 198 network can focus on more diverse features. As shown in
 199 Figure 2d, compared to Figure 2b, the network pays more
 200 attention to object-relevant features.

201 F. Combination with state-of-the-art methods

202 In this section, we combine and compare our proposed
 203 framework with different state-of-the-art (SOTA) methods
 204 to further validate the effectiveness of our approach. In the
 205 Section 4 of the main paper, we mentioned that we did not
 206 compare PLM with SOTA methods. This is because these
 207 methods incorporate a lot of robust learning strategies to
 208 achieve better empirical performance, while our method fo-
 209 cuses solely on enhancing the noisy class posterior estima-
 210 tion to assist in building a classifier-consistent algorithm.
 211 To explore whether our method can flexibly combine with
 212 these robust learning strategies to achieve improved classifica-
 213 tion performance, we combine and compare our proposed
 214 method with SOTA methods using different strategies, as

215 detailed in Table 3 and 4. Specifically, in Table 3, we com-
 216 pare our method with following two representative SOTA
 217 methods on the CIFAR-10 and CIFAR-100 datasets: (1)
 218 DivideMix [4], a method that combines data augmentation,
 219 label selection, co-training, semi-supervised learning, and
 220 pseudo-labeling strategies; (2) CTRR [12], a method based
 221 on contrastive learning strategies that constructs a regular-
 222 ization function. Besides, in Table 4, we compare PLM with
 223 following two methods designed for instance-dependent
 224 noise on the CIFAR-10 dataset: (1) BLTM [9], a method
 225 using deep neural networks and bayes optimal labels to es-
 226 timate transition matrix; (2) CausalNL [11], a method based
 227 on a structural causal framework. We combine PLM with
 228 these methods to verify its complementary effects. In the
 229 combination with DivideMix (named PLM.DivideMix), we
 230 introduce an additional loss term based on Eq. (4) of the
 231 main paper to the selected labeled samples. In the combina-
 232 tion with CTRR and CausalNL (named PLM.CTRR and
 233 PLM.CausalNL), we modify the cross-entropy loss term of
 234 their loss function to Eq. (4) of the main paper. In the com-

Table 3. The average classification accuracy and standard deviation (expressed in percentage) across five trials under various synthetic noisy label settings. The better classification accuracy is indicated in **bold**.

	CIFAR-10				CIFAR-100			
	Sym-20%	Sym-50%	Pair-20%	Pair-45%	Sym-20%	Sym-50%	Pair-20%	Pair-45%
CTRR	93.02 ± 0.12	84.96 ± 1.12	92.81 ± 0.27	68.54 ± 1.81	71.61 ± 0.64	65.53 ± 0.46	69.94 ± 0.36	46.17 ± 0.70
PLM_CTRR	93.16 ± 0.36	86.59 ± 0.26	92.96 ± 0.36	78.56 ± 1.71	72.06 ± 0.26	66.00 ± 0.56	70.44 ± 0.33	46.84 ± 0.52
DivideMix	95.71 ± 0.47	94.77 ± 0.06	92.65 ± 0.38	68.67 ± 1.95	76.72 ± 0.31	73.12 ± 0.30	76.62 ± 0.25	47.01 ± 1.02
PLM_DivideMix	96.06 ± 0.19	95.12 ± 0.18	96.01 ± 0.07	76.27 ± 1.13	79.96 ± 0.14	74.20 ± 0.40	79.93 ± 0.18	47.19 ± 0.88

Table 4. The average classification accuracy and standard deviation (expressed in percentage) across five trials on the CIFAR-10 dataset with instance-dependent noise settings. The better classification accuracy is indicated in **bold**.

	IDN-20%	IDN-30%	IDN-40%	IDN-50%
BLTM	76.70 ± 0.55	72.12 ± 0.59	65.44 ± 1.01	56.77 ± 0.75
PLM_BLM	89.73 ± 0.22	87.47 ± 0.50	84.40 ± 0.84	76.28 ± 3.80
CausalNL	79.66 ± 0.38	76.58 ± 0.46	72.86 ± 0.43	67.75 ± 1.15
PLM_CausalNL	81.44 ± 0.38	78.86 ± 0.92	75.52 ± 0.38	73.20 ± 1.06

235 combination with BLTM (named PLM.BLM), we use the tran-
 236 sition matrix estimated in BLTM. The remaining settings
 237 are consistent with those in Section 4 of the main paper.

238 The experimental results demonstrate that our approach
 239 can integrate with SOTA methods which rely on com-
 240 plex robust learning strategies, improving their classifica-
 241 tion performance. The improvement occurs even though
 242 these methods do not explicitly require the explicit estima-
 243 tion of noisy class posteriors. This could be because, during
 244 the process of estimating noisy class posteriors using PLM,
 245 the richer supervised information assists the model in learn-
 246 ing more reasonable representations that reflect the instance
 247 characteristic.

248 G. Analysis of transition matrix estimation

249 In this section, we aim to validate the assistance of PLM in
 250 transition matrix estimation error. In the experiments pre-
 251 sented in Table 5, we modified the strategy of noisy class
 252 posterior estimation used in the most basic transition ma-
 253 trix estimation method Forward [6] to PLM’s strategy. The
 254 results indicate that PLM can help in transition matrix esti-
 255 mation by reducing the error in estimating noisy class pos-
 256 terior.

Table 5. The average errors and standard deviation of transition matrix estimation across five trials on the CIFAR-10 dataset. The lower error is indicated in **bold**.

	Sym-20%	Sym-50%	Pair-20%	Pair-45%
Forward	0.35 ± 0.01	0.60 ± 0.12	0.27 ± 0.01	0.74 ± 0.03
PLM_Forward	0.16 ± 0.04	0.32 ± 0.06	0.23 ± 0.01	0.62 ± 0.04

H. Experiments on real-world dataset 257

258 In the main paper, we compared the experimental results
 259 on the real-world dataset Clothing1M. To further illustrate
 260 the performance of PLM on real-world datasets, we com-
 261 pared the results on the Animal-10N [7] dataset in the Table
 262 6. Animal-10N consists of 50,000 noisy samples for train-
 263 ing and 5,000 clean samples for testing. We selected 10%
 264 of the training set as the validation set. We use the SGD
 265 optimizer and cosine learning rate decay strategy to train
 266 the network, with an initial learning rate 10^{-2} , weight
 267 decay of 10^{-2} , and momentum of 0.9. The backbone and
 268 other settings are the same as InstanceGM [3]. The results
 269 further demonstrate that PLM can more effectively handle
 270 real-world noise.

Table 6. Accuracy on the Animal-10N benchmark. The baseline results and experimental settings refer to InstanceGM. The better classification accuracy is indicated in **bold**.

Method	CE	Dropout	SELFIE	PLC	Nested	InstanceGM	PLM
Acc. (%)	79.4	81.3	81.8	83.4	81.3	84.6	85.08

References 271

- 272 [1] De Cheng, Tongliang Liu, Yixiong Ning, Nannan Wang,
 273 Bo Han, Gang Niu, Xinbo Gao, and Masashi Sugiyama.
 Instance-dependent label-noise learning with manifold-
 274 regularized transition matrix estimation. In *Proceedings of
 275 the IEEE/CVF Conference on Computer Vision and Pattern
 276 Recognition*, pages 16630–16639, 2022. 2 277
- 278 [2] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of
 279 single-layer networks in unsupervised feature learning. In
 280 *Proceedings of the fourteenth international conference on
 281 artificial intelligence and statistics*, pages 215–223. JMLR
 282 Workshop and Conference Proceedings, 2011. 3 282
- 283 [3] Arpit Garg, Cuong Nguyen, Rafael Felix, Thanh-Toan Do,
 284 and Gustavo Carneiro. Instance-dependent noisy label
 285 learning via graphical modelling. In *Proceedings of the
 286 IEEE/CVF winter conference on applications of computer
 287 vision*, pages 2288–2298, 2023. 5 287
- 288 [4] Junnan Li, Richard Socher, and Steven CH Hoi. Dividemix:
 289 Learning with noisy labels as semi-supervised learning.
 290 In *International Conference on Learning Representations*,
 291 pages 1–13, 2019. 4 291

- 292 [5] Xuefeng Li, Tongliang Liu, Bo Han, Gang Niu, and Masashi
293 Sugiyama. Provably end-to-end label-noise learning with-
294 out anchor points. In *International Conference on Machine*
295 *Learning*, pages 6403–6413. PMLR, 2021. 2
- 296 [6] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon,
297 Richard Nock, and Lizhen Qu. Making deep neural net-
298 works robust to label noise: A loss correction approach. In
299 *Proceedings of the IEEE conference on computer vision and*
300 *pattern recognition*, pages 1944–1952, 2017. 1, 2, 5
- 301 [7] Hwanjun Song, Minseok Kim, and Jae-Gil Lee. SELFIE:
302 Refurbishing unclean samples for robust deep learning. In
303 *ICML*, 2019. 5
- 304 [8] Xiaobo Xia, Tongliang Liu, Nannan Wang, Bo Han, Chen
305 Gong, Gang Niu, and Masashi Sugiyama. Are anchor points
306 really indispensable in label-noise learning? *Advances in*
307 *Neural Information Processing Systems*, 32:1–12, 2019. 2
- 308 [9] Shuo Yang, Erkun Yang, Bo Han, Yang Liu, Min Xu, Gang
309 Niu, and Tongliang Liu. Estimating instance-dependent
310 bayes-label transition matrix using a deep neural network.
311 In *International Conference on Machine Learning*, pages
312 25302–25312. PMLR, 2022. 4
- 313 [10] Yu Yao, Tongliang Liu, Bo Han, Mingming Gong, Jiankang
314 Deng, Gang Niu, and Masashi Sugiyama. Dual t: Reduc-
315 ing estimation error for transition matrix in label-noise learn-
316 ing. *Advances in neural information processing systems*, 33:
317 7260–7271, 2020. 2
- 318 [11] Yu Yao, Tongliang Liu, Mingming Gong, Bo Han, Gang Niu,
319 and Kun Zhang. Instance-dependent label-noise learning un-
320 der a structural causal model. *Advances in Neural Informa-*
321 *tion Processing Systems*, 34:4409–4420, 2021. 4
- 322 [12] Li Yi, Sheng Liu, Qi She, A Ian McLeod, and Boyu Wang.
323 On learning contrastive representations for learning with
324 noisy labels. In *Proceedings of the IEEE/CVF conference*
325 *on computer vision and pattern recognition*, pages 16682–
326 16691, 2022. 4