

---

# Aligned Objective for Soft-Pseudo-Label Generation in Supervised Learning

---

Ning Xu<sup>1,2</sup> Yihao Hu<sup>1,2</sup> Congyu Qiao<sup>1,2</sup> Xin Geng<sup>1,2</sup>

## Abstract

Soft pseudo-labels, generated by the softmax predictions of the trained networks, offer a probabilistic rather than binary form, and have been shown to improve the performance of deep neural networks in supervised learning. Most previous methods adopt classification loss to train a classifier as the soft-pseudo-label generator and fail to fully exploit their potential due to the misalignment with the target of soft-pseudo-label generation, aimed at capturing the knowledge in the data rather than making definitive classifications. Nevertheless, manually designing an effective objective function for a soft-pseudo-label generator is challenging, primarily because datasets typically lack ground-truth soft labels, complicating the evaluation of the soft pseudo-label accuracy. To deal with this problem, we propose a novel framework that alternately trains the predictive model and the soft-pseudo-label generator guided by a meta-network-parameterized label enhancement objective. The parameters of the objective function are optimized based on the feedback from both the performance of the predictive model and the soft-pseudo-label generator in the learning task. Additionally, the framework offers versatility across different learning tasks by allowing direct modifications to the task loss. Experiments on the benchmark datasets validate the effectiveness of the proposed framework. Source code is available at <https://github.com/palm-ml/SEAL>

## 1. Introduction

Soft pseudo-labels, i.e., the soft labels that are not originally included in the training dataset but generated by the softmax predictions of the trained networks, are of paramount importance for state-of-the-art deep learning methods. Numerous previous researches (Zhang et al., 2019; Hinton et al., 2015) have demonstrated the success of soft pseudo-labels in improving the generalization performance of deep neural networks (DNNs) in supervised learning. As pseudo-labels could improve the predictive performance of deep neural networks, they have been successfully applied across various domains, including computer vision (Lukov et al., 2022; Algan & Ulusoy, 2022), natural language processing (Sun et al., 2019; Ngo et al., 2022), and data mining (Li et al., 2015; Xu et al., 2020).

Self-knowledge distillation (Zhang et al., 2019) utilizes a DNNs classifier’s predictions as the soft pseudo-labels to train the model itself to improve the test accuracy. Online distillation (Zhang et al., 2018) adopts an ensemble of classifiers’ predictions as the soft pseudo-label to learn mutually from each other. Label enhancement (Xu et al., 2023) adopts Graph Convolutional Network (GCN) to generate soft pseudo-labels to deal with the multi-label problem. In addition, soft pseudo-labels are adopted to learn with weak supervision. PENCIL (Yi & Wu, 2019) utilizes back-propagation to probabilistically update and correct soft pseudo-labels via updating the network parameters to deal with noisy labels in the training dataset. To identify the true label, partial label learning (Lv et al., 2020) generates soft pseudo-labels to strengthen the weight of the correct label among the candidate label set for training in each epoch.

Most previous methods adopt classification loss to train a classifier as the soft-pseudo-label generator, which is not aligned with the target of soft-pseudo-label generation. The objective of training a classifier is maximizing the generalization of the classifier itself while training a soft-pseudo-label generator is aimed at capturing the knowledge in the data to generate the appropriate soft pseudo-labels for improving the performance of the student models. Therefore, despite making an effort to tune the generated soft pseudo-labels, previous methods usually do not unleash the full potential of soft pseudo-labels for improving the performance of deep neural networks. Nevertheless, manually designing

---

<sup>1</sup>School of Computer Science and Engineering, Southeast University, Nanjing, China <sup>2</sup>Key Laboratory of New Generation Artificial Intelligence Technology and Its Interdisciplinary Applications (Southeast University), Ministry of Education, China. E-mail: {xning, yhhu, qiaocy, xgeng}@seu.edu.cn. Correspondence to: Xin Geng <xgeng@seu.edu.cn>.

an effective objective function for a soft-pseudo-label generator is challenging, primarily because datasets typically lack ground-truth soft labels, complicating the evaluation of the soft pseudo-label accuracy.

To deal with the problem, instead of focusing on manually designing an explicit loss for a classifier to output soft pseudo-labels, we propose a novel framework named SEAL, i.e., Soft-psEudo-lAbeL generator with a learnable label enhancement (Xu et al., 2019; 2020) objective, to alternately train the predictive model and the soft-pseudo-label generator guided by a meta-network-parameterized objective function. Specifically, a meta-network dynamically adjusts the parameters of the loss function used for a soft-pseudo-label generator. This optimization is informed by feedback from both the predictive model and the soft-pseudo-label generator, based on their performance in the learning task, creating a feedback loop that refines the training process. Additionally, our framework offers versatility across different learning tasks by allowing direct modifications to the task loss, making it suitable for a range of applications. Our contributions can be summarized as follows:

- We propose a novel framework named SEAL to generate soft pseudo-labels for improving the performance of deep neural networks in supervised learning via inducing a soft-pseudo-label generator optimized by a loss function parameterized by a meta-network. Extensive experiments validate the effectiveness of SEAL.
- We theoretically demonstrate that the soft pseudo-labels could enable the predictive model to achieve a larger sample margin for classes. This expansion, in turn, results in a more robust and tight generalization of the predictive model when trained using soft pseudo-labels.
- SEAL is flexible to learning tasks as the task loss designed for the supervised learning problems can be modified directly. We empirically show that the framework can be applied to partial label learning problems and achieve state-of-the-art learning performance.

## 2. Related Work

In supervised learning, the generation process of soft pseudo-labels exists in many mainstream algorithms, including knowledge distillation, self-knowledge distillation, label enhancement, and label smoothing.

In knowledge distillation (Hinton et al., 2015), a large teacher model is trained on the dataset with the original hard label and generates soft pseudo-labels to teach a lightweight student model. Self-distillation (Zhang et al., 2019) generates soft pseudo-labels by the model itself, which claims that the soft pseudo-labels can be regarded as the model’s deepest section’s output used to guide the training of shallow sections. Xu et al. (Xu & Liu, 2019) leverage distorted

versions of images to generate soft pseudo-labels for the current batch images. Yun et al. (Yun et al., 2020) introduce class-wise self-distillation that randomly chooses a different sample of the same class to generate a soft pseudo-label for the current sample. Building on this, Zhang et al. (Zhang et al., 2021a) propose to generate soft pseudo-labels based on the statistics of the model prediction for the target category. Taking a different approach, Kim et al. (Kim et al., 2021) suggest generating soft pseudo-labels adaptively by combining the ground-truth and past predictions from the model itself. Some methods also have low training costs for generating soft pseudo-labels. Liang et al. (Liang et al., 2022) propose an efficient self-distillation method that uses the on-the-fly prediction of a network to generate soft pseudo-labels that conform to a distribution. Additionally, Shen et al. (Shen et al., 2022) rearrange the sequential sampling by constraining half of each mini-batch coinciding with the previous iteration. The soft pseudo-labels are generated from the previous iteration, which is computationally efficient.

Label smoothing (Szegedy et al., 2016) prevents the network from becoming over-confident by regulating the model training by substituting one-hot labels with smoothed alternatives, which can also be regarded as soft pseudo-labels. Rather than using a uniform distribution, low-rank adaptive label smoothing (Ghoshal et al., 2021) adopts a more informative noise distribution to generate the soft pseudo-label for each class. To address the challenges in multi-label learning (Zhang & Zhou, 2014), Label enhancement (LE) (Xu et al., 2020) recovers label distribution for training multi-label classifiers, which could also be regarded as soft pseudo-labels. The Graph-Laplacian-based LE method (Xu et al., 2019) utilizes a local similarity matrix to maintain the structural integrity of the feature space, thereby converting discrete labels into soft labels. Zhang et al. (Zhang et al., 2021b) implemented a label propagation technique to disseminate labeling-importance information throughout the network, facilitating the generation of soft pseudo-labels. The manifold-based LE (Hou et al., 2016) leverages the locally linear embedding technique to derive soft pseudo-labels. Low-rank representation LE method (Tang et al., 2020) captures the global relationships among samples and predicts implicit label correlations for generating soft pseudo-labels. Furthermore, Zhu et al. (Zhu et al., 2020) incorporate both the structural relationships between instances and privileged information to generate soft pseudo-labels.

Soft pseudo-labels are commonly employed to serve as a kind of supervision refinement technique in weakly-supervised learning, such as partial label learning (PLL) (Lv et al., 2020; Zhang & Yu, 2015) where the correct label is hidden in candidate label set. Yao et al. (Yao et al., 2020) generate soft pseudo-labels by assembling the model predictions of different epochs and regard them as the auxiliary

supervision information for the next epoch. Progressive identification-based methods (Lv et al., 2020; Feng et al., 2020) normalize the softmax output on candidate labels to create soft pseudo-labels. Xu et al. (Xu et al., 2021) treat the soft pseudo-label as a latent variable and utilize variational inference to approximate the poster of the soft labels. Wang et al. (Wang et al., 2022) generate soft pseudo-labels by assessing the similarity between contrastive embeddings and each class prototype. Meanwhile, the consistency-based method (Wu et al., 2022) proposes bi-level optimization for generating soft pseudo-label. In the further analysis of our experiments, we will demonstrate that our framework could deal with weakly-supervised data by taking the PLL experiment as an example.

### 3. Proposed Method

#### 3.1. Preliminaries

First of all, we briefly introduce some necessary notations. Let  $\mathcal{X} = \mathbb{R}^q$  be the  $q$ -dimensional instance space and  $\mathcal{Y} = \{1, 2, \dots, c\}$  be the label space with  $c$  class labels. Given the training set  $\mathcal{D} = \{(\mathbf{x}_i, y_i) | 1 \leq i \leq n\}$  where  $\mathbf{x}_i$  denotes the  $q$ -dimensional instance and  $y_i \in \mathcal{Y}$  denotes the ground-truth label associated with  $\mathbf{x}_i$ . In this paper, we aim to find a multi-class classifier  $f : \mathcal{X} \mapsto \Delta^{c-1}$  in the function class  $\mathcal{F}$  with the training set  $\mathcal{D}$ , where  $\Delta^{c-1}$  represents the  $c$ -dimensional simplex. For each training example  $(\mathbf{x}_i, y_i)$ , we use the logical label vector  $\mathbf{l}_i = [l_i^1, l_i^2, \dots, l_i^c]^\top \in \{0, 1\}^c$  to represent whether  $j$  is the ground-truth label, i.e.,  $l_i^j = 1$  if  $j = y_i$ , otherwise  $l_i^j = 0$ . The soft pseudo-label of  $\mathbf{x}_i$  is denoted by  $\mathbf{s}_i = [s_i^1, s_i^2, \dots, s_i^c]^\top \in [0, 1]^c$  where  $\sum_{j=1}^c s_i^j = 1$ . Then  $\mathbf{L} = [\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_n]$  and  $\mathbf{S} = [\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_n]$  represent the initial label matrix and soft pseudo-label matrix, respectively.

#### 3.2. Overview

SEAL alternately trains the predictive model and the soft-pseudo-label generator guided by a meta-network-parameterized objective function. The meta-network dynamically adjusts the parameters of the objective function used for a soft-pseudo-label generator. This optimization is informed by feedback from both the predictive model and the soft-pseudo-label generator, based on their performance in the learning task, creating a feedback loop that refines the training process. Additionally, our framework offers versatility across different learning tasks by allowing direct modifications to the task loss.

SEAL forges a link between the predictive model’s parameters and the meta-network’s parameters by utilizing the soft-pseudo-label generator as a conduit within the meta-learning framework at each iteration. The process is initiated by cloning the parameters from both the predictive model and

soft-pseudo-label generator only for meta-learning purposes. Subsequently, the meta-network optimizes the duplicate parameters of the pseudo-label model, making the updated duplicate parameters depend on the meta-network’s parameters. This is followed by optimizing the predictive model’s parameters using the soft pseudo-labels produced by the soft-pseudo-label generator by with the updated duplicate parameters, thereby making the predictive model’s updated duplicate parameters also depend on the meta-network’s parameters. The task loss then evaluates the output of the predictive model using the updated duplicate parameters, and the gradient is backpropagated towards the meta-network. Upon finishing the meta-learning step, the refined meta-network is further employed to optimize the soft-pseudo-label generator to improve the predictive model.

#### 3.3. The SEAL Framework

To train the predictive DNNs  $f$  parameterized by  $\phi$ , we minimize the following empirical risk estimator  $\widehat{R}(f)$  by leveraging the original label  $\mathbf{l}_i$  and soft pseudo-label  $\mathbf{s}_i$  of each instance  $\mathbf{x}_i$ :

$$\widehat{R}(f) = \frac{1}{n} \sum_{i=1}^n (\ell(f(\mathbf{x}_i; \phi), \mathbf{l}_i) + \lambda \ell(f(\mathbf{x}_i; \phi), \mathbf{s}_i)), \quad (1)$$

where  $\ell$  is a cross-entropy function, and the multiplicative factor  $\lambda$  is used to balance the supervision signal provided by the original label  $\mathbf{l}_i$  and soft pseudo-label  $\mathbf{s}_i$ . The soft pseudo-label  $\mathbf{s}_i$  is generated by a soft-pseudo-label generator  $g$  parameterized by  $\psi$ , i.e.,

$$\mathbf{s}_i = g(\mathbf{x}_i; \psi). \quad (2)$$

where the architecture of  $g$  is the same as  $f$ .

Instead of manually designing an explicit loss for a classifier to output soft pseudo-labels, SEAL proposed a meta-network-parameterized objective function for training the soft-pseudo-label generator  $g$  to align with the target of soft-pseudo-label generation. The parameterized objective function is represented by  $h$  parameterized by a meta-network  $\omega$  (Bechtle et al., 2020). Then the soft-pseudo-label generator  $g$  is optimized by minimizing

$$\mathcal{L}(g) = \frac{1}{n} \sum_{i=1}^n h(g(\mathbf{x}_i; \psi), \mathbf{l}_i; \omega). \quad (3)$$

SEAL alternately trains the predictive model  $f$ , the soft-pseudo-label generator  $g$ , and the parameterized objective  $h$  in a meta-learning process.

Firstly, we begin with copying the parameters of the predictive model and soft-pseudo-label generator with  $\phi' = \phi$  and  $\psi' = \psi$ , which are only used for the meta-learning process. Then, we optimize the parameters  $\psi'$  of the soft-pseudo-label generator  $g$  via using the following parametric

objective function  $\mathcal{L}^m$ :

$$\mathcal{L}^m(g(\mathbf{X}; \psi'), \mathbf{L}) = \frac{1}{n} \sum_{i=1}^n h(g(\mathbf{x}_i; \psi'), \mathbf{l}_i; \omega), \quad (4)$$

After the parameter  $\psi'$  is optimized to the optimal parameter  $\psi'^*$ ,  $\psi'^*$  will always depend on the parameters  $\omega$  of the meta-network  $h$ , which we could explicitly express the dependency as  $\psi'^*(\omega)$ .

Next, let  $\mathbf{S}' = [s'_1, s'_2, \dots, s'_n]$  denote soft pseudo-labels generated by the soft-pseudo-label generator  $g(\cdot; \psi'^*)$  according to Eq. (2). We optimize the parameters  $\phi'$  of the predictive model  $f$  via the loss function  $\mathcal{L}^{\text{pl}}$  as follows:

$$\mathcal{L}^{\text{pl}}(f(\mathbf{X}; \phi'), \mathbf{S}') = \frac{1}{n} \sum_{i=1}^n \ell^{\text{pl}}(f(\mathbf{x}_i; \phi'), s'_i), \quad (5)$$

where  $\ell^{\text{pl}}$  serves as a bridge between the predictive model  $f$  and soft-pseudo-label generator  $g$ . Note that  $\mathcal{L}^{\text{pl}}(f(\mathbf{X}; \phi'), \mathbf{S}')$  also perform back-propagation on  $\mathbf{S}'$ , which allows the optimal parameters  $\phi'^*$  after updating to be dependent on the parameters  $\psi'$  of the soft-pseudo-label generator  $g$ , i.e.,  $\phi'^*(\psi'^*)$ , and further dependent on the parameters  $\omega$  of the meta-network  $h$ , i.e.,  $\phi'^*(\psi'^*(\omega))$ .

Finally, upon finishing the optimization of the parameters  $\phi'$  of the predictive model  $f'$ , we use the following task loss  $\mathcal{L}^t$  to optimize the parameters  $\omega$  of the meta-network  $h$ :

$$\mathcal{L}^t(f(\mathbf{X}; \phi'^*), \mathbf{L}) = \frac{1}{n} \sum_{i=1}^n \ell^t(f(\mathbf{x}_i; \phi'^*), \mathbf{l}_i), \quad (6)$$

where  $\ell^t$  is a specified loss function to some specified task.

Overall, the meta-learning objective can be formulated as the optimization problem as follows:

$$\begin{aligned} & \min_{\omega} \mathcal{L}^t(f'(\mathbf{X}; \phi'^*), \mathbf{L}) \\ \text{s.t. } & \phi'^* = \arg \min_{\phi'} \mathcal{L}^{\text{pl}}(f'(\mathbf{X}; \phi'), \mathbf{S}') \\ & \psi'^* = \arg \min_{\psi'} \mathcal{L}^m(g'(\mathbf{X}; \psi'), \mathbf{L}) \end{aligned} \quad (7)$$

Here, to solve the optimization of Eq. (7), we adopt a batch-style strategy to update  $\psi'$ ,  $\phi'$  and  $\omega$  through a single optimization loop, respectively, to guarantee the efficiency of the algorithm. We employ stochastic gradient descent (SGD) optimization to optimize Eq. (4) (5) and (6). Specifically, in each iteration  $k$  of training, the training set  $\mathcal{D}$  is shuffled into  $I$  mini-batches, containing  $m$  training samples  $\{(\mathbf{x}_i, \mathbf{y}_i) | 1 \leq i \leq m\}$ . Firstly, the updating equation of the soft-pseudo-label generator parameters  $\psi'^{[k-1]}$  can be formulated by moving along the direction of the objective loss in Eq. (4) into the new  $\psi'^{[k]}$  on a mini-batch training

data as follows:

$$\begin{aligned} \psi'^{[k]} &= \psi'^{[k-1]} \\ &- \frac{\alpha}{m} \sum_{i=1}^m \frac{\partial h(g'(\mathbf{x}_i; \psi'^{[k-1]}), \mathbf{l}_i; \omega^{[k-1]})}{\partial \psi'^{[k-1]}}, \end{aligned} \quad (8)$$

where  $\alpha$  is the step size. Secondly, the updating equation for the predictive model parameters  $\phi'^{[k-1]}$  can be expressed by adjusting them in the direction of the objective loss in Eq. (5) on a mini-batch of training data, resulting in the new parameter set  $\phi'^{[k]}$ , as follows:

$$\phi'^{[k]} = \phi'^{[k-1]} - \frac{\beta}{m} \sum_{i=1}^m \frac{\partial \ell^{\text{pl}}(f'(\mathbf{x}_i; \phi'^{[k-1]}), s'_i)}{\partial \phi'^{[k-1]}}, \quad (9)$$

where  $\beta$  is the step size and  $s'_i = g(\mathbf{x}_i; \psi'^{[k]})$ . Thirdly, the iterative update for the predictive model parameters  $\omega^{[k-1]}$  can be expressed by aligning them with the direction of the objective loss in Eq. (6) on a mini-batch of training data, yielding the refined parameter set  $\omega^{[k]}$ , as outlined below:

$$\omega^{[k]} = \omega^{[k-1]} - \frac{\kappa}{m} \sum_{i=1}^m \frac{\partial \ell^t(f'(\mathbf{x}_i; \phi'^{[k-1]}), \mathbf{l}_i)}{\partial \omega^{[k-1]}}, \quad (10)$$

where  $\kappa$  is the step size. Finally, we update the soft-pseudo-label generator parameters  $\psi'^{[k-1]}$  via the updated meta-network  $h$  on a mini-batch training data as follows:

$$\begin{aligned} \psi'^{[k]} &= \psi'^{[k-1]} \\ &- \frac{\alpha}{m} \sum_{i=1}^m \frac{\partial h(g(\mathbf{x}_i; \psi'^{[k-1]}), \mathbf{l}_i; \omega^{[k]})}{\partial \psi'^{[k-1]}}, \end{aligned} \quad (11)$$

where we keep the step size the same as that in Eq. (8).

In this way, as the soft-pseudo-label generator parameters  $\psi$  and meta-network parameters  $\omega$  are updated iteratively, the soft pseudo-labels  $\mathbf{S}$  produced by the soft-pseudo-label generator  $g$  is also refined to contribute to the training of the predictive model step by step. The algorithmic description of SEAL is presented in Algorithm 1.

### 3.4. Practical Implementation

**Gradual increase of  $\lambda$ .** In the risk estimator Eq. (1), we use a balancing factor  $\lambda$  during the whole training procedure, which is suggested not to be fixed. At the beginning of the training phase, the meta-network  $h$  is not prepared to train the soft-pseudo-label generator  $g$  to produce reliable soft pseudo-labels, which will lead to a performance drop of the classifier. Hence, we apply a linear ramp-up function to increase  $\lambda$  from a small weight to the given  $\lambda_e$  during the first  $T'$  epochs:

$$\lambda(t) = \min\left\{\frac{t}{T'} \lambda_e, \lambda_e\right\}. \quad (12)$$

**Algorithm 1** SEAL Algorithm

**Require:** The training set  $\mathcal{D} = \{(\mathbf{x}_i, l_i) | 1 \leq i \leq n\}$ , epoch  $T$ , iteration  $I$ ;

- 1: **for**  $t = 1, \dots, T$  **do**
- 2: Shuffle the training set  $\mathcal{D} = \{(\mathbf{x}_i, l_i) | 1 \leq i \leq n\}$  into  $I$  mini-batches;
- 3: **for**  $k = 1, \dots, I$  **do**
- 4: Copy the parameters of the predictive model and soft-pseudo-label generator with  $\phi^{[k-1]} = \phi^{[k-1]}$  and  $\psi^{[k-1]} = \psi^{[k-1]}$ ;
- 5: Update  $\psi^{[k-1]}$  to  $\psi^{[k]}$  by Eq. (8);
- 6: Update  $\phi^{[k-1]}$  to  $\phi^{[k]}$  by Eq. (9);
- 7: Update  $\omega^{[k-1]}$  to  $\omega^{[k]}$  by Eq. (10);
- 8: Update  $\psi^{[k-1]}$  to  $\psi^{[k]}$  by Eq. (11);
- 9: Obtain the soft pseudo-labels  $s_i$  for each example  $\mathbf{x}_i$  by Eq. (2);
- 10: Update the predictive model parameters  $\phi$  by forward computation and back-propagation with the empirical risk estimator in Eq. (1);
- 11: **end for**
- 12: **end for**

**Ensure:** The predictive model  $f(\cdot; \phi)$ .

This dynamic strategy gradually makes the second term in Eq. (1) with soft pseudo-labels dominate in the whole risk estimator, as the meta-network  $h$  is being trained during the first  $T'$  epochs,

**Choices of  $\ell^{\text{pl}}$ .** We use the bidirectional Kullback-Leibler (KL) divergence loss as  $\ell^{\text{pl}}$  to optimize  $\phi'$  in Eq. (5):

$$\ell^{\text{pl}}(f(\mathbf{x}_i; \phi'), s'_i) = \text{KL}(s'_i || f(\mathbf{x}_i; \phi')) + \text{KL}(f(\mathbf{x}_i; \phi') || s'_i), \quad (13)$$

**Choices of  $\ell^t$ .** For supervised learning, considering the top-k prediction, we use the following loss as the task loss  $\ell^t$  to optimize  $\omega$  in Eq. (6):

$$\ell^t(f(\mathbf{x}_i; \phi'^*), l_i) = -\log \left( \sum_{k=1}^c P(j) \left( \sum_{j=1}^k Q_{j, y_i}(f(\mathbf{x}_i; \phi'^*)) \right) \right), \quad (14)$$

where  $P(j)$  is the given prior of the top- $j$ , and  $Q_{j, y_i}$  is the predicted probability of  $y_i$  being the  $j$ -th best prediction for  $\mathbf{x}_i$ . We employ the library provided by (Petersen et al., 2022) to estimate  $Q_{j, y_i}$ .

SEAL is flexible to learning tasks as the task loss designed for the supervised learning problems can be modified directly. For example, we could utilize the following loss via considering the consistency between the feature and label spaces (Wu et al., 2022) as the task loss  $\ell^t$  in SEAL to deal

with partial label learning (Zhang et al., 2017) :

$$\ell^t(f(\mathbf{x}_i; \phi'^*), l_i) = \sum_{k=1}^K \text{KL}(\mathbf{u}_i || f(\mathbf{x}_i^k; \phi'^*)), \quad (15)$$

where  $\{\mathbf{x}_i^1, \mathbf{x}_i^2, \dots, \mathbf{x}_i^K\}$  denotes a  $K$ -augmentation set for the instance  $\mathbf{x}_i$ , and  $\mathbf{u}_i = [u_{i1}, u_{i2}, \dots, u_{ic}]$  with

$$u_{ij} = \frac{(\prod_{k=1}^K f_j(\mathbf{x}_i^k; \phi'^*))^{\frac{1}{K}}}{\sum_{o=1}^c (\prod_{k=1}^K f_o(\mathbf{x}_i^k; \phi'^*))^{\frac{1}{K}}}. \quad (16)$$

### 3.5. Theoretical Analysis

We will theoretically demonstrate that the soft pseudo-labels could enable the predictive model to achieve a larger sample margin for classes, which results in a more robust and tight generalization of the predictive model when trained using soft pseudo-labels.

Let  $\mathbf{Z} = [z_1^\top, z_2^\top, \dots, z_n^\top] \in \mathbb{R}$  denote the feature matrix extracted by the backbone of the classifier,  $\Theta = [\theta_1^\top, \theta_2^\top, \dots, \theta_c^\top]$  represent the parameters of the last classifier layer, implemented with a linear layer. Let  $P_j(\mathbf{x})$  denote the class-conditional distribution, i.e.,  $P_j(\mathbf{x}) = P(\mathbf{x}|y = j)$ , and  $[j] = \{i|y_i = j\}$  denote the sample indices corresponding to class  $j$ . We define the margin for a class  $j$  as follows:

$$\gamma_j = \min_{i \in [j]} \theta_j^\top z_i - \max_{j' \neq j} \theta_{j'}^\top z_i. \quad (17)$$

Let  $G(\mathcal{F}, \mathcal{X}, \mathcal{Y})$  denote the generalization error bound induced from the class margin  $\gamma$  (Cao et al., 2019). Then under the assumption about the extracted features in Appendix A, we could obtain the following theorem:

**Theorem 3.1.** *Let  $\Theta$  be the parameters of the last model layer trained by Gradient Descend starting from random initialization. Suppose that at epoch  $T$ ,  $\forall k \in \mathcal{Y}$  with  $k \neq p$ ,  $\theta_k$  arrives at the optimal prototype  $\mathbf{v}_k$ . Then, we fix the extracted features  $\{z_i | 1 \leq i \leq n\}$  and the prototype except  $\theta_p$  to continue the training process for  $\theta_p$ . Let  $G^s(\mathcal{F}, \mathcal{X}, \mathcal{Y}), G^h(\mathcal{F}, \mathcal{X}, \mathcal{Y})$  denote the generalization error bound based on the class margin derived from the empirical risk estimators with soft pseudo-labels and initial hard labels, respectively. At epoch  $T' > T$ , we could have*

$$G^s(\mathcal{F}, \mathcal{X}, \mathcal{Y}) < G^h(\mathcal{F}, \mathcal{X}, \mathcal{Y}).$$

The proof of Theorem 3.1 is provided in Appendix A. Theorem 3.1 shows that the expected generalization error of the model trained on proper soft pseudo-labels could be bounded by a smaller upper bound than the model trained on initial hard labels.

Table 1. Classification accuracy of each comparing approach (mean  $\pm$  std). The best performance is shown in boldface.  $\bullet$ / $\circ$  indicates whether the performance of SEAL is statistically superior/inferior to the comparing algorithm on each dataset (pairwise t-test at 0.05 significance level).

| Methods               | CIFAR-10                            | CIFAR-100                           | TinyImageNet                        |
|-----------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| CE                    | 91.36 $\pm$ 0.33% $\bullet$         | 68.11 $\pm$ 0.45% $\bullet$         | 54.28 $\pm$ 0.27% $\bullet$         |
| LSR                   | 91.62 $\pm$ 0.28%                   | 68.96 $\pm$ 0.60% $\bullet$         | 54.74 $\pm$ 0.28% $\bullet$         |
| TF-KD <sub>self</sub> | 91.66 $\pm$ 0.04%                   | 70.44 $\pm$ 0.26%                   | 54.34 $\pm$ 0.34% $\bullet$         |
| TF-KD <sub>reg</sub>  | 91.14 $\pm$ 0.35% $\bullet$         | 69.47 $\pm$ 0.56% $\bullet$         | 54.11 $\pm$ 0.37% $\bullet$         |
| CS-KD                 | 91.38 $\pm$ 0.14% $\bullet$         | 68.63 $\pm$ 0.17% $\bullet$         | 55.67 $\pm$ 0.58% $\bullet$         |
| PS-KD                 | 91.23 $\pm$ 0.22% $\bullet$         | 68.39 $\pm$ 0.81% $\bullet$         | 54.35 $\pm$ 0.75% $\bullet$         |
| DLB                   | 91.29 $\pm$ 0.23% $\bullet$         | 68.48 $\pm$ 0.40% $\bullet$         | 53.61 $\pm$ 0.19% $\bullet$         |
| USKD                  | 91.47 $\pm$ 0.15% $\bullet$         | 69.33 $\pm$ 0.66% $\bullet$         | 55.33 $\pm$ 0.32% $\bullet$         |
| SEAL                  | <b>91.79 <math>\pm</math> 0.11%</b> | <b>70.72 <math>\pm</math> 0.17%</b> | <b>56.76 <math>\pm</math> 0.22%</b> |

## 4. Experiments

### 4.1. Datasets

We employ three benchmark datasets for multi-class classification including CIFAR-10, CIFAR-100 (Krizhevsky et al., 2009), and TinyImageNet (Le & Yang, 2015), to evaluate the proposed approach. Both CIFAR-10 and CIFAR-100 consist of a total of 60,000 images with resolutions of  $32 \times 32$  pixels, distributed across 10 classes and 100 classes, respectively. The TinyImageNet dataset is a subset of ILSVRC-2012 (Russakovsky et al., 2015) and contains 200 classes. Each class is represented by 500 training samples and 50 testing samples, all resized to  $64 \times 64$ . In our preprocessing routine, training images from all datasets underwent random cropping and resizing to a uniform  $32 \times 32$  pixel format post-normalization, whereas test images were solely normalized. We allocated 10% of the training data from each dataset for validation purposes.

### 4.2. Baselines

We compare the performance of SEAL with seven soft-pseudo-label approaches:

- LSR (Szegedy et al., 2016): A label-smoothing regularization approach which replaces the one-hot encoded true label with a smoothed distribution that assigns some probabilities to incorrect labels.
- TF-KD<sub>self</sub> (Yuan et al., 2020): A self-training knowledge distillation approach which trains the student model in a normal way to obtain a pre-trained model and uses it to generate soft labels to train itself.
- TF-KD<sub>reg</sub> (Yuan et al., 2020): A teacher-free knowledge distillation approach which uses manually designed soft labels that the probability of a correct class is much higher than that of an incorrect one.
- CS-KD (Yun et al., 2020): A class-wise self-distillation approach which matches or distills the predictive distribu-

tion of the model between different samples of the same label.

- PS-KD (Kim et al., 2021): A progressive self-distillation approach which adjusts the training targets progressively by combining the ground-truth and past predictions from the model itself.
- DLB (Shen et al., 2022): An efficient self-distillation approach which distills the on-the-fly generated smooth labels in the previous iteration after rearranging the sampling sequence.
- USKD (Yang et al., 2023): A self-distillation approach generates customized soft labels for both target and non-target classes without a teacher.

In addition, we also compare the performance of SEAL against CE, which directly uses initial binary labels to train the model with cross-entropy loss.

We employ ResNet-20, ResNet-44 and ResNet-18 for CIFAR-10, CIFAR-100 and TinyImageNet as backbone respectively. For the meta-network, we follow (Bechtle et al., 2020; Raymond et al., 2023) and use a small feedforward neural network with two hidden layers and 40 hidden units. Smooth leaky ReLU activations are used in the hidden layers and smooth Softplus activation is used in the output layer. We configure the total number of epochs as 200 and set the batch size to 128. We employ the SGD optimizer with a momentum of 0.9 and a weight decay of  $1 \times 10^{-4}$ , where the initial learning rate is established at 0.1 with a decay factor of 10%. Additionally, we incorporate standard data augmentation techniques, including Random Horizontal Flipping and Random Cropping.

### 4.3. Experimental Results

Table 1 shows the classification accuracy of each comparative approach for different benchmark datasets. We perform 5 trials with different random seeds and the reported metrics

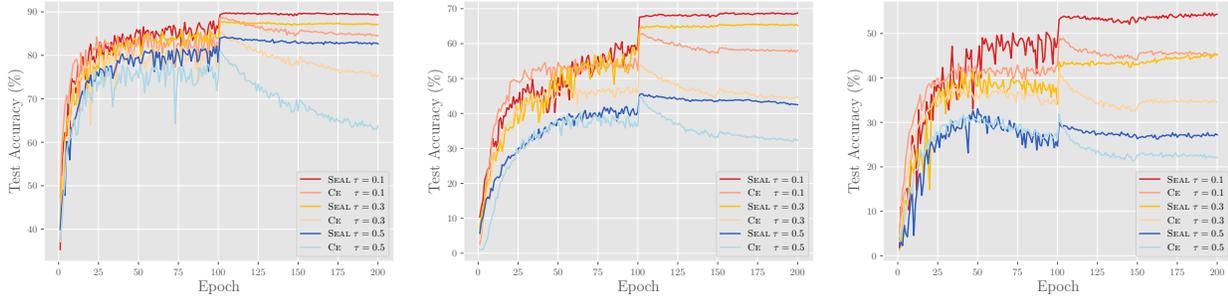


Figure 1. The test accuracy of SEAL and CE on corrupted benchmark datasets with noisy rate  $\tau \in \{0.1, 0.3, 0.5\}$

include both the mean and standard deviation. The best performance is shown in boldface. ●/○ indicates whether the performance of SEAL is statistically superior/inferior to the comparing algorithm on each dataset (pairwise t-test at 0.05 significance level). As shown in table 1, it is impressive to observe that:

- SEAL achieves the best performance against all the comparing approaches on all benchmark datasets.
- SEAL achieves superior or at least comparable performance to other approaches on CIFAR-10 and CIFAR-100.
- SEAL achieves superior to other approaches on TinyImageNet and exceeds the performance of the second-best algorithm by 1.09%.
- With the growing complexity of datasets, our approach demonstrates progressively superior performance.

Additionally, to prove the robustness to label noise of our method, we manually corrupt the benchmark datasets with symmetric noisy labels under the noisy ratio  $\tau \in \{0.1, 0.3, 0.5\}$ . Subsequently, for each selected training instance, we replace its correct label with another possible label to create a noisy label. Figure 1 illustrates the curves of test accuracy comparing the baseline CE with SEAL under different noisy rate  $\tau$  on the benchmark datasets. It can be observed from the figure that SEAL achieves higher test accuracy against CE in all scenarios and the performance of SEAL is robust in the last half of the training stage. These observations demonstrate that SEAL is robust to label noise.

#### 4.4. Further Analysis

##### 4.4.1. ABLATION AND CONVERGENCE STUDY

To show the helpfulness of the meta-network-parameterized objective function in SEAL, the vanilla variant of SEAL, i.e., SEAL-NM, is adopted to conduct the ablation study. For SEAL-NM, the meta-network-parameterized objective

Table 2. Classification accuracy (mean  $\pm$  std) of SEAL and its variant on the benchmark datasets. The best performance is shown in boldface. ●/○ indicates whether the performance of SEAL is statistically superior/inferior to the comparing algorithm on each dataset (pairwise t-test at 0.05 significance level).

| Datasets     | SEAL                                | SEAL-NM            |
|--------------|-------------------------------------|--------------------|
| CIFAR-10     | <b>91.79 <math>\pm</math> 0.11%</b> | 91.41 $\pm$ 0.33%● |
| CIFAR-100    | <b>70.72 <math>\pm</math> 0.17%</b> | 68.63 $\pm$ 0.42%● |
| TinyImageNet | <b>56.76 <math>\pm</math> 0.22%</b> | 55.31 $\pm$ 0.64%● |

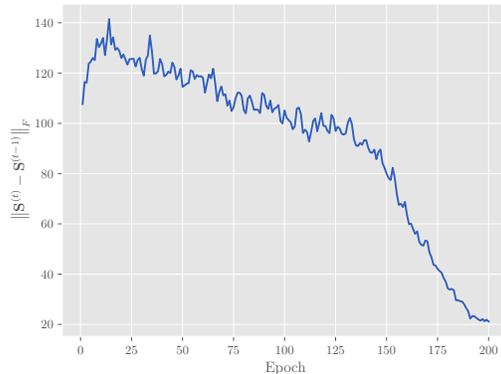


Figure 2. Convergence of the generated soft pseudo-label matrix  $S$  on CIFAR-100.

function is replaced by the cross-entropy loss. Table 2 shows that SEAL achieves superior performance to SEAL-NM on all datasets, which demonstrates the usefulness of the meta-network-parameterized objective function in SEAL for improving the performance of the classifier.

Besides, Figure 2 illustrates the soft pseudo-labels generated by soft-pseudo-label generator in SEAL converges with the number of epochs on CIFAR-100, which shows that the soft pseudo-labels in SEAL could converge efficiently.

Table 3. Classification accuracy (mean<sub>std</sub>) of each comparing algorithm in terms of the different proportion of false positive candidate labels. The best performance (the larger the better) is shown in boldface. ●/○ indicates whether the performance of SEAL is statistically superior/inferior to the comparing algorithm on each dataset (pairwise t-test at 0.05 significance level).

| Methods | CIFAR-10                     |                              |                              |                              | CIFAR-100                    |                              |                              |
|---------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|
|         | 0.3                          | 0.5                          | 0.7                          | 0.9                          | 0.03                         | 0.05                         | 0.1                          |
| PRODEN  | 91.27 <sub>0.33</sub> ●      | 90.42 <sub>0.14</sub> ●      | 89.26 <sub>0.33</sub> ●      | 84.02 <sub>1.16</sub> ●      | 65.03 <sub>0.69</sub> ●      | 64.28 <sub>0.50</sub> ●      | 64.04 <sub>0.24</sub> ●      |
| CC      | 89.75 <sub>0.50</sub> ●      | 87.57 <sub>0.35</sub> ●      | 84.88 <sub>0.25</sub> ●      | 79.96 <sub>0.54</sub> ●      | 64.75 <sub>0.37</sub> ●      | 63.84 <sub>0.89</sub> ●      | 61.66 <sub>0.24</sub> ●      |
| VALEN   | 89.19 <sub>0.49</sub> ●      | 88.36 <sub>0.30</sub> ●      | 87.29 <sub>0.43</sub> ●      | 44.75 <sub>1.96</sub> ●      | 66.77 <sub>0.99</sub> ●      | 65.97 <sub>0.87</sub> ●      | 65.27 <sub>0.18</sub> ●      |
| CAVL    | 89.23 <sub>3.64</sub> ●      | 88.26 <sub>4.06</sub> ●      | 53.38 <sub>2.81</sub> ●      | 13.19 <sub>2.33</sub> ●      | 57.75 <sub>2.35</sub> ●      | 47.07 <sub>3.19</sub> ●      | 25.83 <sub>1.71</sub> ●      |
| IDGP    | 92.07 <sub>0.32</sub> ●      | 91.04 <sub>0.13</sub> ●      | 88.37 <sub>2.50</sub> ●      | 47.76 <sub>1.51</sub> ●      | 68.19 <sub>0.02</sub> ●      | 67.68 <sub>0.29</sub> ●      | 62.39 <sub>1.74</sub> ●      |
| PICO    | 88.14 <sub>0.25</sub> ●      | 86.00 <sub>0.24</sub> ●      | 76.47 <sub>1.65</sub> ●      | 25.63 <sub>6.25</sub> ●      | 60.50 <sub>0.39</sub> ●      | 58.69 <sub>0.52</sub> ●      | 37.91 <sub>1.56</sub> ●      |
| PLCR    | 90.90 <sub>0.14</sub> ●      | 90.39 <sub>0.21</sub> ●      | 89.31 <sub>0.39</sub> ●      | 81.39 <sub>2.79</sub> ●      | 67.01 <sub>0.33</sub> ●      | 65.78 <sub>0.64</sub> ●      | 62.46 <sub>0.37</sub> ●      |
| SEAL    | <b>92.93</b> <sub>0.29</sub> | <b>91.98</b> <sub>0.28</sub> | <b>89.88</b> <sub>0.34</sub> | <b>85.36</b> <sub>0.41</sub> | <b>71.27</b> <sub>0.40</sub> | <b>68.75</b> <sub>0.67</sub> | <b>67.21</b> <sub>1.31</sub> |

#### 4.4.2. LEARNING WITH PARTIAL LABELS

In this subsection, we utilize SEAL to deal with partial label learning (PLL) (Zhang et al., 2016) by modifying the task loss as Eq. (15). Adopting the same experimental settings in previous PLL work (Lv et al., 2020; Wang et al., 2022; Feng et al., 2020), we manually corrupt CIFAR-10 and CIFAR-100 into partially labeled versions by flipping negative labels to false positive labels with the probability {0.3, 0.5, 0.7, 0.9} and {0.03, 0.05, 0.1}, respectively. The performance of SEAL is compared against seven deep PLL methods:

- PRODEN (Lv et al., 2020): A progressive identification approach which approximately minimizes a risk estimator and identifies the true labels in a seamless manner;
- CC (Feng et al., 2020): A classifier-consistent approach which also uses the loss correction strategy to learn the classifier that approaches the optimal one;
- VALEN (Xu et al., 2021): An instance-dependent PLL approach which recovers the latent label distribution via variational inference methods;
- CAVL (Zhang et al., 2022): A progressive identification approach which exploits the class activation value to identify the true label in candidate label sets.
- IDGP (Qiao et al., 2023): A disambiguation approach which builds the model upon a decompositional generation process of candidate labels.
- PICO (Wang et al., 2022): A data-augmentation-based method which identifies the true label via contrastive learning with learned prototypes for image datasets.
- PLCR (Wu et al., 2022): A data-augmentation-based method which identifies the true label via consistency regularization with random augmented instances.

Here, we employ a ResNet-32 as the backbone. The total number of epochs is set to 200. We use SGD optimizer with a momentum of 0.9 and weight decay of

$1 \times 10^{-3}$ . Learning rates are chosen from the orders of magnitude  $\{10^{-2}, 10^{-3}, 10^{-4}\}$  guided by the performance on the validation dataset. Common data augmentations are applied, including Random Horizontal Flipping, Random Cropping, Cutout (Devries & Taylor, 2017), and Auto Augment (Cubuk et al., 2019).

Table 3 illustrates the classification accuracy comparing the PLL baselines with SEAL under different flipping probability of negative labels to false positive labels on CIFAR-10 and CIFAR-100. By changing the task loss in our framework SEAL, we successfully adapt to the PLL task and demonstrate significantly superior performances against all the PLL baselines on all settings.

## 5. Conclusion

In this paper, we propose a novel framework SEAL that trains the predictive model using soft pseudo-labels generated by the specialized soft-pseudo-label generator with a meta-network-parameterized objective. SEAL alternately trains the predictive model and the soft-pseudo-label generator guided by a meta-network-parameterized objective function. The meta-network dynamically adjusts the parameters of the objective function used for a soft-pseudo-label generator. This optimization is informed by feedback from both the predictive model and the soft-pseudo-label generator, based on their performance in the learning task. SEAL is flexible to learning tasks as the task loss designed for the supervised learning problems can be modified directly. Experiments validate the effectiveness of the proposed method.

## Acknowledgments

This research was supported by the National Science Foundation of China (62206050, 62125602, and 62076063), China Postdoctoral Science Foundation (2021M700023), Jiangsu Province Science Foundation for

Youths (BK20210220), Young Elite Scientists Sponsorship Program of Jiangsu Association for Science and Technology (TJ-2022-078), and the Big Data Computing Center of Southeast University.

## Impact Statement

The potential broader impact of our work might be that the need for accurately annotated data would be significantly reduced. As a result, the rate of unemployment for data annotation specialists might be increased.

## References

- Algan, G. and Ulusoy, I. Metalabelnet: Learning to generate soft-labels from noisy-labels. *IEEE Transactions on Image Processing*, 31:4352–4362, 2022.
- Bechtle, S., Molchanov, A., Chebotar, Y., Grefenstette, E., Righetti, L., Sukhatme, G. S., and Meier, F. Meta learning via learned loss. In *25th International Conference on Pattern Recognition, ICPR 2020, Virtual Event / Milan, Italy*, pp. 4161–4168, 2020.
- Cao, K., Wei, C., Gaidon, A., Aréchiga, N., and Ma, T. Learning imbalanced datasets with label-distribution-aware margin loss. In *Advances in Neural Information Processing Systems, NeurIPS 2019, Vancouver, BC, Canada*, volume 32, pp. 1565–1576, 2019.
- Cubuk, E. D., Zoph, B., Mané, D., Vasudevan, V., and Le, Q. V. Autoaugment: Learning augmentation strategies from data. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA*, pp. 113–123, 2019.
- Devries, T. and Taylor, G. W. Improved regularization of convolutional neural networks with cutout. *CoRR*, abs/1708.04552, 2017.
- Feng, L., Lv, J., Han, B., Xu, M., Niu, G., Geng, X., An, B., and Sugiyama, M. Provably consistent partial-label learning. In *Advances in Neural Information Processing Systems, NeurIPS 2020, virtual*, volume 33, pp. 10948–10960, 2020.
- Ghoshal, A., Chen, X., Gupta, S., Zettlemoyer, L., and Mehdad, Y. Learning better structured representations using low-rank adaptive label smoothing. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, 2021*.
- Hinton, G. E., Vinyals, O., and Dean, J. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531, 2015.
- Hou, P., Geng, X., and Zhang, M. Multi-label manifold learning. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, Phoenix, Arizona, USA*, pp. 1680–1686, 2016.
- Kim, K., Ji, B., Yoon, D., and Hwang, S. Self-knowledge distillation with progressive refinement of targets. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada*, pp. 6547–6556, 2021.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Le, Y. and Yang, X. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- Li, Y. and Liang, Y. Learning overparameterized neural networks via stochastic gradient descent on structured data. In Bengio, S., Wallach, H. M., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems, NeurIPS 2018, Montréal, Canada*, volume 31, pp. 8168–8177, 2018.
- Li, Y., Zhang, M., and Geng, X. Leveraging implicit relative labeling-importance information for effective multi-label learning. In *2015 IEEE International Conference on Data Mining, ICDM 2015, Atlantic City, NJ, USA*, pp. 251–260, 2015.
- Liang, J., Li, L., Bing, Z., Zhao, B., Tang, Y., Lin, B., and Fan, H. Efficient one pass self-distillation with zipf’s label smoothing. In *17th European Conference on Computer Vision, ECCV 2022, Tel Aviv, Israel*, volume 13671, pp. 104–119, 2022.
- Lukov, T., Zhao, N., Lee, G. H., and Lim, S. Teaching with soft label smoothing for mitigating noisy labels in facial expressions. In *17th European Conference on Computer Vision, ECCV 2022, Tel Aviv, Israel*, volume 13672, pp. 648–665, 2022.
- Lv, J., Xu, M., Feng, L., Niu, G., Geng, X., and Sugiyama, M. Progressive identification of true labels for partial-label learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, Virtual Event*, volume 119, pp. 6500–6510, 2020.
- Ngo, N. T., Van, L. N., and Nguyen, T. H. Unsupervised domain adaptation for text classification via meta self-paced learning. In *Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea*, pp. 4741–4752, 2022.
- Petersen, F., Kuehne, H., Borgelt, C., and Deussen, O. Differentiable top-k classification learning. In *International Conference on Machine Learning, ICML 2022, Baltimore, Maryland, USA*, volume 162, pp. 17656–17668, 2022.

- Qiao, C., Xu, N., and Geng, X. Decompositional generation process for instance-dependent partial label learning. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, 2023*.
- Raymond, C., Chen, Q., Xue, B., and Zhang, M. Online loss function learning. *CoRR*, abs/2301.13247, 2023.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M. S., Berg, A. C., and Fei-Fei, L. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- Shen, Y., Xu, L., Yang, Y., Li, Y., and Guo, Y. Self-distillation from the last mini-batch for consistency regularization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA*, pp. 11933–11942, 2022.
- Sun, S., Cheng, Y., Gan, Z., and Liu, J. Patient knowledge distillation for BERT model compression. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China*, pp. 4322–4331, 2019.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. Rethinking the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA*, pp. 2818–2826, 2016.
- Tang, H., Zhu, J., Zheng, Q., Wang, J., Pang, S., and Li, Z. Label enhancement with sample correlations via low-rank representation. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA*, pp. 5932–5939, 2020.
- Wang, H., Xiao, R., Li, Y., Feng, L., Niu, G., Chen, G., and Zhao, J. Pico: Contrastive label disambiguation for partial label learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, 2022*.
- Wang, M. and Ma, C. Early stage convergence and global convergence of training mildly parameterized neural networks. In *Advances in Neural Information Processing Systems, NeurIPS 2022, New Orleans, LA, USA*, volume 35, pp. 743–756, 2022.
- Wu, D., Wang, D., and Zhang, M. Revisiting consistency regularization for deep partial label learning. In *International Conference on Machine Learning, ICML 2022, Baltimore, Maryland, USA*, volume 162, pp. 24212–24225, 2022.
- Xu, N., Liu, Y.-P., and Geng, X. Label enhancement for label distribution learning. *IEEE Transactions on Knowledge and Data Engineering*, 33(4):1632–1643, 2019.
- Xu, N., Shu, J., Liu, Y., and Geng, X. Variational label enhancement. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, Virtual Event*, volume 119, pp. 10597–10606, 2020.
- Xu, N., Qiao, C., Geng, X., and Zhang, M. Instance-dependent partial label learning. In *Advances in Neural Information Processing Systems, NeurIPS 2021, virtual*, volume 34, pp. 27119–27130, 2021.
- Xu, N., Shu, J., Zheng, R., Geng, X., Meng, D., and Zhang, M. Variational label enhancement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):6537–6551, 2023.
- Xu, T. and Liu, C. Data-distortion guided self-distillation for deep neural networks. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, Honolulu, Hawaii, USA*, volume 33, pp. 5565–5572, 2019.
- Yang, Z., Zeng, A., Li, Z., Zhang, T., Yuan, C., and Li, Y. From knowledge distillation to self-knowledge distillation: A unified approach with normalized loss and customized soft labels. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 17185–17194, 2023.
- Yao, Y., Deng, J., Chen, X., Gong, C., Wu, J., and Yang, J. Deep discriminative CNN with temporal ensembling for ambiguously-labeled image classification. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA*, pp. 12669–12676, 2020.
- Yi, K. and Wu, J. Probabilistic end-to-end noise correction for learning with noisy labels. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA*, pp. 7017–7025, 2019.
- Yuan, L., Tay, F. E. H., Li, G., Wang, T., and Feng, J. Revisiting knowledge distillation via label smoothing regularization. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA*, pp. 3902–3910, 2020.
- Yun, S., Park, J., Lee, K., and Shin, J. Regularizing class-wise predictions via self-knowledge distillation. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA*, pp. 13873–13882, 2020.
- Zhang, C., Jiang, P., Hou, Q., Wei, Y., Han, Q., Li, Z., and Cheng, M. Delving deep into label smoothing. *IEEE Transactions on Image Processing*, 30:5984–5996, 2021a.

- Zhang, F., Feng, L., Han, B., Liu, T., Niu, G., Qin, T., and Sugiyama, M. Exploiting class activation value for partial-label learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, 2022*.
- Zhang, L., Song, J., Gao, A., Chen, J., Bao, C., and Ma, K. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South)*, pp. 3712–3721, 2019.
- Zhang, M. and Yu, F. Solving the partial label learning problem: An instance-based approach. In *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina*, pp. 4048–4054, 2015.
- Zhang, M., Zhou, B., and Liu, X. Partial label learning via feature-aware disambiguation. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA*, pp. 1335–1344, 2016.
- Zhang, M., Zhang, Q., Fang, J., Li, Y., and Geng, X. Leveraging implicit relative labeling-importance information for effective multi-label learning. *IEEE Transactions on Knowledge and Data Engineering*, 33(5):2057–2070, 2021b.
- Zhang, M.-L. and Zhou, Z.-H. A review on multi-label learning algorithms. *IEEE Transactions on Knowledge and Data Engineering*, 26(8):1819–1837, 2014.
- Zhang, M.-L., Yu, F., and Tang, C.-Z. Disambiguation-free partial label learning. *IEEE Transactions on Knowledge and Data Engineering*, 29(10):2155–2167, 2017.
- Zhang, Y., Xiang, T., Hospedales, T. M., and Lu, H. Deep mutual learning. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA*, pp. 4320–4328, 2018.
- Zhou, X., Liu, X., Zhai, D., Jiang, J., Gao, X., and Ji, X. Prototype-anchored learning for learning with imperfect annotations. In *International Conference on Machine Learning, ICML 2022, Baltimore, Maryland, USA*, volume 162, pp. 27245–27267, 2022.
- Zhu, W., Jia, X., and Li, W. Privileged label enhancement with multi-label learning. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020, Yokohama, Japan*, pp. 2376–2382, 2020.

## A. Appendix

We start with introducing some necessary notations. Let  $\mathbf{Z} = [\mathbf{z}_1^\top, \mathbf{z}_2^\top, \dots, \mathbf{z}_n^\top] \in \mathbb{R}^{2 \times n}$  denote the feature matrix extracted by the backbone of the classifier,  $\Theta(t) = [\theta_1^\top(t), \theta_2^\top(t), \dots, \theta_c^\top(t)] \in \mathbb{R}^{a \times c}$  represent the parameters at the  $t$ -th epoch of the last classifier layer, implemented with a linear layer. Following (Zhou et al., 2022), we refer to  $\theta_j$  as the prototype of the class  $j$ . Let  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n]^\top \in \mathbb{R}^{n \times c}$  denote the label matrix for supervision, i.e.,  $\mathbf{W} = \mathbf{L}$  if we use hard labels, otherwise  $\mathbf{W} = \mathbf{S}$ . For multi-class learning, we use softmax operation and cross-entropy loss function. Hence, the empirical risk estimator at the  $t$ -th epoch is as follows:

$$\mathcal{L}(t) = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c w_i^j \log \frac{e^{\theta_j^\top(t) \mathbf{z}_i}}{\sum_{k=1}^c e^{\theta_k^\top(t) \mathbf{z}_i}}. \quad (18)$$

Let  $P_j(\mathbf{x})$  denote the class-conditional distribution, i.e.,  $P_j(\mathbf{x}) = P(\mathbf{x}|y = j)$ , and  $[j] = \{i|y_i = j\}$  denote the sample indices corresponding to class  $j$ . We define the margin for a class  $j$  as follows:

$$\gamma_j = \min_{i \in [j]} \theta_j^\top \mathbf{z}_i - \max_{j' \neq j} \theta_{j'}^\top \mathbf{z}_i. \quad (19)$$

Let  $\mathcal{L}_{\gamma_j, j}$  denote the margin loss on examples from class  $j$ :

$$\mathcal{L}_{\gamma_j, j}[f] = \mathbb{E}_{\mathbf{x} \sim P_j(\mathbf{x})} \mathbb{I}[\max_{j' \neq j} f_{j'}(\mathbf{x}) > f_j(\mathbf{x}) - \gamma_j], \quad (20)$$

and  $\hat{\mathcal{L}}_{\gamma_j, j}$  denote its empirical variant:

$$\hat{\mathcal{L}}_{\gamma_j, j}[f] = \frac{1}{|[j]|} \sum_{i \in [j]} \mathbb{I}[\max_{j' \neq j} f_{j'}(\mathbf{x}_i) > f_j(\mathbf{x}_i) - \gamma_j]. \quad (21)$$

Here, for convenience, we let  $f_j(\mathbf{x}_i) = \theta_j^\top \mathbf{z}_i$  denote the logits of the model without the softmax operation. For a hypothesis class  $\mathcal{F}$ , let  $\hat{\mathfrak{R}}_j(\mathcal{F})$  denote the empirical Rademacher complexity of its class  $j$  margin:

$$\hat{\mathfrak{R}}_j(\mathcal{F}) = \frac{1}{|[j]|} \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \sum_{i \in [j]} \sigma_i [f_j(\mathbf{x}_i) - \max_{j' \neq j} f_{j'}(\mathbf{x}_i)] \right] \quad (22)$$

where  $\sigma$  is a vector of i.i.d. uniform from  $\{+1, -1\}$  bits. According to (Cao et al., 2019), we have the following theorem of class-balanced generalization error bound:

**Theorem 1.** *With probability  $1 - \delta$  over the randomness of the training data, for class sample margins  $\gamma_1, \dots, \gamma_c > 0$ , let*

$$G(\mathcal{F}, \mathcal{X}, \mathcal{Y}) = \frac{1}{c} \sum_{j=1}^c \left( \hat{\mathcal{L}}_{\gamma_j, j}[f] + \frac{4}{\gamma_j} \hat{\mathfrak{R}}_j(\mathcal{F}) + \epsilon_j(\gamma_j) \right), \quad (23)$$

where  $\epsilon_j(\gamma_j) = \sqrt{\frac{\log \log_2 \left( \frac{2 \max_{\mathbf{x} \in \mathcal{X}, f \in \mathcal{F}} |f(\mathbf{x})|}{\gamma_j} \right) + \log \frac{2\epsilon}{\delta}}{|[j]|}}$  is typically a low-order term in  $|[j]|$ . For all hypotheses  $f \in \mathcal{F}$ , we will have balanced-class generalization bounded by:

$$P_{\mathbf{x}, y}[f_y(\mathbf{x}) < \max_{j \neq y} f_j(\mathbf{x})] \leq G(\mathcal{F}, \mathcal{X}, \mathcal{Y}, \{\gamma_1, \dots, \gamma_c\}). \quad (24)$$

**Lemma A.1.** *Suppose that the  $\gamma_p$  in  $\{\gamma_1, \dots, \gamma_c\}$  is replace by  $\gamma'_p$  to obtain  $G'(\mathcal{F}, \mathcal{X}, \mathcal{Y})$ , which satisfies  $\gamma'_p > \gamma_p$ , then we have  $G'(\mathcal{F}, \mathcal{X}, \mathcal{Y}) < G(\mathcal{F}, \mathcal{X}, \mathcal{Y})$ .*

*Proof.* According to the definition of the class margin in Eq. (19), for the first term of Eq. (23) we have

$$\begin{aligned}
 \hat{\mathcal{L}}_{\gamma_j, j}[f] &= \frac{1}{|[j]|} \sum_{i \in [j]} \mathbb{I}[\max_{j' \neq j} f_{j'}(\mathbf{x}_i) > f_j(\mathbf{x}_i) - \gamma_j] \\
 &= \frac{1}{|[j]|} \sum_{i \in [j]} \mathbb{I}[f_j(\mathbf{x}_i) - \max_{j' \neq j} f_{j'}(\mathbf{x}_i) < \min_{i \in [j]} \boldsymbol{\theta}_j^\top \mathbf{z}_i - \max_{j' \neq j} \boldsymbol{\theta}_{j'}^\top \mathbf{z}_i] \\
 &= \frac{1}{|[j]|} \sum_{i \in [j]} \mathbb{I}[\frac{1}{|[j]|} \sum_{i \in [j]} \mathbb{I}[\min_{i \in [j]} \boldsymbol{\theta}_j^\top \mathbf{z}_i - \max_{j' \neq j} \boldsymbol{\theta}_{j'}^\top \mathbf{z}_i < \min_{i \in [j]} \boldsymbol{\theta}_j^\top \mathbf{z}_i - \max_{j' \neq j} \boldsymbol{\theta}_{j'}^\top \mathbf{z}_i]] \\
 &= 0.
 \end{aligned} \tag{25}$$

For the second term, the Rademacher complexity  $\hat{\mathfrak{R}}_j(\mathcal{F})$  will typically scale as  $\sqrt{\frac{C(\mathcal{F})}{|[j]|}}$ , which  $C(\mathcal{F})$  denotes some complexity measure of  $\mathcal{F}$ , only related to  $\mathcal{F}$  and the number of the instances belonging to class  $j$ . Hence,  $\frac{4}{\gamma_j} \hat{\mathfrak{R}}_j(\mathcal{F})$  will increase as  $\gamma_j$  decreases. So does the third term  $\epsilon_j(\gamma_j)$ .

Then, we have

$$\begin{aligned}
 G(\mathcal{F}, \mathcal{X}, \mathcal{Y}) &= \frac{1}{c} \left( \sum_{j \neq p} \hat{\mathcal{L}}_{\gamma_j, j}[f] + \frac{4}{\gamma_j} \hat{\mathfrak{R}}_j(\mathcal{F}) + \epsilon_j(\gamma_j) + (\hat{\mathcal{L}}_{\gamma_p, p}[f] + \frac{4}{\gamma_p} \hat{\mathfrak{R}}_p(\mathcal{F}) + \epsilon_p(\gamma_p)) \right) \\
 &= \frac{1}{c} \left( \sum_{j \neq p} \hat{\mathcal{L}}_{\gamma_j, j}[f] + \frac{4}{\gamma_j} \hat{\mathfrak{R}}_j(\mathcal{F}) + \epsilon_j(\gamma_j) + (\frac{4}{\gamma_p} \hat{\mathfrak{R}}_p(\mathcal{F}) + \epsilon_p(\gamma_p)) \right) \\
 &> \frac{1}{c} \left( \sum_{j \neq p} \hat{\mathcal{L}}_{\gamma_j, j}[f] + \frac{4}{\gamma_j} \hat{\mathfrak{R}}_j(\mathcal{F}) + \epsilon_j(\gamma_j) + (\frac{4}{\gamma'_p} \hat{\mathfrak{R}}_p(\mathcal{F}) + \epsilon_p(\gamma'_p)) \right) \\
 &= \frac{1}{c} \left( \sum_{j \neq p} \hat{\mathcal{L}}_{\gamma_j, j}[f] + \frac{4}{\gamma_j} \hat{\mathfrak{R}}_j(\mathcal{F}) + \epsilon_j(\gamma_j) + (\hat{\mathcal{L}}_{\gamma'_p, p}[f] + \frac{4}{\gamma'_p} \hat{\mathfrak{R}}_p(\mathcal{F}) + \epsilon_p(\gamma'_p)) \right) \\
 &= G'(\mathcal{F}, \mathcal{X}, \mathcal{Y})
 \end{aligned} \tag{26}$$

Hence, we obtain  $G'(\mathcal{F}, \mathcal{X}, \mathcal{Y}) < G(\mathcal{F}, \mathcal{X}, \mathcal{Y})$ , and complete the proof.

We make an assumption on the distribution of the extracted features as follows:

**Assumption 1.** For  $p \in \mathcal{Y}$ , let  $\mathbf{v}_p = \frac{\frac{1}{|[p]|} \sum_{i \in [p]} \mathbf{z}_i}{\|\frac{1}{|[p]|} \sum_{i \in [p]} \mathbf{z}_i\|}$ . The extracted features  $\{\mathbf{z}_i | 1 \leq i \leq n\}$  satisfy: (1)  $\forall 1 \leq i \leq n$ ,  $\|\mathbf{z}_i\| = 1$  and  $\forall i, j \in [p]$ ,  $\mathbf{z}_i^\top \mathbf{z}_j > 0$ . (2)  $\forall i \in [p]$ ,  $\exists j \in [p]$ ,  $\mathbf{z}_j = \|\mathbf{z}_j + \mathbf{v}_p\| \mathbf{v}_p - \mathbf{z}_i$ . (3)  $\forall p, q \in \mathcal{Y}$  with  $p \neq q$ ,  $\forall i \in [p]$  and  $j \in [q]$ ,  $\mathbf{v}_p^\top \mathbf{z}_i \geq \mathbf{v}_q^\top \mathbf{z}_i$  and  $\mathbf{v}_p^\top \mathbf{z}_i \geq \mathbf{v}_p^\top \mathbf{z}_j$ . (4)  $\forall p \in \mathcal{Y}$ ,  $\forall \mathbf{d} \in \Delta^{c-1}$  with  $p = \arg \max_k d^k$ ,  $\sum_{i \in [p]} \sum_{j \in [p]} d^m \mathbf{z}_i^\top \mathbf{z}_j > \sum_{i \notin [p]} \sum_{j \in [p]} (1 - d^p) \mathbf{z}_i^\top \mathbf{z}_j$ .

Here,  $\mathbf{v}_p$  is a normalized center of the extracted features belongs to class  $p$ . (1) means that the extracted feature  $\mathbf{z}_i$  has been normalized to a hypersphere with a radius of 1 (Zhou et al., 2022), and for two extracted features belongs to the same class, they do not go into opposite directions (Wang & Ma, 2022). (2) means that for each extracted feature of each class  $p$ , there is a symmetric point  $\mathbf{z}_j$  about the vector  $\mathbf{v}_p$ . (3) means that  $\mathbf{v}_q$  can be considered as the optimal prototype of each class  $q$ , and the extracted features  $\{\mathbf{z}_i | 1 \leq i \leq n\}$  have the separability (Li & Liang, 2018). (4) is an assumption that the correlation between the extracted features belonging to the same class is stronger than that ones belonging to different classes.

Let  $\angle(\cdot, \cdot)$  denote the angle between two vectors. Let  $o_p = \arg \min_{i \in [p]} \mathbf{v}_p^\top \mathbf{z}_i$ , and  $\mathbf{z}_{o_p}$  can be seen as the edge sample of the class  $p$ . According to Assumption 1, there also exists  $o'_p \in [p]$ , such that  $\mathbf{v}_p^\top \mathbf{z}_{o_p} = \mathbf{v}_p^\top \mathbf{z}_{o'_p}$ . Then, we make the following assumption about the prototypes:

**Assumption 2.**  $\forall k \in \mathcal{Y}$  with  $k \neq p$ ,  $\boldsymbol{\theta}_k = \mathbf{v}_k$ , and the normalized prototype  $\boldsymbol{\theta}_p$  satisfy: (1)  $\|\boldsymbol{\theta}_p\| = 1$ ; (2)  $\angle(\boldsymbol{\theta}_p, \mathbf{z}_{o_p}) + \angle(\boldsymbol{\theta}_p, \mathbf{v}_p) = \angle(\mathbf{z}_{o_p}, \mathbf{v}_p)$ ; (3)  $\forall q = \arg \min_{q \neq p} \angle(\mathbf{v}_p, \mathbf{v}_q)$ ,  $\angle(\mathbf{v}_q, \mathbf{z}_{o'_p}) < \angle(\mathbf{v}_q, \mathbf{z}_{o_p})$ . (4)  $\forall q \in \mathcal{Y}$  with  $q \neq p$ ,  $\forall i \in [q]$ , for  $q' = \arg \max_{q' \neq q} \boldsymbol{\theta}_{q'}^\top \mathbf{z}_i$ ,  $p \neq q'$ .

Here, (1) means that the final prototypes used to classify the extracted features are also on a hypersphere with a radius of 1 (Zhou et al., 2022). (2) means that the prototype  $\theta_p$  lies between the  $v_p$  and  $z_{o'_p}$ . (3) means that the prototype  $v_q$ , which is nearest to  $v_p$ , is closer to  $z_{o'_p}$  than  $z_{o_p}$ . (4) means that the change of  $\theta_p$  will only have effect on  $\gamma_p$ .

Then, we could obtain the following lemma:

**Lemma A.2.** *Under Assumption 2, suppose that  $\Theta$  with  $\theta_p$  and  $\Theta'$  with  $\theta'_p$  satisfy Assumption 2, and  $\theta_p^\top v_p - \theta'_p^\top v_p > 0$ , we have  $\gamma_p - \gamma'_p > 0$ .*

*Proof.* The class margin  $\gamma_p$  and  $\gamma'_p$  could be computed as  $\gamma_p = \theta_p^\top z_{o'_p} - \theta_q^\top z_{o'_p}$  and  $\gamma'_p = \theta'_p{}^\top z_{o'_p} - \theta_q^\top z_{o'_p}$  with some  $q \neq p$ , according to Assumption 2. Then we have:

$$\begin{aligned} \gamma_p - \gamma'_p &= \theta_p^\top z_{o'_p} - \theta'_p{}^\top z_{o'_p} \\ &= \cos \angle(\theta_p, z_{o'_p}) - \cos \angle(\theta'_p, z_{o'_p}) \\ &= \cos \left( \angle(\theta_p, v_p) + \angle(v_p, z_{o'_p}) \right) - \cos \left( \angle(\theta'_p, v_p) + \angle(v_p, z_{o'_p}) \right) \end{aligned} \quad (27)$$

Since  $\theta_p^\top v_p - \theta'_p{}^\top v_p > 0$ , we could get  $\angle(\theta_p, v_p) < \angle(\theta'_p, v_p)$ . Hence,  $\gamma_p - \gamma'_p > 0$ .

**Theorem 2.** *Let  $\Theta$  be the parameters of the last model layer trained by Gradient Descend starting from random initialization. Suppose that at epoch  $T$ ,  $\forall k \in \mathcal{Y}$  with  $k \neq p$ ,  $\theta_k$  arrives at the optimal prototype  $v_k$ . Then, we fix the extracted features  $\{z_i | 1 \leq i \leq n\}$  and the prototype except  $\theta_p$  to continue the training process for  $\theta_p$ . Let  $G^s(\mathcal{F}, \mathcal{X}, \mathcal{Y})$ ,  $G^h(\mathcal{F}, \mathcal{X}, \mathcal{Y})$  denote the generalization error bound based on the class margin derived from the empirical risk estimators with soft pseudo-labels and initial hard labels, respectively. At epoch  $T' > T$ , we could have*

$$G^s(\mathcal{F}, \mathcal{X}, \mathcal{Y}) < G^h(\mathcal{F}, \mathcal{X}, \mathcal{Y}).$$

*Proof.*  $\forall p \in \mathcal{Y}$ , the gradient of  $\mathcal{L}(t)$  with respect to  $\theta_p^\top(t)$  can be derived as follows:

$$\frac{\partial \mathcal{L}(t)}{\partial \theta_p^\top(t)} = \frac{1}{n} \sum_{i=1}^n \left( \sum_{j \neq p} w_i^j \frac{e^{(\theta_p(t) - \theta_j(t))^\top z_i}}{1 + \sum_{k \neq j} e^{(\theta_k(t) - \theta_j(t))^\top z_i}} - w_i^p \frac{\sum_{k \neq p} e^{(\theta_k(t) - \theta_j(t))^\top z_i}}{1 + \sum_{k \neq p} e^{(\theta_k(t) - \theta_j(t))^\top z_i}} \right) z_i^\top \quad (28)$$

Let  $\mathcal{L}^h$  and  $\mathcal{L}^s$  denote the risk estimators using the hard labels and soft labels, respectively. The gradient difference  $\Delta(t)$  with respect to the prototype  $\theta_p(t)$  between them can be derived as follows:

$$\begin{aligned} \Delta(t) &= \frac{\partial \mathcal{L}^h(t)}{\partial \theta_p^\top(t)} - \frac{\partial \mathcal{L}^s(t)}{\partial \theta_p^\top(t)} \\ &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{j \neq p} (l_i^j - s_i^j) \frac{e^{(\theta_p(t) - \theta_j(t))^\top z_i}}{1 + \sum_{k \neq j} e^{(\theta_k(t) - \theta_j(t))^\top z_i}} - (l_i^p - s_i^p) \frac{\sum_{k \neq p} e^{(\theta_k(t) - \theta_j(t))^\top z_i}}{1 + \sum_{k \neq p} e^{(\theta_k(t) - \theta_j(t))^\top z_i}} \right) z_i^\top \\ &= \frac{1}{n} \left( \sum_{i \in [m]} \left( \sum_{j \neq p} (l_i^j - s_i^j) \frac{e^{(\theta_p(t) - \theta_j(t))^\top z_i}}{1 + \sum_{k \neq j} e^{(\theta_k(t) - \theta_j(t))^\top z_i}} - (l_i^p - s_i^p) \frac{\sum_{k \neq p} e^{(\theta_k(t) - \theta_j(t))^\top z_i}}{1 + \sum_{k \neq p} e^{(\theta_k(t) - \theta_j(t))^\top z_i}} \right) z_i^\top \right. \\ &\quad \left. + \sum_{i \notin [m]} \left( \sum_{j \neq p} (l_i^j - s_i^j) \frac{e^{(\theta_p(t) - \theta_j(t))^\top z_i}}{1 + \sum_{k \neq j} e^{(\theta_k(t) - \theta_j(t))^\top z_i}} - (l_i^p - s_i^p) \frac{\sum_{k \neq p} e^{(\theta_k(t) - \theta_j(t))^\top z_i}}{1 + \sum_{k \neq p} e^{(\theta_k(t) - \theta_j(t))^\top z_i}} \right) z_i^\top \right) \\ &= \frac{1}{n} \left( \sum_{i \in [p]} (s_i^p - 1) z_i^\top + \sum_{i \notin [p]} s_i^p z_i^\top \right) \end{aligned} \quad (29)$$

We investigate the difference between the prototype  $\theta_p(T')$  and  $\theta_p(T)$  when the parameters are trained by Gradient Decent

with the step size  $\eta > 0$ :

$$\begin{aligned}
 \boldsymbol{\theta}_p(T') &= \boldsymbol{\theta}_p(T' - 1) - \eta \frac{\partial \mathcal{L}(T' - 1)}{\partial \boldsymbol{\theta}_p(T' - 1)} \\
 &= \boldsymbol{\theta}_p(T' - 2) - \eta \frac{\partial \mathcal{L}(T' - 1)}{\partial \boldsymbol{\theta}_p(T' - 1)} - \eta \frac{\partial \mathcal{L}(T' - 2)}{\partial \boldsymbol{\theta}_p(T' - 2)} \\
 &= \boldsymbol{\theta}_p(T) - \sum_{r=T}^{T'-1} \eta \frac{\partial \mathcal{L}(r)}{\partial \boldsymbol{\theta}_p(r)}.
 \end{aligned} \tag{30}$$

Hence, we could obtain:

$$\boldsymbol{\theta}_p(T') - \boldsymbol{\theta}_p(T) = \sum_{r=T}^{T'-1} \eta \left( -\frac{\partial \mathcal{L}(r)}{\partial \boldsymbol{\theta}_p(r)} \right). \tag{31}$$

We transpose both sides of Eq. (31) and then right-multiply with the vector  $\mathbf{v}_p$ :

$$\boldsymbol{\theta}_p^\top(T') \mathbf{v}_p - \boldsymbol{\theta}_p^\top(T) \mathbf{v}_p = \sum_{r=T}^{T'-1} \eta \left( -\frac{\partial \mathcal{L}(r)}{\partial \boldsymbol{\theta}_p(r)} \right)^\top \mathbf{v}_p. \tag{32}$$

Let  $\boldsymbol{\theta}_p^h$  and  $\boldsymbol{\theta}_p^s$  denote the prototypes updated by  $\mathcal{L}^h$  and  $\mathcal{L}^s$ , respectively. Based on Eq. (32), we could obtain:

$$\boldsymbol{\theta}_p^{s\top}(T') \mathbf{v}_p - \boldsymbol{\theta}_p^{h\top}(T') \mathbf{v}_p = \sum_{r=T}^{T'-1} \eta \left( \frac{\partial \mathcal{L}^h(r)}{\partial \boldsymbol{\theta}_p^h(r)} - \frac{\partial \mathcal{L}^s(r)}{\partial \boldsymbol{\theta}_p^s(r)} \right)^\top \mathbf{v}_p. \tag{33}$$

Due to that  $\mathbf{s}_i \in \Delta^{c-1}$ , we could obtain the following according to Assumption 1.(4)

$$\boldsymbol{\theta}_p^{s\top}(T') \mathbf{v}_p - \boldsymbol{\theta}_p^{h\top}(T') \mathbf{v}_p > 0. \tag{34}$$

According to Lemma A.1 and A.2, we have  $\gamma_p^s > \gamma_p^h$  and finally obtain  $G^s(\mathcal{F}, \mathcal{X}, \mathcal{Y}) < G^h(\mathcal{F}, \mathcal{X}, \mathcal{Y})$ .