
DP-SSL: Towards Robust Semi-supervised Learning with A Few Labeled Samples

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The scarcity of labeled data is a critical obstacle to deep learning. Semi-supervised
2 learning (SSL) provides a promising way to leverage unlabeled data by pseudo
3 labels. However, when the size of labeled data is very small (say a few labeled
4 samples per class), SSL performs poorly and unstably, possibly due to the low
5 quality of learned pseudo labels. In this paper, we propose a new SSL method called
6 DP-SSL that adopts an innovative data programming (DP) scheme to generate
7 probabilistic labels for unlabeled data. Different from existing DP methods that
8 rely on human experts to provide initial labeling functions (LFs), we develop a
9 multiple-choice learning (MCL) based approach to automatically generate LFs
10 from scratch in SSL style. With the noisy labels produced by the LFs, we design
11 a label model to resolve the conflict and overlap among the noisy labels, and
12 finally infer probabilistic labels for unlabeled samples. Extensive experiments
13 on four standard SSL benchmarks show that DP-SSL can provide reliable labels
14 for unlabeled data and achieve better classification performance on test sets than
15 existing SSL methods, especially when only a small number of labeled samples
16 are available. Concretely, for CIFAR-10 with only 40 labeled samples, DP-SSL
17 achieves 93.82% annotation accuracy on unlabeled data and 93.46% classification
18 accuracy on test data, which are higher than the SOTA results.

19 1 Introduction

20 The de-facto approaches to deep learning achieve phenomenal success with the release of huge labeled
21 datasets. However, large manually-labeled datasets are time-consuming and expensive to acquire,
22 especially when expert labelers are required. Nowadays, many techniques are proposed to alleviate
23 the burden of manual labeling and help to train models from scratch, such as active learning [1],
24 crowd-labeling [2], distant supervision [3], semi [4]/weak [5]/self-supervision [6]. Among them,
25 semi-supervised learning (SSL) is one of the most popular techniques to cope with the scarcity of
26 labeled data. Two major strategies of SSL are pseudo labels [7] and consistency regularization [8].
27 Pseudo labels (also called self-training [9]) utilize a model's predictions as the labels to train the
28 model again, while consistency of regularization forces a model to make the same prediction under
29 different transformations. However, when the size of labeled data is small, SSL performance degrades
30 drastically in both accuracy and robustness. Fig. 1 shows the change of prediction error rate with
31 the number of labeled samples of CIFAR-10. When the number of labeled samples reduces from
32 250 to 40, error rates of major existing SSL methods increase from 4.74% (USADTM) to 36.49%
33 (MixMatch). One possible reason of performance deterioration is due to quality degradation of learnt
34 pseudo labels when labeled data size is small. Therefore, in this paper we address this problem by
35 developing sophisticated labeling techniques for unlabeled data to boost SSL even when the number
36 of labeled samples is very small (e.g. a few labeled samples per class).

37 Recently, *data programming* (DP) was proposed as a new
 38 paradigm of weak supervision [10]. In DP, human experts
 39 are required to transform the decision-making process into
 40 a series of small functions (called *labeling functions*, abbrevi-
 41 ated as LFs), thus data can be labeled programmatically.
 42 Besides, a label model is applied to determining the correct
 43 labels based on consensus from the noisy and conflicting
 44 labels assigned by the LFs. Such a paradigm achieves
 45 considerable success in NLP tasks [11–14]. In addition,
 46 DP has also been applied to computer vision tasks [15, 16].
 47 However, current DP methods require human experts to
 48 provide initial LFs, which is time-consuming and expen-
 49 sive, and it is not easy to guarantee the quality of LFs.
 50 Furthermore, LFs specifically defined for one task usually
 51 cannot be re-used for other tasks.

52 In this paper, we propose a new SSL method called DP-
 53 SSL that is effective and robust even when the number of
 54 labeled samples is very small. In DP-SSL, an innovative
 55 data programming (DP) scheme is developed to generate
 56 probabilistic labels for unlabeled data. Different from
 57 existing DP methods, we develop a *multiple-choice learn-*
 58 *ing* (MCL) based approach to automatically generate LFs from scratch in SSL style. To remedy the
 59 over-confidence problem with existing MCL methods, we assign an additional option as abstention for
 60 each LF. After that, we design a label model to resolve the conflict and overlap among the noisy labels
 61 generated by LFs, and infer a probabilistic label for each unlabeled sample. Finally, the probabilistic
 62 labels are used to train the end model for classifying unlabeled data. Our experiments validate the
 63 effectiveness and advantage of DP-SSL. As shown in Fig. 1, DP-SSL performs best, and only 1.76%
 64 increase of error rate when the size of labeled samples decreases from 250 to 40 in CIFAR-10.

65 Note that the pseudo labels used in existing SSL methods is quite different from the probabilistic
 66 labels in DP-SSL, which may explain the advantage of DP-SSL over existing SSL methods. On the
 67 one hand, pseudo labels are “hard” labels that indicate an unlabeled sample belonging to a certain
 68 class or not, while probabilistic labels are “soft” labels that indicate the class distributions of unlabeled
 69 samples. Obviously, the latter should be more flexible and robust. On the other hand, pseudo labels
 70 are actually generated by a single model for all unlabeled samples, while probabilistic labels are
 71 generated from a number of diverse and specialized LFs (due to the MCL mechanism), which makes
 72 the latter more powerful in generalization as a whole.

73 In summary, the contributions of this paper are as follows: 1) We propose a new SSL method DP-SSL
 74 that employs an innovative data programming method to generate probabilistic labels for unlabeled
 75 data, which makes DP-SSL effective and robust even when there are only a few labeled samples per
 76 class. 2) We develop a multiple choice learning based approach to automatically generate diverse and
 77 specialized LFs from scratch for unlabeled data in SSL manner. 3) We design a label model with a
 78 novel potential and an unsupervised quality guidance regularizer to infer probabilistic labels from the
 79 noisy labels generated by LFs. 4) We conduct extensive experiments on four standard benchmarks,
 80 which show that DP-SSL outperforms the state-of-the-art methods, especially when only a small
 81 number of labeled samples are available, DP-SSL is still effective and robust.

82 2 Related Work

83 Here we briefly review the latest advances in multiple choice learning, semi-supervised learning, and
 84 data programming, which are related to our work. Detailed information is available in [17–20].

85 2.1 Multiple Choice Learning

86 Multiple choice learning (MCL) [21] was proposed to overcome the low diversity problem of models
 87 trained independently in ensemble learning. For example, stochastic multiple choice learning [22] is
 88 for training diverse deep ensemble models. However, a crucial problem with MCL is that each model
 89 tends to be overconfident. which results in poor final prediction. To solve this problem, [23] forces

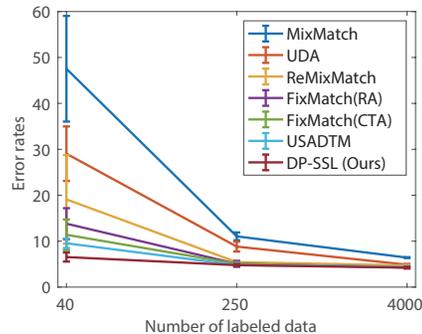


Figure 1: Error rate vs. #labeled samples (CIFAR-10). Results of existing methods are from the original papers. When only 40 labeled samples are given, all existing SSL methods are substantially degraded and more unstably, while our method is still effective and robust.

90 the predictions of non-specialized models to meet a uniform distribution, so that the final decision is
91 summed over diverse outputs. [24] proposes an additional network to estimate the weight of each
92 specialist’s output. In this paper, we develop an improved MCL based scheme to automatically
93 generate diverse and specialized labeling functions (LFs) from scratch in an SSL manner. These LFs
94 are used to generate preliminary (usually noisy) labels for unlabeled data.

95 2.2 Semi-supervised Learning

96 Semi-supervised learning (SSL) has been extensively studied in image classification [25], object
97 detection [26], and semantic segmentation [27]. Two popular SSL strategies for image classification
98 are pseudo labels [7] and consistency regularization [8]. Pseudo-label methods generate artificial
99 labels for some unlabeled images and then train the model with these artificial labels, while consistency
100 regularization tries to obtain an artificial distribution/label and applied it as a supervision signal with
101 other augmentations/views. These two strategies have been adopted by a number of recent SSL
102 works [4, 8, 28–37]. For example, FixMatch [4] proposes a simple combination of pseudo labels
103 and consistency regularization. [35] employs unsupervised learning and clustering to determine the
104 pseudo labels. In this paper, we propose a new SSL method that is effective and robust even when
105 the size of labeled data is very small. Our method employs an innovative data programming alike
106 method to automatically generate probabilistic labels for unlabeled data.

107 2.3 Data Programming

108 Data programming [10] is a weak supervision paradigm proposed to infer correct labels based on
109 the consensus among noisy labels from labeling functions (LFs), which are modules embedded with
110 decision-making processes for generating labels programmatically. Following the DP paradigm,
111 Snorkel [12] and Snuba [38] were proposed as a rapid training data creation system. Their LFs are
112 built with various weak supervision sources, like pattern regexes, heuristics, and external knowledge
113 base etc. Recently, more works are reported in the literature [11, 13–16, 20, 39–45]. Among
114 them, [11, 13, 14, 43–45] focus on the adaption of label model in DP. For example, [20] aims to
115 reduce the computational cost and proposes a closed-formed solution for training the label model.
116 [15, 16, 39–41] apply DP to computer vision. Concretely, [16, 40, 41] heavily rely on the pretrained
117 models. [39] combines crowdsourcing, data augmentation, and DP to create weak labels for image
118 classification. [15] presents a novel view for resolving infrequent data in scene graph prediction
119 training datasets via image-agnostic features in LFs. However, all these methods cannot directly
120 applied to training models from scratch with a small number of labeled samples. Thus, in this paper
121 we extend DP by exploring both MCL and SSL to generate arbitrary labeling functions.

122 3 Method

123 For a C -class SSL classification problem, assume that all training data X are divided into labeled data
124 X_l and unlabeled X_u , and test data are denoted as X_t . Following the notation in [4, 35], $\{x_l, x_l^w\} \in$
125 X_l are the paired labeled samples with labels $y_l \in \{1, \dots, C\}$, and $\{x_u, x_u^w, x_u^s\} \in X_u$ are the triple
126 unlabeled samples. Here, x_l and x_u represent the raw images without any transformations. $x_l^w, x_u^w,$
127 and x_u^s are the images based on the weak and strong augmentation strategies, respectively. In this
128 paper, weak augmentation uses a standard flip-and-shift strategy, and strong augmentation is the
129 RandAugment[46] strategy with Cutout [47] augmentation operation.

130 3.1 Framework

131 Fig. 2 shows the framework of our DP-SSL method, which works in three major steps as follows:

- 132 • Step 1. We employ an MCL based approach to automatically generate LFs from scratch
133 in an SSL style. Here, each LF is trained on a subset of C classes in training set based on
134 MCL. As shown in Fig 2, the 2nd LF is trained with samples of classes “horse” and “dog”,
135 and abstains from prediction when facing monkey images.
- 136 • Step 2. A graphical model is developed as the label model to aggregate the noisy labels and
137 produce probabilistic labels for unlabeled training data. The label model is learned in an
138 SSL manner with an additional regularizer.

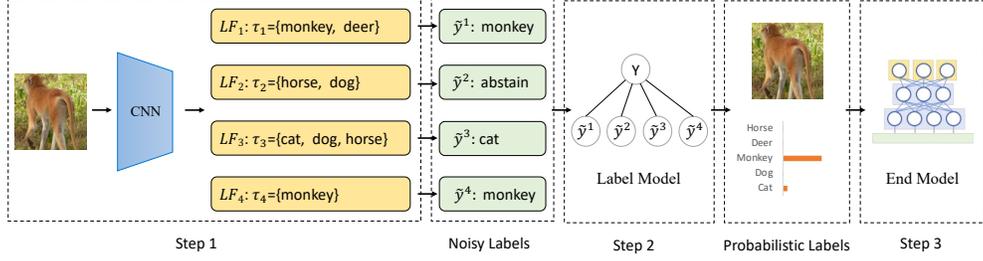


Figure 2: Framework of the DP-SSL method with four LFs.

- Step 3. The end model is trained with both provided labels and probabilistic labels generated from Step 2. Finally, we verify the performance of the end model on test data.

3.2 Labeling Function

In Step 1 of our method, LFs are exploited to generate noisy labels for each unlabeled image. In previous DP works for computer vision, LFs are built via external image-agnostic knowledge [15] or pretrained models [16, 40, 41]. However, it is difficult to explicitly describe the rules of image classification. Instead, here we innovatively explore MCL and SSL for automatic LF generation.

As shown in Fig. 2, we share the same backbone (Wide ResNet [48] in this paper) to extract features of images for multiple prediction heads (called LFs in this paper). To promote the diversity of LFs, we transform the features and feed each LF with different transformed features as follows:

$$\tilde{f}_k = \sum_{j=1}^{HW} \frac{e^{-\beta_k \Lambda(f_j, c_k)}}{\sum_{k'=1}^K e^{-\beta_{k'} \Lambda(f_j, c_{k'})}} (f_j - c_k). \quad (1)$$

In this paper, $f \in \mathbb{R}^{H \times W \times C}$ denotes the feature before global average pooling in the backbone, $f_j \in \mathbb{R}^C$ is the feature vector at position j of f . $c_k \in \mathbb{R}^C$ is the learnable clustering center in the k -th LF, β_k is the learnable variable of the k -th cluster, $\Lambda(A, B)$ is the distance metric for A and B . Thus, \tilde{f}_k corresponds to the feature fed to the k -th LF, and describes the k -th aggregated pattern of f among K centers c_k , it can also be considered as a learnable weighted average pooling for feature f . In our experiments, $e^{-\beta_k \Lambda(f_j, c_k)}$ is approximately implemented in the form as in [49].

As depicted in [22], the classifiers lack diversity of prediction even trained with different protocols. Therefore, we adopt MCL to assign a subset of labeled data for each classifier automatically to improve diversity, which is formulated as

$$\begin{aligned} \mathcal{L}_1(x_l^w, y_l) &= \sum_{k=1}^K u^k * H(y_l, \tilde{p}_l^{w,k}) \\ s.t. \quad \sum_{k=1}^K u^k &= \rho * K, u^k \in \{0, 1\}, \end{aligned} \quad (2)$$

with $H(y_l, \tilde{p}_l^{w,i}) \leq H(y_l, \tilde{p}_l^{w,j})$ for $\forall i, j$ when $u^i = 1, u^j = 0$. Above, $H(y_l, \tilde{p}_l^{w,k})$ denotes the cross entropy between the ground truth label y_l and the predicted distribution $\tilde{p}_l^{w,k}$ of the k -th LF, u^k is the indicator variable in MCL to indicate whether the k -th LF is specialized in labeling x_l^w . K is the number of LFs in DP and ρ is the ratio of specialist LFs. When ρ is equal to 1, MCL deteriorates to the basic ensemble learning, where all K classifiers are trained with the same data.

Based on MCL, each LF is a specialist for some classes, so it can get high accuracy for samples in these classes. While for samples from other classes not specialized by the LF, it fails to predict due to over-confidence. Thus, we allow each LF to abstain from some samples in the dataset. Formally, we denote 0 as the abstention label, and the specialized classes of the k -th LF as $\tau_k \in \{\tau_k^1, \dots, \tau_k^{|\tau_k|}\}$. Then, the output of the k -th LF \tilde{y}^k satisfies $\tilde{y}^k \in \{0\} \cup \tau_k$. For example, the output of the 1st LF in

168 Fig. 2 is among “monkey”, “deer” and “abstention” for its specialized category set is $\tau_1 = \{\text{monkey},$
 169 $\text{deer}\}$. The objective function over labeled samples with abstention option is

$$\mathcal{L}_2(x_l^w, y_l) = \sum_{k=1}^K (\mathbb{1}(y_l \in \tau_k)H(y_l, \tilde{q}_l^{w,k}) + \mathbb{1}(y_l \notin \tau_k)H(0, \tilde{q}_l^{w,k})), \quad (3)$$

170 where $\tilde{q}_l^{w,k}$ is the normalized probability distribution of $\text{CONCAT}(\tilde{p}_l^{w,k}, \tilde{p}_{l,a}^{w,k})$ and $\tilde{p}_{l,a}^{w,k}$ is the
 171 probability of abstention. Then, for the unlabeled training data, we follow the settings in FixMatch [4],
 172 where unlabeled data are supervised by the pseudo labels $\tilde{y}_u^{w,k}$ of weak augmentation data x_u^w . Thus,

$$\mathcal{L}(x_u^w, x_u^s) = \sum_{k=1}^K \mathbb{1}(\max(\tilde{q}_u^{w,k}) \geq \epsilon) \left(\mathbb{1}(\tilde{y}_u^{w,k} \in \tau_k)H(\tilde{y}_u^{w,k}, \tilde{q}_u^{s,k}) + \mathbb{1}(\tilde{y}_u^{w,k} \notin \tau_k)H(0, \tilde{q}_u^{s,k}) \right), \quad (4)$$

173 where $\tilde{y}_u^{w,k} = \text{arg max}(\tilde{q}_u^{w,k})$. Specifically, we keep only samples whose largest probability (includ-
 174 ing the abstention option) is above the predefined threshold ϵ (0.95 in our paper), and train the model
 175 on the kept data with pseudo label $\tilde{y}_u^{w,k}$. Accordingly, the training in this step is to minimize the
 176 objective function as follows:

$$\mathcal{L}(x_l^w, y_l, x_u^w, x_u^s) = \mu_1 \mathcal{L}_1(x_l^w, y_l) + \mu_2 \mathcal{L}_2(x_l^w, y_l) + \mu_u \mathcal{L}(x_u^w, x_u^s), \quad (5)$$

177 where μ_1 , μ_2 and μ_u are hyper-parameters. In our implementation, we first set $\mu_1 = 1$ and
 178 $\mu_2 = \mu_u = 0$, then adjust μ_1 to 0 and $\mu_2 = \mu_u = 1$ after the convergence of \mathcal{L}_1 .

179 Generally, in Step 1, MCL is expected to generate specialized class sets τ for LFs, with which samples
 180 are more easily discriminated by SSL classifiers even there are a few labeled samples. Besides, the
 181 abstention option is for addressing the over-confidence problem of samples from non-specialized sets.

182 3.3 Label Model

183 In Step 2 of our method, we utilize a graphical model to specify a single prediction by integrating
 184 noisy labels provided by K LFs. For simplification, we assume that the K LFs are independent (as
 185 shown in Fig. 2). Then, suppose that $\tilde{\mathbf{y}} \in \mathbb{R}^K$ is the vectorized form of the predictions from K LFs,
 186 the joint distribution of the label model can be described as:

$$P(y, \tilde{\mathbf{y}}) = \frac{1}{Z} \prod_{k=1}^K \phi(y, \tilde{y}^k) \quad (6)$$

187 where Z is the normalizer of the joint distribution, ϕ is the potential that couples the target y and noisy
 188 label \tilde{y}^k . In this paper, we extend the dimension of parameters θ in label model to $K \times C$ to support
 189 multi-class classification. Set $e_{ky} := \exp(\theta_{ky})$, which is the exponent of parameters θ_{ky} . Now we
 190 are to construct the potential function ϕ . Due to the specialized LFs, the potential ϕ should benefit
 191 the final prediction when a noisy label agrees with the target. That is, we should have $\phi(y, \tilde{y}^k) > 1$.
 192 Thus, we set ϕ as $1 + e_{ky}$ for this case. On the contrary, the potential ϕ should negatively impact the
 193 final prediction when a noisy label conflicts with the target label in the specialized category set, i.e.,
 194 we should have $\phi(y, \tilde{y}^k) < 1$. Therefore, for this case we set ϕ to $1/(1 + e_{ky})$. For the other cases,
 195 we follow the design in [14]. In summary, the potential ϕ is defined as follows:

$$\phi(y, \tilde{y}^k) = \begin{cases} 1 + e_{ky}, & \text{if } y \in \tau_k, \tilde{y}^k \in \tau_k, \tilde{y}^k = y \\ 1/(1 + e_{ky}), & \text{if } y \in \tau_k, \tilde{y}^k \in \tau_k, \tilde{y}^k \neq y \\ e_{ky}, & \text{if } y \notin \tau_k, \tilde{y}^k \in \tau_k, \tilde{y}^k \neq y \\ 1, & \text{otherwise} \end{cases} \quad (7)$$

196 With the potential above, the normalizer Z of the joint distribution in Eq. (6) can be obtained by
 197 summarizing over y and \tilde{y}^k :

$$\begin{aligned} Z &= \sum_{y \in \mathcal{Y}} \prod_{k=1}^K \sum_{\tilde{y}^k \in \{0\} \cup \tau_k} \phi(y, \tilde{y}^k) \\ &= \sum_{y \in \mathcal{Y}} \prod_{k=1}^K \left(\mathbb{1}(y \in \tau_k) (2 + e_{ky} + \frac{|\tau_k| - 1}{1 + e_{ky}}) + \mathbb{1}(y \notin \tau_k) (1 + |\tau_k| e_{ky}) \right). \end{aligned} \quad (8)$$

198 Then, the objective function of the label model can be expressed in an SSL manner as follows:

$$\mathcal{L}(\tilde{\mathbf{y}}_l, y_l, \tilde{\mathbf{y}}_u) = \underbrace{\sum_{x_l} H(y_l, P(y, \tilde{\mathbf{y}}_l))}_{\text{labeled samples}} + \underbrace{\left(-\sum_{x_u} \log \sum_{y \in \mathcal{Y}} P(y, \tilde{\mathbf{y}}_u)\right)}_{\text{unlabeled samples}} + R(\theta, \tilde{\mathbf{y}}_u), \quad (9)$$

199 where the first part is the cross-entropy loss, the second is the negative log marginal likelihood on the
 200 observed noisy labels \tilde{y}_u , and the third is a regularizer. In our method, the regularizer is utilized to
 201 guide the label model with statistical information (the accuracy of each LF). However, the accuracy
 202 of each LF on noisy labels is unavailable, while the accuracy on labeled training is almost 100% due
 203 to over-fitting. Thus, we have to estimate the accuracy of each LF with the observable noisy labels \tilde{y} ,
 204 which will be presented in Sec. 3.4. After training, the label model produces probabilistic labels π by
 205 computing the joint distribution in Eq. (6) with the noisy labels \tilde{y} .

206 3.4 Accuracy Estimation

207 Now, we formally describe our method for estimating the accuracy of LFs. We transform the multi-
 208 class problem into C one-versus-all tasks. For the i -th one-versus-all task, we denote the unobserved
 209 true labels as $Y_i \in \{\pm 1\}$ (+1 is the positive label, -1 represents that of the other categories), noisy
 210 labels of the k -th LF as $\lambda_i^k \in \{\pm 1, 0\}$ (0 for abstention). Then, we can write $\mathbb{E}[\lambda_i^k Y_i]$ as

$$\begin{aligned} \mathbb{E}[\lambda_i^k Y_i] &= P(\lambda_i^k Y_i = 1) - P(\lambda_i^k Y_i = -1) \\ &= P(\lambda_i^k Y_i = 1) - (1 - P(\lambda_i^k Y_i = 1) - P(\lambda_i^k Y_i = 0)) \\ &= 2P(\lambda_i^k = Y_i) + P(\lambda_i^k = 0) - 1. \end{aligned} \quad (10)$$

211 Assume that $\lambda_i^j \perp \lambda_i^k | Y_i$ for distinct j and k , then

$$\mathbb{E}[\lambda_i^j \lambda_i^k] = \mathbb{E}[\lambda_i^j Y_i^2 \lambda_i^k] = \mathbb{E}[\lambda_i^j Y_i] \mathbb{E}[\lambda_i^k Y_i] \quad (11)$$

212 with the fact that $Y_i^2 = 1$. In Eq. (11), $\hat{\mathbb{E}}[\lambda_i^j \lambda_i^k] = \frac{1}{|x_u|} \sum_{x_u} \lambda_i^j \lambda_i^k$ is observable, which can be derived
 213 from the noisy labels of the j -th and k -th LFs, while $\mathbb{E}[\lambda_i^j Y_i]$ and $\mathbb{E}[\lambda_i^k Y_i]$ remains to be solved due
 214 to true label Y_i is unavailable. Next, we introduce a third labeling result from the l -th LF as λ_i^l , such
 215 that $\hat{\mathbb{E}}[\lambda_i^j \lambda_i^l]$ and $\hat{\mathbb{E}}[\lambda_i^k \lambda_i^l]$ are observable. Now, $|\hat{\mathbb{E}}[\lambda_i^j Y_i]|$, $|\hat{\mathbb{E}}[\lambda_i^k Y_i]|$, $|\hat{\mathbb{E}}[\lambda_i^l Y_i]|$ can be solved by a
 216 triplet method as follows:

$$\begin{aligned} |\hat{\mathbb{E}}[\lambda_i^j Y_i]| &= \sqrt{|\hat{\mathbb{E}}[\lambda_i^j \lambda_i^k] \cdot \hat{\mathbb{E}}[\lambda_i^j \lambda_i^l] / \hat{\mathbb{E}}[\lambda_i^k \lambda_i^l]|}, \\ |\hat{\mathbb{E}}[\lambda_i^k Y_i]| &= \sqrt{|\hat{\mathbb{E}}[\lambda_i^j \lambda_i^k] \cdot \hat{\mathbb{E}}[\lambda_i^k \lambda_i^l] / \hat{\mathbb{E}}[\lambda_i^j \lambda_i^l]|}, \\ |\hat{\mathbb{E}}[\lambda_i^l Y_i]| &= \sqrt{|\hat{\mathbb{E}}[\lambda_i^j \lambda_i^l] \cdot \hat{\mathbb{E}}[\lambda_i^k \lambda_i^l] / \hat{\mathbb{E}}[\lambda_i^j \lambda_i^k]|}. \end{aligned} \quad (12)$$

217 We can obtain the estimated accuracy of each LF by resolving the sign of $\mathbb{E}[\lambda_i^k Y_i]$. Let $\hat{a}_i^k := \hat{P}(\lambda_i^k =$
 218 $Y_i | \lambda_i^k \neq 0)$ be the estimated accuracy of the k -th LF on the i -th category. Therefore, the regularizer
 219 of $R(\theta, \tilde{y}_u)$ can be formulated as

$$R(\theta, \tilde{\mathbf{y}}_u) = \sum_{i=1}^C \sum_k^K \hat{a}_i^k \log P_\theta(\lambda_i^k = Y_i | \tilde{y}_u^k \neq 0) + (1 - \hat{a}_i^k) \log(1 - P_\theta(\lambda_i^k = Y_i | \tilde{y}_u^k \neq 0)) \quad (13)$$

220 where $P_\theta(\lambda_i^k = Y_i | \tilde{y}_u^k \neq 0)$ can be computed in closed form by marginalizing over all the other
 221 variables in the model in Eq. (6) without noisy labels \tilde{y} . Details of P_θ can be referred to **Appendix**.

222 3.5 End Model

223 In Step 3, probabilistic labels are used to train an end model under any network architecture. We
 224 utilize noise-aware empirical risk expectation as the objective function to take annotation errors into
 225 account. Accordingly, the final objective function is as follows:

$$\mathcal{L}(x_l, y_l, x_u, \pi) = \underbrace{\sum_{x_l} H(y_l, p_l)}_{\text{labeled samples}} + \underbrace{\sum_{x_u} \mathbb{E}_{y \sim \pi} H(y, p_u)}_{\text{unlabeled samples with probabilistic label}} \quad (14)$$

226 where p_l and p_u are the predicted distributions of x_l and x_u , π is the distribution produced by the
 227 label model in Sec. 3.3. Actually, $H(y_l, p_l)$ can also be rewritten as $\mathbb{E}_{y \sim P(y_l)} H(y, p_l)$ where $P(y_l)$
 228 is the one-hot distribution of y_l .

229 4 Experiments

230 4.1 Implementation Details

231 In the training phase, we follow the settings of previous works [4, 34, 35], augment data in weak (a
 232 standard flip-and-shift strategy) and strong forms (RandAugment [46] followed by Cutout [47]
 233 operation), and utilize a Wide ResNet as the end model for a fair comparison. In our framework,
 234 the batch size for labeled data and unlabeled data is set to 64 and 448, respectively. Besides, we use
 235 the same hyperparameters ($K = 50$, $\rho = 0.2$, $\epsilon = 0.95$) for all datasets. We compare DP-SSL with
 236 major existing methods on CIFAR-10 [50], CIFAR-100 [50], SVHN [51] and STL-10 [52]. We also
 237 analyze the effect of annotation and conduct ablation study in Sec. 4.4 and Sec. 4.5 respectively. All
 238 experiments are implemented in Pytorch v1.7 and conducted on 16 NVIDIA RTX3090s.

239 4.2 Datasets

240 **CIFAR-10 and CIFAR-100** [50] contain 50,000 training examples and 10,000 validation examples.
 241 All images are of 32x32 pixel size and fall in 10 or 100 classes, respectively.

242 **SVHN** [51] is a digital image dataset that consists of 73,257, 26,032 and 531,131 samples in the
 243 train, test, and extra folders. It has the same image resolution and category number as CIFAR-10.

244 **STL-10** [52] is a dataset for evaluating unsupervised and semi-supervised learning. It consists of
 245 5000 labeled images and 8000 validation samples of 96x96 size from 10 classes. Besides, there are
 246 100,000 unlabeled images available, including odd samples.

247 4.3 Comparison with Existing SSL Methods

248 For a fair comparison, we conduct experiments with the codebase of FixMatch and cite the results on
 249 CIFAR-10, CIFAR-100, SVHN and STL-10 from [4, 35]. We utilize the same network architecture
 250 (a Wide ResNet-28-2 for CIFAR-10 and SVHN, WRN-28-8 for CIFAR-100, and WRN-37-2 for
 251 STL-10) and training protocol of FixMatch, such as optimizer and learning rate schedule. Unlabeled
 252 data are generated by the scripts in FixMatch. Results of DP-SSL and existing methods in Tab. 1 and
 253 Tab. 2 are presented with the mean and standard deviation (STD) of accuracy on 5 pre-defined folds.

254 As shown in Tab. 1, our method achieves the best performance in most cases, especially when
 255 there are only 4 labeled samples per class. Specifically, our method achieves a 93.46% accuracy on
 256 CIFAR-10 with 4 labeled samples per category, which is 3.3% higher than that of USADTM — the
 257 state-of-the-art method. Again on STL-10, our method surpasses USADTM and achieves the best
 258 performance when there are 4 and 25 labeled samples per class.

259 On CIFAR-100, our method performs the best for 40 labels case and the 2nd for 2500 and 10,000 labels
 260 cases. We also notice that DP-SSL has relatively large STDs for 2500 and 10,000 labels cases, which is
 261 due to the coarse accuracy estimation. In fact, even if triplet mean is adopted in estimation, the triplet
 262 selection in Eq. (12) still impacts accuracy estimation and regularizer a lot, especially when $\mathbb{E}[\lambda_i^k Y_i]$
 263 is close to 0 or sign recovery of $\mathbb{E}[\lambda_i^k Y_i]$ is wrong. Actually, there are some advanced approaches to
 264 unsupervised accuracy estimation [53–55] that can replace the naive triplet mean estimation. Ideally,
 265 if we can obtain the exact accuracy of each class $\hat{b}_i^k := \hat{P}(\lambda_i^k = Y_i | \lambda_i^k = 1)$ and regularize it as
 266 $R(\theta, \tilde{\mathbf{y}}_u) = \sum_{i=1}^C \sum_k^K \hat{b}_i^k \log P_\theta(\lambda_i^k = Y_i | \tilde{y}_u^k = i) + (1 - \hat{b}_i^k) \log(1 - P_\theta(\lambda_i^k = Y_i | \tilde{y}_u^k = i))$, we
 267 will get an end model with $(27.92 \pm \mathbf{0.23})\%$ error rate for 2500 labeled samples.

268 Comparing with USADTM, our method does not perform well enough when more labeled data
 269 available. For USADTM, apart from the proxy label generator, unsupervised representation learning
 270 contributes a lot for its performance. As shown in the ablation study of [35], USADTM without
 271 unsupervised representation learning achieves around 5.73% and 4.99% error rate for 250 and 4000
 272 labeled samples in CIFAR-10, while our method DP-SSL obtains 4.78% and 4.23% error rate.

Table 1: Results of error rate on CIFAR-10, CIFAR-100 and SVHN for different existing SSL methods (II-Model [28], Pseudo-Labeling [7], Mean Teacher [31], MixMatch [30], UDA [33], ReMixMatch [34], FixMatch [4] and USADTM [35]) and our DP-SSL method.

Method	CIFAR-10			CIFAR-100			SVHN		
	40 labels	250 labels	4000 labels	400 labels	2500 labels	10000 labels	40 labels	250 labels	1000 labels
II -Model	-	54.26±3.97	14.01±0.38	-	57.25±0.48	37.88±0.11	-	18.96±1.92	7.54±0.36
Pseudo-Labeling	-	49.78±0.43	16.09±0.28	-	57.38±0.46	36.21±0.19	-	20.21±1.09	9.94±0.61
Mean Teacher	-	32.32±2.30	9.19±0.19	-	53.91±0.57	35.83±0.24	-	3.57±0.11	3.42±0.07
MixMatch	47.54±11.50	11.05±0.86	6.42±0.10	67.61±1.32	39.94±0.37	28.31±0.33	42.55±14.53	3.98±0.23	3.50±0.28
UDA	29.05±5.93	8.82±1.08	4.88±0.18	59.28±0.88	33.13±0.22	24.50±0.25	52.63±20.51	5.69±2.76	2.46±0.24
ReMixMatch	19.10±9.64	5.44±0.05	4.72±0.13	44.28±2.06	27.43 ±0.31	23.03±0.56	3.34±0.20	2.92±0.48	2.65±0.08
FixMatch (RA)	13.81±3.37	5.07±0.65	4.26±0.05	48.85±1.75	28.29±0.11	22.60±0.12	3.96±2.17	2.48±0.38	2.28±0.11
FixMatch (CTA)	11.39±3.35	5.07±0.33	4.31±0.15	49.95±3.01	28.64±0.24	23.18±0.11	7.65±7.65	2.64±0.64	2.36±0.19
USADTM	9.54±1.04	4.80±0.32	4.40±0.15	43.36±1.89	28.11±0.21	21.35 ±0.17	3.01±1.97	2.11 ±0.65	1.96 ±0.05
DP-SSL (ours)	6.54 ±0.98	4.78 ±0.26	4.23 ±0.20	43.17 ±1.29	28.00±0.79	22.24±0.31	2.98 ±0.86	2.16±0.36	1.99±0.18
Fully Supervised		2.74			16.84			1.48	

Table 2: Results of error rate on STL-10.

STL-10							
Method	1000 labels	Method	1000 labels	Method	40 labels	250 labels	1000 labels
II -Model	26.23±0.82	UDA	7.66±0.56	USADTM	9.63±1.35	6.85±1.09	4.01 ±0.59
Pseudo-Labeling	27.99±0.80	ReMixMatch	5.23±0.45	DP-SSL (ours)	9.32 ±0.91	6.83 ±0.71	4.97±0.42
Mean Teacher	21.43±2.39	FixMatch (RA)	7.98±1.50	Fully Supervised		1.48	
MixMatch	10.41±0.61	FixMatch (CTA)	5.17±0.63				

273 4.4 Analysis

274 **Annotation performance.** Intuitively, the holistic performance of the end model in our method
275 highly depends on the quality of annotation results. Thus, we present the macro precision/recall/F1
276 score and coverage of the annotated labels of our method on CIFAR-10, CIFAR-100, and SVHN
277 in Tab. 3. We can see that our method achieves over 99% coverage, which means that it produces
278 probabilistic labels for almost all unlabeled data. Comparing to the results in [35], the label model
279 with 40 labeled samples outperforms the proxy label generator, FixMatch and USADTM get 88.51%
280 and 89.48% accuracy, respectively. Furthermore, our method achieves 97.36% accuracy for unlabeled
281 data with the top-500 highest probabilities in each category. Meanwhile, we also present results of
282 Majority Voting and FlyingSquid [20] in Tab. 3 based on the noisy labels from Step 1 of our method
283 for comparison. Majority Voting gets bad performance because the number of LFs triggered for
284 different categories is not equal. For FlyingSquid, we implement it with C one-versus-all models to
285 support multi-class tasks, and the large C in CIFAR-100 results in the worst performance.

286 **Barely supervised learning.** We conduct experiments to test the performance (accuracy and STD)
287 of our method on CIFAR-10 for some extreme cases (10, 20 and 30 labeled samples) to verify
288 the effectiveness of our method. Here, we select the labeled data through the scripts of FixMatch
289 with 5 different random seeds. As claimed in FixMatch, it reaches between 48.58% and 85.32%
290 test accuracy with a median of 64.28% for 10 labeled samples, while our method obtains accuracy
291 from 61.32% to 83.7%. As for 20 and 30 labeled samples, our method gets $(85.29 \pm 3.14)\%$ and
292 $(89.81 \pm 1.59)\%$ accuracy respectively, which have much smaller STDs than that reported in [36].

Table 3: The macro Precision/Recall/F1 Score/Coverage of the annotated labels on CIFAR-10, CIFAR-100, and SVHN for our method and two typical existing label models.

Method	Metrics	CIFAR-10			CIFAR-100			SVHN		
		40 labels	250 labels	4000 labels	400 labels	2500 labels	10000 labels	40 labels	250 labels	1000 labels
Majority Vote	F1 Score	85.96	94.23	95.77	49.97	69.81	76.03	90.86	95.38	96.14
FlyingSquid[20]	F1 Score	90.25	94.99	95.85	48.90	69.73	74.12	93.92	97.24	97.70
DP-SSL (ours)	Precision	93.47	95.30	95.89	55.62	71.91	75.12	95.20	97.65	97.79
	Recall	93.82	95.33	95.91	56.86	72.01	78.35	96.78	97.64	97.94
	F1 Score	93.61	95.19	95.90	54.42	71.89	76.36	95.95	97.59	97.81
	Coverage	99.35	99.79	99.91	99.33	99.87	99.94	99.15	99.67	99.93

293 **4.5 Ablation Study**

294 In DP-SSL, LFs and the label model are the core com-
 295 ponents to assign probabilistic labels for training the end
 296 model. Here, we check the effects of the following factors
 297 in the process of producing probabilistic labels by taking
 298 CIFAR-10 as the example. For ease of exposition, only the
 299 accuracy of predicted labels is presented in Tab. 4.

300 **MCL.** Feature transformation (FT) described in Eq. (1)
 301 can be regarded as a weighted spatial pooling for extracted
 302 features. It is proposed to boost the diversity of generated
 303 LFs. We conduct comparative experiments for three con-
 304 figurations: 1) *Exp1*: w.o. MCL, 2) *Exp2*: MCL w.o. FT,
 305 3) *Exp3*: MCL w. FT. The results are presented in Tab. 4. It
 306 is interesting to see that *Exp1* is better than *Exp2* but worse
 307 than *Exp3*. In fact, *Exp1* is a simple ensemble model with
 308 a shared backbone, where each LF is trained independently
 309 and predicts the labels within C categories. In *Exp2*, we
 310 observe that some classifiers have never been optimized in
 311 the training phase and thus have an empty specialized set
 312 when only a few labeled samples per class are available.
 313 Moreover, the specialized sets of many LFs are duplicate,
 314 which incurs a negative impact on the performance. How-
 315 ever, MCL with FT addresses the drawbacks and helps our method obtain versatile LFs. Some
 316 detailed examples are presented in **Appendix**

317 **Hyperparameters.** K and ρ are the number of LFs and the ratio of specialists in Eq. (2). *Exp4-13*
 318 in Tab. 4 present the variance of performance for different K and ρ . In *Exp4-8*, *Exp6* with $K=50$
 319 performs the best when 40 labeled samples are available, while *Exp8* with $K=100$ wins the others
 320 when 250 labeled samples are provided. To trade off the cost of computation and performance, we
 321 set $K=50$ as the default setting in our paper. On the other hand, performance reaches the best from
 322 *Exp9* to *Exp13* when $\rho=0.2$. Actually, when ρ is 1.0, MCL becomes a naive ensemble model.

323 **Regularizer.** The regularizer is proposed to impose a global guidance and improve the robustness of
 324 the label model. As shown in Tab. 4, the regularizer does boost the accuracy, especially when facing
 325 less labeled samples. Besides, as mentioned in Sec. 4.3, the high-quality guidance of the regularizer
 326 also reduces the label model’s performance variance, thus improves its robustness.

327 **5 Conclusion**

328 In this paper, we explore the data programming idea to boost SSL when only a small number of labeled
 329 samples available by providing more accurate labels for unlabeled data. To this end, we propose
 330 a new SSL method DP-SSL that employs an innovative DP mechanism to automatically generate
 331 labeling functions. To make the labeling functions diverse and specialized, a multiple choice learning
 332 based approach is developed. Furthermore, we design an effective label model by incorporating a
 333 novel potential and a regularizer with estimated accuracy. With this model, probabilistic labels are
 334 inferred by resolving the conflict and overlap among noisy labels from the labeling functions. Finally,
 335 an end model is trained under the supervision of the probabilistic labels. Extensive experiments show
 336 that DP-SSL can produce high-quality probabilistic labels, and outperforms the existing methods to
 337 achieve a new SOTA, especially when only a small number of labeled samples are available.

338 **6 Limitations of This Work**

339 In this work, we use coarse accuracy estimation as the statistic information to guide the label model
 340 for simplicity. As described in Sec. 3.4, we estimate the accuracy $P_\theta(\lambda_i^k = Y_i | \lambda_i^k \neq 0)$, rather than
 341 class-wise accuracy $P_\theta(\lambda_i^k = Y_i | \lambda_i^k = 1)$. Besides, we do not consider the dependency between C
 342 one-vs-all tasks.

Table 4: Annotation performance for different configurations on CIFAR-10 with 40 and 250 labels. K and ρ are set to 50 and 0.2 by default.

Experiments	40 labels	250 labels
Exp1: w.o. MCL	92.46	95.02
Exp2: MCL w.o. FT	91.61	94.98
Exp3: MCL w. FT	93.82	95.33
Exp4: $K=20, \rho=0.2$	91.27	94.75
Exp5: $K=40, \rho=0.2$	93.46	95.03
Exp6: $K=50, \rho=0.2$	93.82	95.33
Exp7: $K=60, \rho=0.2$	93.54	95.28
Exp8: $K=100, \rho=0.2$	93.68	95.35
Exp9: $K=50, \rho=0.1$	92.95	94.90
Exp10: $K=50, \rho=0.2$	93.82	95.33
Exp11: $K=50, \rho=0.3$	93.55	95.23
Exp12: $K=50, \rho=0.5$	93.07	94.86
Exp13: $K=50, \rho=1.0$	92.46	95.02
Exp14: w.o. Regularizer	93.19	94.94
Exp15: Regularizer	93.82	95.33

References

- 343
- 344 [1] M. Gao, Z. Zhang, G. Yu, S. Ö. Arık, L. S. Davis, and T. Pfister, “Consistency-based semi-supervised
345 active learning: Towards minimizing labeling cost,” in *ECCV*. Springer, 2020, pp. 510–526.
- 346 [2] C. Vondrick, D. Patterson, and D. Ramanan, “Efficiently scaling up crowdsourced video annotation,” *IJCV*,
347 vol. 101, no. 1, pp. 184–204, 2013.
- 348 [3] Y. Yao, A. Zhang, X. Han, M. Li, C. Weber, Z. Liu, S. Wermter, and M. Sun, “Visual distant supervision
349 for scene graph generation,” *arXiv preprint arXiv:2103.15365*, 2021.
- 350 [4] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Kurakin, and C.-L.
351 Li, “Fixmatch: Simplifying semi-supervised learning with consistency and confidence,” *NeurIPS*, vol. 33,
352 2020.
- 353 [5] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, “Is object localization for free?-weakly-supervised learning
354 with convolutional neural networks,” in *CVPR*, 2015, pp. 685–694.
- 355 [6] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual
356 representations,” in *ICML*. PMLR, 2020, pp. 1597–1607.
- 357 [7] D.-H. Lee *et al.*, “Pseudo-label: The simple and efficient semi-supervised learning method for deep neural
358 networks,” in *Workshop on ICML*, vol. 3, no. 2, 2013.
- 359 [8] M. Sajjadi, M. Javanmardi, and T. Tasdizen, “Regularization with stochastic transformations and perturba-
360 tions for deep semi-supervised learning,” in *NeurIPS*, 2016.
- 361 [9] B. Zoph, G. Ghiasi, T.-Y. Lin, Y. Cui, H. Liu, E. D. Cubuk, and Q. Le, “Rethinking pre-training and
362 self-training,” *NeurIPS*, vol. 33, 2020.
- 363 [10] A. Ratner, C. De Sa, S. Wu, D. Selsam, and C. Ré, “Data programming: Creating large training sets,
364 quickly,” *NeurIPS*, vol. 29, p. 3567, 2016.
- 365 [11] A. Awasthi, S. Ghosh, R. Goyal, and S. Sarawagi, “Learning from rules generalizing labeled exemplars,”
366 *ICLR*, 2020.
- 367 [12] A. J. Ratner, S. H. Bach, H. R. Ehrenberg, and C. Ré, “Snorkel: Fast training set generation for information
368 extraction,” in *SIGMOD*, 2017, pp. 1683–1686.
- 369 [13] A. Ratner, B. Hancock, J. Dunnmon, F. Sala, S. Pandey, and C. Ré, “Training complex models with
370 multi-task weak supervision,” in *AAAI*, vol. 33, no. 01, 2019, pp. 4763–4771.
- 371 [14] O. Chatterjee, G. Ramakrishnan, and S. Sarawagi, “Data programming using continuous and quality-guided
372 labeling functions,” *AAAI*, 2020.
- 373 [15] V. S. Chen, P. Varma, R. Krishna, M. Bernstein, C. Re, and L. Fei-Fei, “Scene graph prediction with limited
374 labels,” in *ICCV*, 2019, pp. 2580–2590.
- 375 [16] S. Hooper, M. Wornow, H. S. Ying, H. Kellman, Peter and Xue, F. Sala, C. Langlotz, and C. Ré, “Cut out
376 the annotator, keep the cutout: better segmentation with weak supervision,” *ICLR*, 2021.
- 377 [17] N. C. Garcia, S. A. Bargal, V. Ablavsky, P. Morerio, V. Murino, and S. Sclaroff, “Distillation multiple
378 choice learning for multimodal action recognition,” in *WACV*, 2021, pp. 2755–2764.
- 379 [18] J. E. Van Engelen and H. H. Hoos, “A survey on semi-supervised learning,” *Machine Learning*, vol. 109,
380 no. 2, pp. 373–440, 2020.
- 381 [19] M. F. Chen, D. Y. Fu, F. Sala, S. Wu, R. T. Mullapudi, F. Poms, K. Fatahalian, and C. Ré, “Train and you’ll
382 miss it: Interactive model iteration with weak supervision and pre-trained embeddings,” *arXiv preprint*
383 *arXiv:2006.15168*, 2020.
- 384 [20] D. Fu, M. Chen, F. Sala, S. Hooper, K. Fatahalian, and C. Ré, “Fast and three-rious: Speeding up weak
385 supervision with triplet methods,” in *ICML*. PMLR, 2020, pp. 3280–3291.
- 386 [21] A. Guzman-Rivera, P. Kohli, D. Batra, and R. Rutenbar, “Efficiently enforcing diversity in multi-output
387 structured prediction,” in *Artificial Intelligence and Statistics*. PMLR, 2014, pp. 284–292.
- 388 [22] S. Lee, S. Purushwalkam, M. Cogswell, V. Ranjan, D. J. Crandall, and D. Batra, “Stochastic multiple
389 choice learning for training diverse deep ensembles,” in *NeurIPS*, 2016.

- 390 [23] K. Lee, C. Hwang, K. Park, and J. Shin, “Confident multiple choice learning,” in *ICML*. PMLR, 2017, pp.
391 2014–2023.
- 392 [24] K. Tian, Y. Xu, S. Zhou, and J. Guan, “Versatile multiple choice learning and its application to vision
393 computing,” in *CVPR*, 2019, pp. 6349–6357.
- 394 [25] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image
395 database,” in *CVPR*. IEEE, 2009, pp. 248–255.
- 396 [26] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft
397 coco: Common objects in context,” in *ECCV*. Springer, 2014, pp. 740–755.
- 398 [27] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and
399 B. Schiele, “The cityscapes dataset for semantic urban scene understanding,” in *CVPR*, 2016, pp. 3213–
400 3223.
- 401 [28] S. Laine and T. Aila, “Temporal ensembling for semi-supervised learning,” in *ICLR*, 2017.
- 402 [29] A. Rasmus, H. Valpola, M. Honkala, M. Berglund, and T. Raiko, “Semi-supervised learning with ladder
403 networks,” in *NeurIPS*, 2015.
- 404 [30] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. Raffel, “Mixmatch: A holistic
405 approach to semi-supervised learning,” in *NeurIPS*, 2019.
- 406 [31] A. Tarvainen and H. Valpola, “Mean teachers are better role models: Weight-averaged consistency targets
407 improve semi-supervised deep learning results,” in *NeurIPS*, 2017.
- 408 [32] T. Miyato, S.-i. Maeda, M. Koyama, and S. Ishii, “Virtual adversarial training: a regularization method for
409 supervised and semi-supervised learning,” *TPAMI*, vol. 41, no. 8, pp. 1979–1993, 2018.
- 410 [33] Q. Xie, Z. Dai, E. Hovy, T. Luong, and Q. Le, “Unsupervised data augmentation for consistency training,”
411 in *NeurIPS*, vol. 33, 2020.
- 412 [34] D. Berthelot, N. Carlini, E. D. Cubuk, A. Kurakin, K. Sohn, H. Zhang, and C. Raffel, “Remixmatch:
413 Semi-supervised learning with distribution alignment and augmentation anchoring,” in *ICLR*, 2020.
- 414 [35] T. Han, J. Gao, Y. Yuan, and Q. Wang, “Unsupervised semantic aggregation and deformable template
415 matching for semi-supervised learning,” in *NeurIPS*, 2020.
- 416 [36] J. Li, C. Xiong, and S. Hoi, “Comatch: Semi-supervised learning with contrastive graph regularization,”
417 *arXiv preprint arXiv:2011.11183*, 2020.
- 418 [37] Z. Hu, Z. Yang, X. Hu, and R. Nevatia, “Simple: Similar pseudo label exploitation for semi-supervised
419 classification,” *arXiv preprint arXiv:2103.16725*, 2021.
- 420 [38] P. Varma and C. Ré, “Snuba: automating weak supervision to label training data,” in *VLDB*, vol. 12, no. 3.
421 NIH Public Access, 2018, p. 223.
- 422 [39] G. Heo, Y. Roh, S. Hwang, D. Lee, and S. E. Whang, “Inspector gadget: A data programming-based
423 labeling system for industrial images,” *VLDB*, 2020.
- 424 [40] A. Pal and V. N. Balasubramanian, “Adversarial data programming: Using gans to relax the bottleneck of
425 curated labeled data,” in *CVPR*, 2018, pp. 1556–1565.
- 426 [41] N. Das, S. Chaba, R. Wu, S. Gandhi, D. H. Chau, and X. Chu, “Goggles: Automatic image labeling with
427 affinity coding,” in *Proceedings of the 2020 ACM SIGMOD International Conference on Management of*
428 *Data*, 2020, pp. 1717–1732.
- 429 [42] B. Boecking, W. Neiswanger, E. Xing, and A. Dubrawski, “Interactive weak supervision: Learning useful
430 heuristics for data labeling,” *ICLR*, 2021.
- 431 [43] P. Varma, B. He, P. Bajaj, I. Banerjee, N. Khandwala, D. L. Rubin, and C. Ré, “Inferring generative model
432 structure with static analysis,” *NeurIPS*, vol. 30, p. 239, 2017.
- 433 [44] S. H. Bach, B. He, A. Ratner, and C. Ré, “Learning the structure of generative models without labeled
434 data,” in *ICML*. PMLR, 2017, pp. 273–282.
- 435 [45] P. Varma, F. Sala, A. He, A. Ratner, and C. Ré, “Learning dependency structures for weak supervision
436 models,” in *ICML*. PMLR, 2019, pp. 6418–6427.

- 437 [46] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, “Randaugment: Practical automated data augmentation
438 with a reduced search space,” in *Workshop on CVPR*, 2020, pp. 702–703.
- 439 [47] T. DeVries and G. W. Taylor, “Improved regularization of convolutional neural networks with cutout,”
440 *arXiv preprint arXiv:1708.04552*, 2017.
- 441 [48] S. Zagoruyko and N. Komodakis, “Wide residual networks,” in *BMVC*. BMVC, 2016.
- 442 [49] R. Girdhar, D. Ramanan, A. Gupta, J. Sivic, and B. Russell, “Actionvlad: Learning spatio-temporal
443 aggregation for action classification,” in *CVPR*, 2017, pp. 971–980.
- 444 [50] A. Krizhevsky, “Learning multiple layers of features from tiny images,” *Master’s thesis, University of
445 Tront*, 2009.
- 446 [51] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, “Reading digits in natural images with
447 unsupervised feature learning,” in *Workshop on NIPS*, 2011.
- 448 [52] A. Coates, A. Ng, and H. Lee, “An analysis of single-layer networks in unsupervised feature learning,” in
449 *Proceedings of the 14th international conference on artificial intelligence and statistics*. JMLR Workshop
450 and Conference Proceedings, 2011, pp. 215–223.
- 451 [53] A. Jaffe, B. Nadler, and Y. Kluger, “Estimating the accuracies of multiple classifiers without labeled data,”
452 in *Artificial Intelligence and Statistics*. PMLR, 2015, pp. 407–415.
- 453 [54] E. Platanios, H. Poon, T. M. Mitchell, and E. J. Horvitz, “Estimating accuracy from unlabeled data: A
454 probabilistic logic approach,” *NeurIPS*, vol. 30, pp. 4361–4370, 2017.
- 455 [55] P. A. Traganitis, A. Pages-Zamora, and G. B. Giannakis, “Blind multiclass ensemble classification,” *IEEE
456 Transactions on Signal Processing*, vol. 66, no. 18, pp. 4737–4752, 2018.

457 Checklist

- 458 1. For all authors...
- 459 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
460 contributions and scope? [Yes] See in Line.5-7 in abstract and the last paragraph in
461 Sec. 1.
- 462 (b) Did you describe the limitations of your work? [Yes] See in Sec. 6.
- 463 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 464 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
465 them? [Yes] Our paper conforms to the ethics review guidelines.
- 466 2. If you are including theoretical results...
- 467 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 468 (b) Did you include complete proofs of all theoretical results? [Yes]
- 469 3. If you ran experiments...
- 470 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
471 imental results (either in the supplemental material or as a URL)? [Yes] See major
472 experimental settings and data in Sec. 4.1 and Sec. 4.2. Code would be included in the
473 supplemental material.
- 474 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
475 were chosen)? [Yes] We follow the data splits of FixMatch [4] in Sec. 4.3. For common
476 hyperparameters, we use the default value in FixMatch [4]. While for characteristic
477 hyperparameters, see in Sec. 4.5.
- 478 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
479 ments multiple times)? [Yes] The error bars are attached in Tab. 1, and Tab. 2.
- 480 (d) Did you include the total amount of compute and the type of resources used (e.g., type
481 of GPUs, internal cluster, or cloud provider)? [Yes] See the last sentence of 4.1.
- 482 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 483 (a) If your work uses existing assets, did you cite the creators? [Yes] We have cited all of
484 them.

- 485 (b) Did you mention the license of the assets? [N/A]
486 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
487 We will attach our codes and some of the pretrained model in supplemental material.
488 (d) Did you discuss whether and how consent was obtained from people whose data you're
489 using/curating? [Yes] Datasets in our paper are publicly available
490 (e) Did you discuss whether the data you are using/curating contains personally identifiable
491 information or offensive content? [N/A]
492 5. If you used crowdsourcing or conducted research with human subjects...
493 (a) Did you include the full text of instructions given to participants and screenshots, if
494 applicable? [N/A]
495 (b) Did you describe any potential participant risks, with links to Institutional Review
496 Board (IRB) approvals, if applicable? [N/A]
497 (c) Did you include the estimated hourly wage paid to participants and the total amount
498 spent on participant compensation? [N/A]