

COMPLEMENTARY LABEL LEARNING WITH POSITIVE LABEL GUESSING AND NEGATIVE LABEL ENHANCEMENT

Anonymous authors

Paper under double-blind review

ABSTRACT

Complementary label learning (CLL) is a weakly supervised learning paradigm that constructs a multi-class classifier only with complementary labels, specifying classes that the instance does not belong to. We reformulate CLL as an inverse problem that infers the full label information from the output space information. To be specific, we propose to split the inverse problem into two subtasks: *positive label guessing* (PLG) and *negative label enhancement* (NLE), collectively called PLNL. Specifically, we use well-designed criteria for evaluating the confidence of the model output, accordingly divide the training instances into three categories: highly-confident, moderately-confident and under-confident. For highly-confident instances, we perform PLG to assign them pseudo labels for supervised training. For moderately-confident and under-confident instances, we perform NLE by enhancing their negative label set with different levels and train them with the augmented negative labels iteratively. In addition, we unify PLG and NLE into a consistent framework, in which we can view all the pseudo-labeling-based methods from the perspective of negative label recovery. We prove that the error rates of both PLG and NLE are upper bounded, and based on that we can construct a classifier consistent with that learned by clean full labels. Extensive experiments demonstrate the superiority of PLNL over the state-of-the-art CLL methods, e.g., on STL-10, we increase the classification accuracy from 34.96% to 55.25%. **The code has been submitted to supplementary material.**

1 INTRODUCTION

Over the past few years, large-scale and accurately labeled data has tremendously boosted the development of deep neural networks. However, collecting accurately labeled data is extremely time-consuming, labor-intensive and sometimes requires specific expertise in real-world tasks. To reduce the dependency on large-scale and accurately labeled datasets, deep learning communities have given increasing attention to weakly supervised learning, including but not limited to partial label learning (Cour et al., 2011; Xie and Huang, 2018; Feng and An, 2019; Lv et al., 2020; Xia et al., 2023; Huang and Cheung, 2024; He et al., 2024; Tian et al., 2024), noisy label learning (Natarajan et al., 2013; Han et al., 2018; Song et al., 2022; Wei et al., 2020; Zhang et al., 2023; Huang et al., 2023), semi-supervised learning (Van Engelen and Hoos, 2020; Sohn et al., 2020; Xie et al., 2020; Yang et al., 2022; Li et al., 2023b; Xie et al., 2023), positive-unlabeled learning (Niu et al., 2016; Kiryo et al., 2017).

Here, we consider a recently proposed weakly supervised learning framework called complementary label learning (CLL) (Ishida et al., 2017; Feng et al., 2020). In CLL, each training instance is associated with one or multiple complementary labels (CLs) which specify one or multiple classes that the instance does not belong to. The goal of CLL is to learn a multi-class classifier only from complementary labeled data. In real-world scenario, if the number of classes is huge, choosing the correct class label from many candidate classes is difficult and laborious, while choosing one or several of the incorrect class labels as CLs would be much easier and thus less costly. Recently, CLL has been applied to online learning (Kaneko et al., 2019), medical image segmentation (Rezaei et al., 2020) and medical molecular imaging (Tapper et al., 2024), etc. Besides, another promising future application scenario of CLL is to ensure privacy security in data collection scenarios. For example,

054 collecting some survey data may require extremely private questions and it would be mentally less
055 demanding if we ask the respondent to provide some incorrect answers as CLs (Dwork, 2008).
056

057 Previous studies on CLL can be roughly divided into two categories: methods that attempt to construct
058 an unbiased risk estimator (URE-based) (Ishida et al., 2017; 2019; Feng et al., 2020) and methods
059 based on feature learning (FL-based) (Chou et al., 2020; Wang et al., 2021; Liu et al., 2022; Jiang
060 et al., 2024). For URE-based methods, Ishida et al. (Ishida et al., 2017) and Feng et al. (Feng et al.,
061 2020) showed that the ordinary classification risk can be recovered by their proposed unbiased risk
062 estimator only from complementary labeled data. Ishida et al. (Ishida et al., 2019) later extended
063 the unbiased risk estimator to arbitrary losses and models. Chou et al. (Chou et al., 2020) proposed
064 a surrogate complementary loss framework, which avoids the extremely noisy gradient problem
065 encountered in unbiased risk estimator. For FL-based methods, Wang et al. (Wang et al., 2021)
066 gave the first attempt to leverage regularization techniques with complementary label by aligning the
067 model output of one instance and its multiple augmented views. Liu et al. (Liu et al., 2022) proposed
068 to integrate self-supervised and self-distillation to complementary learning. Jiang et al. (Jiang et al.,
069 2024) leveraged a contrastive learning framework to facilitate CLL. These methods mainly focus on
070 the design of robust loss functions or the exploration of feature space information, while neglecting
the power of output space information.

071 We propose that CLL can be viewed as solving the multi-class classification problem from two
072 inverse aspects, where one is to infer the positive label and another is to infer the negative labels. To
073 this end, we propose two subtasks: *positive label guessing* (PLG) and *negative label enhancement*
074 (NLE). We use well-designed criteria for evaluating the confidence of the model output, accordingly
075 divide the training instances into three categories: highly-confident, moderately-confident and under-
076 confident in each epoch. We perform PLG by simply pseudo-labeling highly-confident instances
077 for supervised training. Unlike pseudo-labeling methods used in semi-supervised learning (SSL),
078 PLG pseudo-labeling reaches high selected ratio and high precision even without any positive labels
079 available.

080 More importantly, previous SSL methods lack the utilization of untrustworthy instances. They either
081 discard this part or simply employ techniques such as consistency regularization. In this paper,
082 we perform NLE by enhancing the negative label set of moderately-confident and under-confident
083 instances and train them with the augmented negative labels iteratively.

084 Although PLG and NLE will inevitably bring pseudo-labeling errors, we theoretically prove that the
085 error rates are upper bounded. And the generalization error of the learned classifier under PLG and
086 NLE errors is also upper bounded, which means that we can construct a classifier consistent with that
087 learned by clean full labels. We demonstrate that PLNL achieves state-of-the-art performance on five
088 benchmark datasets. Our contributions can be summarized as follows:

- 089 • *A novel method for CLL.* Different from conventional loss design methods, we pioneer a
090 novel method for CLL called PLNL that formulates CLL from output space information and
091 solve it by two subtasks: PLG and NLE.
- 092 • *A unified framework for pseudo-labeling-based methods.* From the perspective of negative
093 label recovery, we construct a unified framework for pseudo-labeling-based methods. We
094 empirically show that PLNL outperforms state-of-the-art SSL methods in terms of pseudo-
095 labeling error, selected ratio and recovered negative labels.
- 096 • *Solid theoretical analysis.* We theoretically prove that both the error rates of PLG and NLE
097 are upper bounded. The generalization error of the learned classifier is also upper bounded.
- 098 • *State-of-the-art performance.* Extensive experiments on five benchmark datasets demon-
099 strate the superiority of PLNL over the state-of-the-art CLL methods.

102 2 PRELIMINARIES

103 **Ordinary Multi-Class Classification.** Let $\mathcal{X} \in \mathbb{R}^d$ denote the feature space with d dimensions
104 and $\mathcal{Y} = \{1, 2, \dots, K\}$ denote the label space with K classes. The precisely labeled instance
105 $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ is sampled from an unknown probability distribution $p(\mathbf{x}, y)$. The goal of ordinary
106 multi-class classification is to learn a parameterized function $f(x) : \mathbb{R}^d \rightarrow \mathbb{R}^K$ that minimizes the
107

classification risk:

$$R(f) = \mathbb{E}_{p(\mathbf{x}, y)} \mathcal{L}(f(\mathbf{x}), y), \quad (1)$$

where $\mathbb{E}_{p(\mathbf{x}, y)}$ refers to the expectation across all possible samples drawn from the distribution $p(\mathbf{x}, y)$, $\mathcal{L} : \mathbb{R}^K \times \mathcal{Y} \rightarrow \mathbb{R}$ is a multi-class classification loss function. In this paper, we consider a common case where the function f is a deep neural network with the softmax output layer, where $f(\mathbf{x})$ is considered as the output prediction confidence of the model on each class. Since the probability distribution $p(\mathbf{x}, y)$ is unknown, we use the empirical risk $\hat{R}(f)$ to approximate $R(f)$. Assuming a dataset $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ is independently drawn from distribution $p(\mathbf{x}, y)$, then we have

$$\hat{R}(f) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i), y_i). \quad (2)$$

Complementary Label Learning. Different from the ordinary multi-class classification, in CLL, let $\{(\mathbf{x}_i, \bar{Y}_i)\}_{i=1}^N$ be the complementary labeled dataset, where N is the dataset size, \bar{Y}_i indicates the complementary (negative) label set of \mathbf{x}_i . Each complementary labeled instance $(\mathbf{x}, \bar{Y}) \in \mathcal{X} \times \mathcal{Y}$ is sampled from an unknown probability distribution $\bar{p}(\mathbf{x}, \bar{y})$. Our goal is to learn a classifier that minimizes the classification risk Eq. (1) only from complementary labeled training instances. Then the empirical risk becomes:

$$\hat{R}(f) = \frac{1}{N} \sum_{i=1}^N \bar{\mathcal{L}}_{CLL}(f(\mathbf{x}_i), \bar{Y}_i), \quad (3)$$

where $\bar{\mathcal{L}}_{CLL}$ is a specially designed loss function for learning from only complementary labeled data.

3 PROPOSED METHOD

The overall framework of PLNL is shown in Fig. 1, and the pseudo-code is presented in Appendix A. We begin by employing weak and strong augmentation to a complementary labeled image \mathbf{x}_i , which leads to two augmented views $\mathbf{x}_i^w, \mathbf{x}_i^s$. These two images are then fed into a two-view network with shared weight $f(\mathbf{x}; \Theta)$ to obtain two prediction confidences $f(\mathbf{x}_i^w)$ and $f(\mathbf{x}_i^s)$. Then we utilize the two-view prediction confidences to select three subsets of training instances mentioned above, i.e., highly-confident, moderately-confident and under-confident. We select these subsets using the historical confidences of the previous training epochs to better alleviate confirmation bias. Finally, different techniques are utilized to conquer individual subsets. In this section, we first explain the well-designed confidence-based instances selection strategy, and then introduce the PLG for the highly-confident instances and two different versions of NLE for the moderately-confident instances and the under-confident instances in detail.

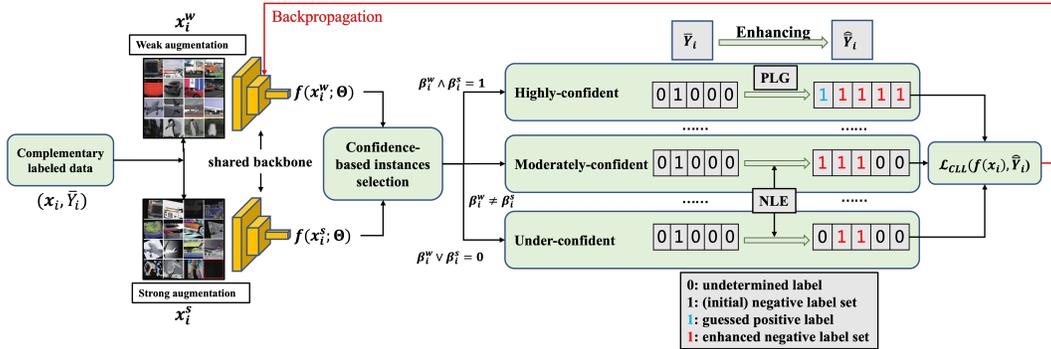


Figure 1: The overall framework of PLNL. We employ a two-view network (shared backbone) to extract features and compute confidences for weak and strong augmentations of one instance respectively. After the selection of highly-confident, moderately-confident and under-confident set, we employ PLG on highly-confident instances and NLE for the rest. The loss is computed on the enhanced labels of both views and the model updates through backpropagation.

3.1 CONFIDENCE-BASED INSTANCES SELECTION

We first maintain two memory banks M^w and M^s for weak and strong augmentation respectively, each with size $t \times N \times K$ to store the historical prediction confidence over the past t epochs. For simplification, we only consider one view here unless otherwise specified.

$$M_i = [f^1(\mathbf{x}_i), \dots, f^t(\mathbf{x}_i)]. \quad (4)$$

where $f^t(\mathbf{x}_i)$ denotes the prediction confidence of the t -th epoch in the memory bank M .

We propose to select subsets based on the following criteria.

$$\omega_{i1} = \forall 1 \leq j \leq t, \arg \max(f^j(\mathbf{x}_i)) \notin \bar{Y}_i, \quad (5)$$

$$\omega_{i2} = \forall 1 \leq j, k \leq t, \arg \max(f^j(\mathbf{x}_i)) = \arg \max(f^k(\mathbf{x}_i)), \quad (6)$$

$$\omega_{i3} = \forall 1 \leq j \leq t, \max(f^j(\mathbf{x}_i)) \geq \lambda, \quad (7)$$

where $\omega_{i1}, \omega_{i2}, \omega_{i3}$ are boolean variables which indicate whether the corresponding criterion is satisfied. ω_{i1} ensures that the label corresponding to the max prediction confidence does not fall on the complementary label set, which excludes the complementary labels from being selected as positive label. ω_{i2} ensures that the max prediction confidence be stable and show no sign of fluctuations over the past t epochs. ω_{i3} ensures the max prediction confidence should be higher than a threshold λ . Note that λ can be either a man-made threshold or a self-adaptive one, which will be discussed in detail later.

Warm up. Before selecting, we warm up the model using the entire training set. The goal of this stage is to reduce the classification risk and obtain some historical prediction confidence to construct the memory bank since we have no historical information at initial epoch. In this paper, we use SCL-LOG algorithm (Chou et al., 2020) to warm up models for 20 epochs.

Instance-aware self-adaptive threshold. The threshold in criterion ω_{i3} , as is mentioned above, can be a fixed high threshold (like 0.95). However, a single global threshold does not consider the fitting difficulties of different instances, i.e., hard instances and easy instances. This will result in very few samples being selected in the early training stages as well as confirmation bias. Therefore, to comprehensively consider historical information, we design an instance-aware self-adaptive threshold for each instance at each epoch t as:

$$\lambda(t) = \alpha\lambda(t-1) + (1-\alpha)f(t), \lambda(0) = \frac{1}{K}, \quad (8)$$

where K is the number of classes, $f(t) = \max f(\mathbf{x})$ and α is the ratio which controls the threshold stability.

The threshold is initialized at a low value $\frac{1}{K}$, which will take more data into account and helps speed up convergence in the early stages. As the prediction confidence increases, the threshold grows higher to filter out wrong pseudo labels to alleviate the confirmation bias. Note that we compute $\lambda^w(t)$ and $\lambda^s(t)$ for two different views respectively according to Eq. (8). We use the momentum average confidence of each instance, computed based on all previous epochs. In this way, the threshold comprehensively considers historical information and remains stable and trustworthy.

Subset Selection. For the two-view network, we perform two independent verifications. Let $\beta_i^w = \omega_{i1} \wedge \omega_{i2} \wedge \omega_{i3}$ be the indicator of satisfying the criteria. Thus, β_i^w and β_i^s indicate whether the weak and strong views meet the criteria respectively. For an instance, if both views meet the criteria, we select it to the highly-confident subset, i.e.,

$$\mathcal{H} = \{\mathbf{x}_i | \beta_i^w \wedge \beta_i^s = 1\}, \quad (9)$$

It means the prediction confidences of both views are stable and high, thus we consider them to be highly-confident. The size of \mathcal{H} is denoted as N_h .

Similarly, the moderately-confident subset consists of instances only one augmented version of which meet the criteria, i.e.,

$$\mathcal{M} = \{\mathbf{x}_i | \beta_i^w \neq \beta_i^s\}, \quad (10)$$

It means only one view’s prediction is trustworthy, the other is not. which shows that the model is moderately-confident about its prediction. The size of \mathcal{M} is denoted as N_m .

Finally, the under-confident subset consists of the rest of the instances, i.e.,

$$\mathcal{U} = \{\mathbf{x}_i | \beta_i^w \vee \beta_i^s = 0\}, \quad (11)$$

It means the prediction confidences of both views do not meet the designed criteria, these instances are considered under-confident. The size of \mathcal{U} is denoted as N_u .

After obtaining \mathcal{H} , \mathcal{M} and \mathcal{U} , we design different techniques to better utilize these different types of training instances.

3.2 POSITIVE LABEL GUESSING

For highly-confident set \mathcal{H} , we consider the label with the max prediction confidence as its positive label. **Conversely, all remaining labels are considered complementary labels.** Let \hat{Y}_i be the enhanced negative label set for instance \mathbf{x}_i , we have:

$$\hat{Y}_i = \{c | c \in Y_i, c \neq \hat{y}_i\} \quad (12)$$

where \hat{y}_i is the guessed positive label and $Y_i = \{1, 2, \dots, K\}$ is the full label set.

For highly-confident set \mathcal{H} , we compute the CLL loss on the negative labels for both views:

$$\mathcal{L}_h = \frac{1}{N} \sum_{i=1}^{N_h} \bar{\mathcal{L}}_{CLL}(f(\mathbf{x}_i^w), \hat{Y}_i) + \bar{\mathcal{L}}_{CLL}(f(\mathbf{x}_i^s), \hat{Y}_i) \quad (13)$$

where $f(\mathbf{x}_i^w)$ and $f(\mathbf{x}_i^s)$ denote model outputs of weak augmentation and strong augmentation respectively.

3.3 NEGATIVE LABEL ENHANCEMENT

For the moderately-confident set \mathcal{M} and the under-confident set \mathcal{U} , guessing the positive labels directly might lead to much more errors due to their relatively lower confidence. Therefore, we employ a different strategy for these instances, called negative label enhancement (NLE).

The rationale of NLE is that more negative labels will bring in additional supervision information for better training. However, whether the enhanced negative labels are correct remains a question. Intuitively, randomly enhancing negative labels will bring in a large number of labeling errors. To better enhance the reliability of NLE, we further exploit information in the output space and design the following solution.

Calculation of k Nearest Neighbor (k -NN) instances. For instance \mathbf{x}_i and its model output prediction confidence y_i , we can compute its k -NN instances in the output space. It is safe to assume that nearby instances in the output space should have the same positive label with a high probability, while their original complementary label sets are likely to vary. The formal definition of this assumption is as follows:

Assumption 1. $\forall(\mathbf{x}_i, \bar{Y}_i) \in \mathcal{D}$ and its k -NN instances $(\mathbf{x}_i^{(j)}, \bar{Y}_i^{(j)})$, the positive label y_i exists in its k -NN instances’ complementary label set $\bar{Y}_i^{(j)}$ with probability no more than α_k , any negative label $y'_i \neq y_i$ exist in its k -NN instances’ complementary label set $\bar{Y}_i^{(j)}$ with probability no less than β_k .

This assumption describes the intrinsic characteristics of CLL in the output space, which can be interpreted in two aspects. First, similarity in the input space will be mapped to similarity in the output space, which has been widely utilized for tackling representation learning problems (He et al., 2020). Second, instances of the same category are likely to be labeled with complementary labels of different categories, which is key to enhancing the negative labels.

k -NN label frequency. For instance \mathbf{x}_i , we propose to calculate the times a negative label appears in its k -NN instances’ complementary label set and then enhance top- τ_i frequent ones, that is, add

them to the complementary label set of \mathbf{x}_i . We define the j -th k -NN label frequency of \mathbf{x}_i as follows:

$$\mathbf{F}_{ij} = \sum_{v=1}^k \mathbb{I}(j \in \bar{Y}_i^{(v)}), \quad (14)$$

where $\bar{Y}_i^{(v)}$ denotes the complementary label set of the v -th nearest instance of \mathbf{x}_i .

Negative label enhancement. We enhance the complementary label set \bar{Y}_i by adding additional labels with top- τ_i label frequency. For each instance \mathbf{x}_i in \mathcal{M} and \mathcal{U} , the enhanced complementary (negative) label set \hat{Y}_i is calculated by:

$$\hat{Y}_i = \{c | c \in \bar{Y}_i \vee c \in \text{top-}\tau_i\text{-max}_j(\mathbf{F}_{ij})\}. \quad (15)$$

However, the prediction confidence of the under-confident is more unreliable than that of the moderately-confident. Therefore, we should be more conservative when enhancing these instances as the k -NN information may be more unreliable. In our work, we set $\tau_i = \lceil \frac{K-s_i}{10} \rceil$ for \mathcal{U} where s_i is the size of \bar{Y}_i . For \mathcal{M} , we set $\tau_i = (1 + \frac{e}{E_{max}}) \lceil \frac{K-s_i}{10} \rceil$ where e is current epoch, E_{max} is total epochs. This provides a linear growing strategy for the moderately-confident because the model’s output becomes increasingly accurate as the training progresses.

For moderately-confident set \mathcal{M} and under-confident set \mathcal{U} , we compute the CLL loss on the negative labels for both views:

$$\mathcal{L}_{m,u} = \frac{1}{N} \sum_{i=1}^{N_m+N_u} \bar{\mathcal{L}}_{CLL}(f(\mathbf{x}_i^w), \hat{Y}_i) + \bar{\mathcal{L}}_{CLL}(f(\mathbf{x}_i^s), \hat{Y}_i) \quad (16)$$

where $f(\mathbf{x}_i^w)$ and $f(\mathbf{x}_i^s)$ denote model outputs of weak augmentation and strong augmentation respectively.

4 A UNIFIED FRAMEWORK FOR PSEUDO-LABELING-BASED METHODS

Pseudo-labeling, which has been widely used in the recent semi-supervised learning (SSL) methods, is employed by giving unlabeled instances pseudo labels and train them in a supervised way. PLNL is an extension of pseudo-labeling. We not only consider pseudo-labeling of highly-confident instances, but also consider enhancing the negative label set of untrustworthy instances. In this way, we actually recover more supervised information than only leveraging pseudo-labeling and further boost the classification performance.

In this section, we construct a unified framework where PLG and NLE are viewed from the perspective of negative label recovery. Let \hat{y}_i be the pseudo-label of \mathbf{x}_i . Let Y_i be the full label space. Let \hat{Y}_i be the reconstructed (PLNL) or imposed (pseudo-labeling) negative label set for \mathbf{x}_i .

For PLNL, PLG is equivalent to reconstructing a negative label set $\hat{Y}_i = \{c | c \in Y_i, c \neq \hat{y}_i\}$ of size $K - 1$, in which only the guessed positive label does not belong. NLE is equivalent to reconstructing a negative label set of size $s_i + \tau_i$, where we add τ_i negative labels to the original negative label set of size s_i .

Similarly, for pseudo-labeling methods, let the pseudo label for instance \mathbf{x}_i be \hat{y}_i as well. The process of pseudo-label highly-confident instances is also equivalent to imposing a negative label set $\hat{Y}_i = \{c | c \in Y_i, c \neq \hat{y}_i\}$ of size $K - 1$ as additional supervised information.

In this paper, we propose two metrics for evaluation of pseudo-labeling-based methods. Firstly, we define selected ratio η :

$$\eta = \frac{N_h}{N}, \quad (17)$$

Obviously, η evaluates the ratio of highly-confident instances selected for pseudo-labeling methods.

Furthermore, we define average size of enhanced negative label set \bar{s} :

$$\bar{s} = \frac{\sum_{i=1}^N |\hat{Y}_i|}{N}, \quad (18)$$

In section 6, we empirically show that PLNL achieves lower error rate ϵ , higher selection ratio η and obviously larger size of negative label set \bar{s} compared with pseudo-labeling method Fixmatch (Sohn et al., 2020) and Freematch (Wang et al., 2022).

5 THEORETICAL ANALYSIS

5.1 GENERALIZATION BOUND

For simplification, we only consider one view network here, which has no influence on the deduction of generalization error bound. Our goal is to learn a classification model $f(\mathbf{x}; \Theta)$ by minimizing the empirical risk $\widehat{R}'(f)$ acquired from data with enhanced negative labels:

$$\widehat{R}'(f) = \frac{1}{N} \sum_{i=1}^N \widehat{\mathcal{L}}_{CLL}(f(\mathbf{x}_i), \widehat{Y}_i), \quad (19)$$

where \widehat{Y}_i denotes the enhanced negative label set of \mathbf{x}_i .

Let the CLL loss function be $\widehat{\mathcal{L}}_{CLL}(f(\mathbf{x}), \widehat{Y}_i) = \sum_{y \notin \widehat{Y}_i} (1/(K - |\widehat{Y}_i|)) \ell(f(\mathbf{x}), y)$ where $|\widehat{Y}_i|$ is the size of the enhanced negative label set. Let $\bar{s} = \frac{\sum_{i=1}^N |\widehat{Y}_i|}{N}$ be the average size of enhanced negative label set. Let $\epsilon_1 = \frac{\sum_{i=1}^{N_h} \mathbb{I}(y_i \in \widehat{Y}_i)}{N_h}$ be the error rate of PLG. Let $\epsilon_2 = \frac{\sum_{i=1}^{N_m+N_u} \mathbb{I}(y_i \in \widehat{Y}_i)}{N_m+N_u}$ be the error rate of NLE. The actual pseudo-labeling error rate $\epsilon = \frac{\sum_{i=1}^N \mathbb{I}(y_i \in \widehat{Y}_i)}{N} = \frac{N_h}{N} \epsilon_1 + \frac{N_m+N_u}{N} \epsilon_2 = \eta \epsilon_1 + (1 - \eta) \epsilon_2$. Moreover, $\ell(f(\mathbf{x}), y)$ is ρ -Lipschitz w.r.t. $f(\mathbf{x})$ where ρ can be any Lipschitz constant (not necessarily the best). Let $\mathfrak{R}_N(\mathcal{F})$ be the expected Rademacher complexity (Mohri et al., 2018) of \mathcal{F} with N training instances. Let $\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \widehat{R}'(f)$ be the empirical risk minimizer, where \mathcal{F} is a function class, and $f^* = \operatorname{argmin}_{f \in \mathcal{F}} R(f)$ be the true risk minimizer. We derive the following theorem, which provides a generalization error bound for the proposed method.

Theorem 1. *Suppose that $\ell(f(\mathbf{x}), y)$ is bounded by B . For pseudo-labeling error rate $\epsilon \in (0, 1)$, for any $\delta > 0$, with probability at least $1 - \delta$, we have*

$$R(\hat{f}) - R(f^*) \leq 2(1 - \frac{1 - \epsilon}{K - \bar{s}})B + 4\rho K \mathfrak{R}_N(\mathcal{F}) + 2KB \sqrt{\frac{\log \frac{2}{\delta}}{2N}}, \quad (20)$$

Remark. Detailed proofs are provided in Appendix B. Theorem 1 shows that as $N \rightarrow \infty$, $\epsilon_1 \rightarrow 0$, $\epsilon_2 \rightarrow 0$, the empirical risk minimizer converges to the true risk minimizer with high probability. It can be observed from Eq. (20) that the generalization bound is influenced by five factors: the number of categories K , the average size of enhanced negative label set \bar{s} and two error rates. This is consistent with the intuition that more categories and less complementary labels will make the CLL problem harder. **In a nutshell, smaller PLNL pseudo-labeling error rates ϵ_1, ϵ_2 and larger size of enhanced negative label set \bar{s} will produce better generalization performance.**

5.2 ERROR BOUND OF POSITIVE LABEL GUESSING

Theorem 2. *Suppose that y_i denote the ground-truth positive label of \mathbf{x}_i and \hat{y}_i denote the guessed positive label which might not be true. PLG error rate ϵ_1 is upper bounded by:*

$$\epsilon_1 = \mathbb{P}(y_i \in \widehat{Y}_i) \leq (K - 1 - s_i)\psi, \quad (21)$$

where K is class number, s_i is the size of \bar{Y}_i and $\psi \in (0, \frac{1}{K-1-s_i})$. Detailed proofs are provided in Appendix C.

5.3 ERROR BOUND OF NEGATIVE LABEL ENHANCEMENT

Theorem 3. *Suppose that y denote the ground-truth positive label of \mathbf{x}_i and y' denote an arbitrary negative label. Let $F_i^{(\tau_i)}$ denote the τ_i -th largest label frequency. Let p denote the probability of the*

ground-truth positive label y_i appearing in its k -NN instance’s complementary label set. Let q denote the probability of the label y' appearing in its k -NN instance’s complementary label set. The NLE error rate ϵ_2 is upper bounded by:

$$\epsilon_2 = \mathbb{P}(y_i \in \widehat{Y}_i) \leq \sum_{j=1}^k \binom{|Y_i| - 1}{|Y_i| - \tau_i} F_{\beta_k}(k - j + 1, j)^{(|Y_i| - \tau_i)} b_{\alpha_k}(k, j), \quad (22)$$

where $F_{\beta_k}(k, j) = \int_0^{\beta_k} p^{k-1}(1-p)^{j-1} dt$ denotes the regularized incomplete beta function, $b_{\alpha_k}(k, j) = \binom{k}{j} \alpha_k^j (1 - \alpha_k)^{k-j}$ is the probability mass function of a binomial distribution $B(k, \alpha_k)$. Detailed proofs are provided in Appendix D.

Remark. Theorem 2 and Theorem 3 show that both PLG error rate ϵ_1 and NLE error rate ϵ_2 are upper bounded under mild condition.

Table 1: Comparison of classification accuracies between different methods on four datasets with a single complementary label per instance. The results (mean \pm std) are reported over 3 random trials. The best results are highlighted in bold (The same applies hereinafter).

Method	STL-10	SVHN	FMNIST	CIFAR-10
UB-EXP	28.84 \pm 0.54%	88.93 \pm 0.17%	87.96 \pm 0.08%	62.90 \pm 0.06%
UB-LOG	20.41 \pm 0.46%	89.59 \pm 0.08%	87.59 \pm 0.14%	70.28 \pm 0.12%
SCL-EXP	31.03 \pm 0.61%	88.66 \pm 0.20%	88.31 \pm 0.09%	72.35 \pm 0.10%
SCL-LOG	30.74 \pm 0.72%	89.26 \pm 0.24%	88.03 \pm 0.10%	79.87 \pm 0.14%
POCR	34.96 \pm 0.32%	96.65 \pm 0.14%	92.29 \pm 0.07%	94.15 \pm 0.09%
SELF-CL	30.87 \pm 0.72%	90.13 \pm 0.23%	84.86 \pm 0.10%	88.95 \pm 0.22%
ComCo	32.43 \pm 0.28%	91.41 \pm 0.35%	85.42 \pm 0.40%	89.36 \pm 0.76%
Ours	55.25\pm0.36%	97.58\pm0.18%	93.38\pm0.06%	94.78\pm0.12%

Table 2: Comparison of classification accuracies between different methods on five datasets with multiple complementary labels per instance. The results (mean \pm std) are reported over 3 random trials.

Method	STL-10	SVHN	FMNIST	CIFAR-10	CIFAR-100
UB-EXP	60.85 \pm 0.12%	95.23 \pm 0.09%	92.34 \pm 0.28%	91.13 \pm 0.23%	34.43 \pm 0.08%
UB-LOG	62.84 \pm 0.17%	94.76 \pm 0.07%	91.84 \pm 0.29%	92.01 \pm 0.21%	52.76 \pm 0.15%
SCL-EXP	62.96 \pm 0.10%	95.28 \pm 0.14%	92.20 \pm 0.27%	91.85 \pm 0.25%	47.81 \pm 0.09%
SCL-LOG	61.60 \pm 0.14%	94.88 \pm 0.16%	91.51 \pm 0.25%	92.67 \pm 0.18%	49.40 \pm 0.19%
POCR	74.51 \pm 0.29%	97.14 \pm 0.09%	94.76 \pm 0.26%	96.09 \pm 0.27%	53.16 \pm 0.11%
SELF-CL	69.85 \pm 0.20%	91.58 \pm 0.30%	94.92 \pm 0.21%	92.23 \pm 0.16%	57.65 \pm 0.25%
ComCo	73.28 \pm 0.19%	95.41 \pm 0.23%	92.01 \pm 0.16%	91.38 \pm 0.73%	57.88 \pm 0.95%
Ours	77.11\pm0.14%	98.13\pm0.11%	95.16\pm0.13%	96.80\pm0.28%	64.33\pm0.43%

6 EXPERIMENT

6.1 EXPERIMENT SETUP

Datasets. We use five commonly used benchmark datasets, STL-10 (Coates et al., 2011), Fashion-MNIST (FMNIST) (Xiao et al., 2017), SVHN (Netzer et al., 2011), CIFAR-10 and CIFAR-100 (Krizhevsky and Hinton, 2009). We conduct experiments by considering both the scenarios of Single CLL (SCLL) and Multiple CLL (MCLL). To generate single complementary label, we randomly select one of the complementary classes per instance. To generate multiple complementary labels, let s be the size of \bar{Y} , we first instantiate $p(s) = \binom{K-1}{s} / (2^K - 2)$, $s \in \{1, 2, \dots, K-1\}$, which represents the possibility of randomly sample a complementary label set whose size is s from all possible complementary label sets which has $2^K - 2$ sets to choose from. Then for each instance

x_i , we first sample s_i from $p(s_i)$, and then sample a complementary label set \bar{Y}_i with size s from $p(\bar{Y}_i) = 1/\binom{K-1}{s_i}$.

Compared methods. We compare the performance of PLNL with seven state-of-the-art CLL methods, including UB-EXP (Feng et al., 2020), UB-LOG (Feng et al., 2020), SCL-EXP (Chou et al., 2020), SCL-LOG (Chou et al., 2020), POCR (Wang et al., 2021), SELF-CL (Liu et al., 2022) and ComCo (Jiang et al., 2024) and two state-of-the-art SSL methods, Fixmatch (Sohn et al., 2020) and Freematch (Wang et al., 2022).

Implementation. Implementation details are provided in Appendix F.

6.2 MAIN RESULTS

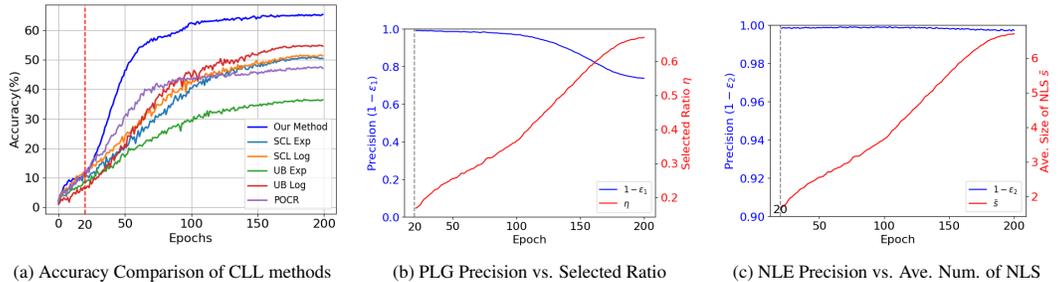


Figure 2: The experiments is conducted on CIFAR-100 with multiple complementary labels (MCLL). (a) The accuracy of PLNL improves tremendously over epochs and achieves the best finally. (b) The precision of PLG decreases slowly, while the selected ratio steadily rises, indicating a growing proportion of selected instances during training. (c) Ave. Size of NLS denotes average number of negative label set \bar{s} . The precision of NLE remains relatively stable with a slight decrease, whereas the average size of negative label set increases significantly, showing a steady recovery of negative labels.

PLNL achieves SOTA results. As shown in Table 1 and Table 9, PLNL outperforms all the compared method by a significant margin across all datasets. Specifically, on STL-10 dataset, we outperform the previous SOTA by **20.29%** and **1.61%** in both SCLL and MCLL settings. Furthermore, PLNL performs even better in harder scenarios where there is larger label space or less supervised information for each class. We challenge this by showing our results on CIFAR-100 datasets. On CIFAR-100 with Multiple CLs, the improvement is **6.45%** compared to previous SOTA. Fig. 2a further demonstrates that PLNL significantly outperforms the compared ones.

PLNL pseudo-labeling achieves excellent performance with extremely high precision. Fig. 2b shows that as the number of epochs increases, PLG will select more and more highly-confident instance, eventually occupying most of the dataset, while the precision only drops slightly in the final stage, which maintains high precision and high selected ratio. Meanwhile, Fig. 2c shows that NLE identifies more and more negative labels with extremely high precision, which maintains above 0.99 throughout training.

Compared with SOTA semi-supervised learning methods, PLNL performs even better. From Fig. 3a and Fig. 3b, we can see that PLNL achieves both higher selected ratio and recovered more negative labels compared with Fixmatch and Freematch. PLNL is both accurate and comprehensive in recovering label information. This highlights PLNL’s enhanced capacity in leveraging moderately-confident and under-confident instances for label recovery, showcasing both stability and scalability.

PLNL can reduce the generalization bound and achieves lower generalization bound compared with SSL methods. Fig. 3c demonstrates that $(1 - \frac{1-\epsilon}{K-\bar{s}})$ decreases as the training progresses. As $(1 - \frac{1-\epsilon}{K-\bar{s}})$ is the variable of the generalization bound derived in Theorem 1, it is safe to conclude that PLNL can continuously and significantly reduce the generalization error.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

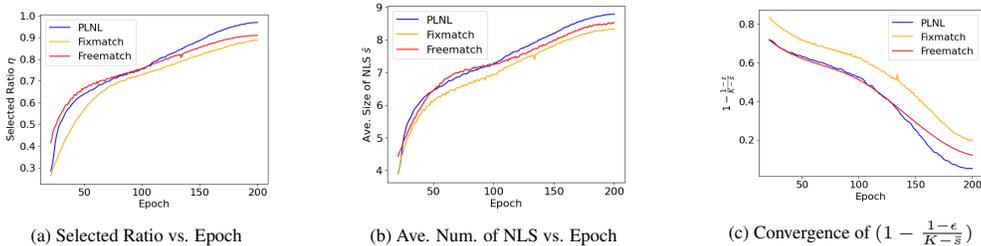


Figure 3: The experiments is conducted on CIFAR-10 with single complementary labels (SCLL). (a) shows that selected ratio of PLNL transcends Fixmatch and Freematch significantly. (b) shows that average size of NLS of PLNL is significantly larger due to specially designed technique NLE for enhancing the untrustworthy negative labels. Nearly all negative labels are revealed at the end of training, almost reaching 9 negative labels for each instances in CIFAR-10 (c) indicates that the value of $(1 - \frac{1-\epsilon}{K-s})$ decreases steadily during training.

6.3 ABLATION STUDY

Two-view networks facilitate increased pseudo-labeling precision. We observe that two-view network significantly boost the performance of PLNL pseudo-labeling, which helps accurately select more highly-confident instances. As shown in Table 3, the η increases **8.75%** and **8.65%** respectively on CIFAR-10 and CIFAR-100 with $1 - \epsilon_1$ increases **6.58%** and **9.22%** respectively. We also compare PLNL with one variant: PLNL v1 where we replace the two-view networks in PLNL with a single network in Table 4, which shows that ours outperforms PLNL v1 by a remarkable margin (e.g. **+6.00%** on STL-10).

All the components contribute to the performance gain. We explore the effectiveness of our proposed PLG and NLE method on three settings STL-10 (SCLL), CIFAR-10 (SCLL), CIFAR-100 (MCLL). Specifically, we compare PLNL with several variants: (1) PLNL v2 which removes the PLG component; (2) PLNL v3 which removes the NLE component. From Table 4, we observe that PLNL outperforms PLNL v2 remarkably (e.g., **+5.43%** on STL-10), which proves the effectiveness of positive label guessing. We also observe that ours outperforms PLNL v3 (e.g., **+1.19%** on CIFAR-100), which proves the effectiveness of NLE.

7 CONCLUSION.

In this paper, we introduce a novel complementary label learning method that reformulates CLL as an inverse problem to infer the full label information from the output space information. To this end, we split this inverse problem into two subtasks (PLG and NLE). A confidence-based instances selection module is proposed for dataset split: highly-confident, moderately-confident and under-confident. Then we perform PLG for highly-confident instances by assigning pseudo-labels to them. For moderately-confident and under-confident instances, we perform NLE by enhancing their negative label set with different levels and train them with the augmented negative labels iteratively. We theoretically prove that when pseudo-labeling error is limited, we can construct a classifier consistent with that learned by clean full labels. The upper bounds of PLG and NLE error rate are deduced and we empirically show that PLNL can infer both positive and negative labels with a high precision. We conducted extensive experiments which demonstrate that PLNL achieves a new state-of-the-art in CLL. In addition, extensive ablation studies have proved the effectiveness of each component.

Table 3: The performance of PLNL with single network on two settings.

Method	CIFAR-10 SCLL		CIFAR-100 MCLL	
	η	$1-\epsilon_1$	η	$1-\epsilon_1$
Single	84.65	90.21	67.43	70.62
Two-view	93.40	96.79	76.08	79.84

Table 4: Classification accuracy of degenerated methods on three settings.

Method	STL-10 SCLL	CIFAR-10 SCLL	CIFAR-100 MCLL
PLNL	55.25	94.78	64.33
PLNL v1	49.25	93.75	63.09
PLNL v2	49.82	92.01	58.94
PLNL v3	53.22	94.28	63.14
POCR	34.96	94.15	53.16

REFERENCES

- 540
541
542 Shoshana Abramovich, Graham Jameson, and Gord Sinnamon. Refining Jensen’s inequality. *Bulletin*
543 *mathématique de la Société des Sciences Mathématiques de Roumanie*, pages 3–14, 2004.
- 544 Yu-Ting Chou, Gang Niu, Hsuan-Tien Lin, and Masashi Sugiyama. Unbiased risk estimators can
545 mislead: A case study of learning with complementary labels. In *International conference on*
546 *machine learning*, pages 1929–1938. PMLR, 2020.
- 547 Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised
548 feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence*
549 *and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011.
- 550
551 Timothee Cour, Ben Sapp, and Ben Taskar. Learning from partial labels. *The Journal of Machine*
552 *Learning Research*, 12:1501–1536, 2011.
- 553 Cynthia Dwork. Differential privacy: A survey of results. In *International conference on theory and*
554 *applications of models of computation*, pages 1–19. Springer, 2008.
- 555
556 Lei Feng and Bo An. Partial label learning with self-guided retraining. In *Proceedings of the AAAI*
557 *conference on artificial intelligence*, volume 33, pages 3542–3549, 2019.
- 558
559 Lei Feng, Takuo Kaneko, Bo Han, Gang Niu, Bo An, and Masashi Sugiyama. Learning with multiple
560 complementary labels. In *International conference on machine learning*, pages 3072–3081. PMLR,
561 2020.
- 562 Xiuwen Gong, Dong Yuan, Wei Bao, and Fulin Luo. A unifying probabilistic framework for partially
563 labeled data learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- 564
565 Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi
566 Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels.
567 *Advances in neural information processing systems*, 31, 2018.
- 568
569 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual net-
570 works. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands,*
October 11–14, 2016, Proceedings, Part IV 14, pages 630–645. Springer, 2016.
- 571
572 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
573 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on*
computer vision and pattern recognition, pages 9729–9738, 2020.
- 574
575 Shuo He, Chaojie Wang, Guowu Yang, and Lei Feng. Candidate label set pruning: A data-centric
576 perspective for deep partial-label learning. In *The Twelfth International Conference on Learning*
577 *Representations*, 2024. URL <https://openreview.net/forum?id=Fk5IzauJ7F>.
- 578
579 Jintao Huang and Yiu-Ming Cheung. Trustworthy partial label learning with out-of-distribution
detection. *arXiv preprint arXiv:2403.06681*, 2024.
- 580
581 Zhizhong Huang, Junping Zhang, and Hongming Shan. Twin contrastive learning with noisy labels.
582 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
11661–11670, 2023.
- 583
584 Takashi Ishida, Gang Niu, Weihua Hu, and Masashi Sugiyama. Learning from complementary labels.
585 *Advances in neural information processing systems*, 30, 2017.
- 586
587 Takashi Ishida, Gang Niu, Aditya Menon, and Masashi Sugiyama. Complementary-label learning for
588 arbitrary losses and models. In *International conference on machine learning*, pages 2971–2980.
PMLR, 2019.
- 589
590 Hiroki Ishiguro, Takashi Ishida, and Masashi Sugiyama. Learning from noisy complementary labels
591 with robust loss functions. *IEICE TRANSACTIONS on Information and Systems*, 105(2):364–376,
592 2022.
- 593
Haoran Jiang, Zhihao Sun, and Yingjie Tian. Comco: Complementary supervised contrastive learning
for complementary label learning. *Neural Networks*, 169:44–56, 2024.

- 594 Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. *IEEE*
595 *Transactions on Big Data*, 7(3):535–547, 2019.
- 596
- 597 Takuo Kaneko, Issei Sato, and Masashi Sugiyama. Online multiclass classification based on prediction
598 margin for partial feedback. *arXiv preprint arXiv:1902.01056*, 2019.
- 599
- 600 Ryuichi Kiryo, Gang Niu, Marthinus C Du Plessis, and Masashi Sugiyama. Positive-unlabeled
601 learning with non-negative risk estimator. *Advances in neural information processing systems*, 30,
602 2017.
- 603 Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.
604 Technical Report 0, University of Toronto, Toronto, Ontario, 2009. URL <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
- 605
- 606 Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- 607
- 608 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image
609 pre-training with frozen image encoders and large language models, 2023a. URL <https://arxiv.org/abs/2301.12597>.
- 610
- 611 Siyuan Li, Weiyang Jin, Zedong Wang, Fang Wu, Zicheng Liu, Cheng Tan, and Stan Z Li. Semireward:
612 A general reward model for semi-supervised learning. *arXiv preprint arXiv:2310.03013*, 2023b.
- 613
- 614 Jiabin Liu, Biao Li, Minglong Lei, and Yong Shi. Self-supervised knowledge distillation for
615 complementary label learning. *Neural Networks*, 155:318–327, 2022.
- 616
- 617 Jiaqi Lv, Miao Xu, Lei Feng, Gang Niu, Xin Geng, and Masashi Sugiyama. Progressive identification
618 of true labels for partial-label learning. In *International conference on machine learning*, pages
619 6500–6510. PMLR, 2020.
- 620
- 621 Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT
622 press, 2018.
- 623
- 624 Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning with
625 noisy labels. *Advances in neural information processing systems*, 26, 2013.
- 626
- 627 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al.
628 Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep
learning and unsupervised feature learning*, volume 2011, page 7. Granada, Spain, 2011.
- 629
- 630 Gang Niu, Marthinus Christoffel Du Plessis, Tomoya Sakai, Yao Ma, and Masashi Sugiyama. Theo-
631 retical comparisons of positive-unlabeled learning against positive-negative learning. *Advances in
neural information processing systems*, 29, 2016.
- 632
- 633 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
634 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward
635 Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,
636 Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning
637 library, 2019.
- 638
- 639 Mina Rezaei, Haojin Yang, and Christoph Meinel. Recurrent generative adversarial network for
640 learning imbalanced medical image semantic segmentation. *Multimedia Tools and Applications*,
79(21):15329–15348, 2020.
- 641
- 642 Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Do-
643 gus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning
644 with consistency and confidence. *Advances in neural information processing systems*, 33:596–608,
645 2020.
- 646
- 647 Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy
labels with deep neural networks: A survey. *IEEE transactions on neural networks and learning
systems*, 2022.

- 648 William Tapper, Gustavo Carneiro, Christos Mikropoulos, Spencer A Thomas, Philip M Evans, and
649 Stergios Boussios. The application of radiomics and ai to molecular imaging for prostate cancer.
650 *Journal of personalized medicine*, 14(3):287, 2024.
- 651 Shiyu Tian, Hongxin Wei, Yiqun Wang, and Lei Feng. Crosel: Cross selection of confident pseudo
652 labels for partial-label learning, 2024.
- 653
- 654 Jesper E Van Engelen and Holger H Hoos. A survey on semi-supervised learning. *Machine learning*,
655 109(2):373–440, 2020.
- 656
- 657 Deng-Bao Wang, Lei Feng, and Min-Ling Zhang. Learning from complementary labels via partial-
658 output consistency regularization. In *IJCAI*, pages 3075–3081, 2021.
- 659 Yidong Wang, Hao Chen, Qiang Heng, Wenxin Hou, Yue Fan, Zhen Wu, Jindong Wang, Marios
660 Savvides, Takahiro Shinozaki, Bhiksha Raj, et al. Freematch: Self-adaptive thresholding for
661 semi-supervised learning. *arXiv preprint arXiv:2205.07246*, 2022.
- 662
- 663 Hongxin Wei, Lei Feng, Xiangyu Chen, and Bo An. Combating noisy labels by agreement: A joint
664 training method with co-regularization. In *Proceedings of the IEEE/CVF conference on computer
665 vision and pattern recognition*, pages 13726–13735, 2020.
- 666
- 667 Shiyu Xia, Jiaqi Lv, Ning Xu, Gang Niu, and Xin Geng. Towards effective visual representations
668 for partial-label learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
669 Pattern Recognition*, pages 15589–15598, 2023.
- 670
- 671 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
672 machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 673
- 674 Kouzhiqiang Yucheng Xie, Jing Wang, Yuheng Jia, Boyu Shi, and Xin Geng. Rankmatch: A novel
675 approach to semi-supervised label distribution learning leveraging inter-label correlations. *arXiv
676 preprint arXiv:2312.06343*, 2023.
- 677
- 678 Ming-Kun Xie and Sheng-Jun Huang. Partial multi-label learning. In *Proceedings of the AAAI
679 conference on artificial intelligence*, volume 32, 2018.
- 680
- 681 Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and Quoc Le. Unsupervised data augmentation
682 for consistency training. *Advances in neural information processing systems*, 33:6256–6268, 2020.
- 683
- 684 Xiangli Yang, Zixing Song, Irwin King, and Zenglin Xu. A survey on deep semi-supervised learning.
685 *IEEE Transactions on Knowledge and Data Engineering*, 35(9):8934–8954, 2022.
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701

A THE PSEUDO-CODE OF PLNL

Algorithm 1 Pseudo-code of PLNL.

Input: Training dataset \mathcal{D} , mini-batch size B , epochs E_{max} , a two-view network with shared parameter Θ and two memory bank matrices M^w, M^s , hyperparameters t, k, τ_m, τ_u .

- 1: Initialize M^w, M^s, Θ by warming up 20 epochs using SCL-LOG (Chou et al., 2020).
- 2: **for** $e = 1, 2, \dots, E_{max}$ **do**
- 3: Shuffle \mathcal{D} into $\frac{|\mathcal{D}|}{B}$ mini-batches;
- 4: Construct M^w, M^s ;
- 5: Compute two-view network outputs of each instances;
- 6: Select $\mathcal{H}, \mathcal{M}, \mathcal{U}$ based on criteria in Eq. (5);
- 7: Pseudo-label the highly-confident set \mathcal{H} ; // PLG
- 8: Enhance negative labels of set \mathcal{M}, \mathcal{U} based on Eq. (14) and Eq. (15); // NLE
- 9: **for** $i = 1, 2, \dots, \frac{|\mathcal{D}|}{B}$ **do**
- 10: Fetch an batch \hat{B}_i with enhanced negative label set;
- 11: Compute the empirical risk $\hat{R}'(f)$ by Eq. (19);
- 12: Update parameters of Θ ;
- 13: **end for**
- 14: **end for**

Output: parameters of Θ .

B PROOF OF THEOREM 1

We first derive the uniform deviation bound between $\hat{R}(f)$ and $R(f)$.

Lemma 1. *Suppose that the binary loss function $\ell(f(\mathbf{x}), y)$ is ρ -Lipschitz continuous w.r.t. $f(\mathbf{x})$. For any $\delta > 0$, with probability at least $1 - \delta$, we have*

$$\left| R(f) - \hat{R}(f) \right| \leq 2\rho K \mathfrak{R}_N(\mathcal{F}) + KB \sqrt{\frac{\log \frac{2}{\delta}}{2N}}. \quad (23)$$

Proof. We firstly define the Rademacher complexity of \mathcal{L} and \mathcal{F} with N training instances as follows:

$$\begin{aligned} & \mathfrak{R}_N(\mathcal{L} \circ \mathcal{F}) \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{y}, \sigma} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^N \sigma_i \mathcal{L}_{CLL}(f(\mathbf{x}_i), \hat{Y}_i) \right]. \end{aligned} \quad (24)$$

Considering that $\mathcal{L}_{CLL}(f(\mathbf{x}), \hat{Y}_i) = \sum_{y \notin \hat{Y}_i} \ell(f(\mathbf{x}), y)$, we have

$$\begin{aligned} \mathfrak{R}_N(\mathcal{L} \circ \mathcal{F}) &\leq K \mathfrak{R}_N(\ell \circ \mathcal{F}), \\ &\leq \rho K \mathfrak{R}_N(\mathcal{F}), \end{aligned} \quad (25)$$

where the second line is based on Lipschitz continuity of $\ell(f(\mathbf{x}), y)$.

Suppose an instance (\mathbf{x}_i, y_i) is replaced by another arbitrary instance (\mathbf{x}'_i, y'_i) , this leads to a change of $\sup_{f \in \mathcal{F}} R(f) - \hat{R}(f)$ no greater than $\frac{KB}{N}$ due to the fact that ℓ is bounded by B . According to McDiarmid's inequality (Mohri et al., 2018), for any $\delta > 0$, with probability at least $1 - \frac{\delta}{2}$, we have

$$\sup_{f \in \mathcal{F}} R(f) - \hat{R}(f) \leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} R(f) - \hat{R}(f) \right] + KB \sqrt{\frac{\log \frac{2}{\delta}}{2N}}. \quad (26)$$

In addition, it is routine (Mohri et al., 2018) that

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} R(f) - \hat{R}(f) \right] \leq 2\bar{\mathfrak{R}}_N(\mathcal{F}). \quad (27)$$

By combining Eq. (26) and Eq. (27), and further taking the other direction of $\sup_{f \in \mathcal{F}} \widehat{R}(f) - R(f)$ into account, with probability at least $1 - \delta$, we have

$$\sup_{f \in \mathcal{F}} |R(f) - \widehat{R}(f)| \leq 2\rho K \mathfrak{R}_N(\mathcal{F}) + KB \sqrt{\frac{\log \frac{2}{\delta}}{2N}}, \quad (28)$$

which concludes the proof.

Due to inevitable error made in PLG and NLE. The problem of CLL actually has been translated into a noisy CLL problem where the true label might be mislabeled as complementary label (Ishiguro et al., 2022). Let ϵ_1 be the error rate of PLG and ϵ_2 be the error rate of NLE. Since NLE error rate $\epsilon_1 = \frac{\sum_{i=1}^{N_h} \mathbb{I}(y_i \in \widehat{Y}_i)}{N_h}$ and NLE error rate $\epsilon_2 = \frac{\sum_{i=1}^{N_m+N_u} \mathbb{I}(y_i \in \widehat{Y}_i)}{N_m+N_u}$, the actual noise rate $\epsilon = \frac{\sum_{i=1}^N \mathbb{I}(y_i \in \widehat{Y}_i)}{N}$ can be calculated as $\epsilon = \frac{N_h}{N} \epsilon_1 + \frac{N_m+N_u}{N} \epsilon_2 = \eta \epsilon_1 + (1 - \eta) \epsilon_2$. We further bound the difference between $\widehat{R}(f)$ and $\widehat{R}'(f)$.

Lemma 2. Suppose that the binary loss function ℓ is bounded by B . For some noise rate $\epsilon \in (0, 1)$ and average complementary label size \bar{s} for any $f \in \mathcal{F}$, we have

$$|\widehat{R}'(f) - \widehat{R}(f)| \leq (1 - \frac{1 - \epsilon}{K - \bar{s}})B. \quad (29)$$

We firstly proved the upper bound of the $\widehat{R}'(f)$:

$$\begin{aligned} \widehat{R}'(f) &= \frac{1}{N} \sum_{i=1}^N \bar{\mathcal{L}}_{CLL}(f(\mathbf{x}_i), \widehat{Y}_i), \\ &= \frac{1}{N} \sum_{i=1}^N \ell(f(\mathbf{x}_i), y_i) + \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \notin \widehat{Y}_i) \left[\sum_{c \notin \widehat{Y}_i, c \neq y_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), c) - \frac{K - |\widehat{Y}_i| - 1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), y_i) \right] \\ &\quad + \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \in \widehat{Y}_i) \left[\sum_{c \notin \widehat{Y}_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), c) - \ell(f(\mathbf{x}_i), y_i) \right], \\ &\leq \widehat{R}(f) + \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \notin \widehat{Y}_i) \sum_{c \notin \widehat{Y}_i, c \neq y_i} \frac{1}{K - \bar{s}} \ell(f(\mathbf{x}_i), c) + \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \in \widehat{Y}_i) \sum_{c \notin \widehat{Y}_i} \frac{1}{K - \bar{s}} \ell(f(\mathbf{x}_i), c), \\ &\leq \widehat{R}(f) + (1 - \epsilon) \frac{K - \bar{s} - 1}{K - \bar{s}} B + \epsilon B, \end{aligned} \quad (30)$$

where the second line holds based on $\epsilon = \frac{\sum_{i=1}^N \mathbb{I}(y_i \in \widehat{Y}_i)}{N}$ and Jensen's inequality (Abramovich et al., 2004) as , and we can prove the lower bound in a similar way:

$$\begin{aligned} \widehat{R}'(f) &= \frac{1}{N} \sum_{i=1}^N \ell(f(\mathbf{x}_i), y_i) + \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \notin \widehat{Y}_i) \left[\sum_{c \notin \widehat{Y}_i, c \neq y_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), c) - \frac{K - |\widehat{Y}_i| - 1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), y_i) \right] \\ &\quad + \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \in \widehat{Y}_i) \left[\sum_{c \notin \widehat{Y}_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), c) - \ell(f(\mathbf{x}_i), y_i) \right], \\ &\geq \widehat{R}(f) - \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \notin \widehat{Y}_i) \sum_{c \notin \widehat{Y}_i, c \neq y_i} \frac{1}{K - \bar{s}} \ell(f(\mathbf{x}_i), c) - \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \in \widehat{Y}_i) \sum_{c \notin \widehat{Y}_i} \frac{1}{K - \bar{s}} \ell(f(\mathbf{x}_i), c), \\ &\geq \widehat{R}(f) - (1 - \epsilon) \frac{K - \bar{s} - 1}{K - \bar{s}} B - \epsilon B, \end{aligned} \quad (31)$$

By combining these two sides, we have:

$$|\widehat{R}'_h(f) - \widehat{R}_h(f)| \leq 1 - \frac{1 - \epsilon}{K - \bar{s}} B, \quad (32)$$

810 which concludes the proof.

811 Then, for any $\delta > 0$, with probability at least $1 - \delta$, we have

$$\begin{aligned}
812 & R(\hat{f}) \\
813 & \leq \widehat{R}(\hat{f}) + 2\rho K \mathfrak{R}_N(\mathcal{F}) + KB \sqrt{\frac{\log \frac{2}{\delta}}{2N}}, \\
814 & \leq \widehat{R}'(\hat{f}) + \left(1 - \frac{1 - \epsilon}{K - \bar{s}}\right)B + 2\rho K \mathfrak{R}_N(\mathcal{F}) + KB \sqrt{\frac{\log \frac{2}{\delta}}{2N}}, \\
815 & \leq \widehat{R}'(f^*) + \left(1 - \frac{1 - \epsilon}{K - \bar{s}}\right)B + 2\rho K \mathfrak{R}_N(\mathcal{F}) + KB \sqrt{\frac{\log \frac{2}{\delta}}{2N}}, \\
816 & \leq \widehat{R}(f^*) + 2\left(1 - \frac{1 - \epsilon}{K - \bar{s}}\right)B + 2\rho K \mathfrak{R}_N(\mathcal{F}) + KB \sqrt{\frac{\log \frac{2}{\delta}}{2N}}, \\
817 & \leq R(f^*) + 2\left(1 - \frac{1 - \epsilon}{K - \bar{s}}\right)B + 4\rho K \mathfrak{R}_N(\mathcal{F}) + 2KB \sqrt{\frac{\log \frac{2}{\delta}}{2N}}, \tag{33}
\end{aligned}$$

818 where the first and fifth line are based on Lemma 1, the second and fourth line are based on Lemma 2,
819 the third line is based on the definition of the empirical risk minimizer $\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \widehat{R}(f)$ which
820 means any other $f \neq \hat{f}$ would lead to a larger risk of $\widehat{R}(f)$.

821 C PROOF OF THEOREM 2

822 Motivated by the formulation of partially labeled data learning in (Gong et al., 2022), we assume
823 that the full label information $Y = \{y, \tilde{Y}\}$ where y is the ground-truth label and \tilde{Y} is a set of the rest
824 labels. And we assume that Q is a hypothetical predictive model with parameters θ , where θ is the
825 model parameter used for prediction. PLG aims to identify the ground-truth label from the label set
826 except complementary label set, which can be implemented by maximizing the conditional likelihood
827 of training dataset with respect to parameters θ . The conditional log-likelihood given all training
828 examples can be expressed as follows.

$$829 f(\theta|\tilde{Y}) = \frac{1}{N} \sum_{i=1}^N \log Q(y_i|\mathbf{x}_i). \tag{34}$$

830 By multiplying and dividing classifier Q by the true distribution of identified ground-truth labels
831 given features $P(y|\mathbf{x})$, we can formulate Eq. (34) as follows.

$$832 f(\theta|\tilde{Y}) = \frac{1}{N} \sum_{i=1}^N \log \frac{Q(y_i|\mathbf{x}_i)}{P(y_i|\mathbf{x}_i)} + \frac{1}{N} \sum_{i=1}^N \log P(y_i|\mathbf{x}_i). \tag{35}$$

833 By multiplying and dividing the probability $P(Y|\mathbf{x})$ to the second term, we can formulate Eq. (34)
834 as follows.

$$835 f(\theta|\tilde{Y}) = \frac{1}{N} \sum_{i=1}^N \log \frac{Q(y_i|\mathbf{x}_i)}{P(y_i|\mathbf{x}_i)} + \frac{1}{N} \sum_{i=1}^N \log \frac{P(y_i|\mathbf{x}_i)}{P(Y_i|\mathbf{x}_i)} + \frac{1}{N} \sum_{i=1}^N \log P(Y_i|\mathbf{x}_i). \tag{36}$$

836 We use $\mathbb{E}_{(\mathcal{X}, \mathcal{Y})}$ operator to calculate the expectation of the random variables $(\mathcal{X}, \mathcal{Y})$, meaning $n \rightarrow \infty$.

$$837 f(\theta|\tilde{Y}) = \mathbb{E}_{(\mathcal{X}, \mathcal{Y})} \left\{ \log \frac{Q(y|\mathbf{x})}{P(y|\mathbf{x})} \right\} + \mathbb{E}_{(\mathcal{X}, \mathcal{Y})} \left\{ \log \frac{P(y|\mathbf{x})}{P(Y|\mathbf{x})} \right\} + \mathbb{E}_{(\mathcal{X}, \mathcal{Y})} \{ \log P(Y|\mathbf{x}) \}. \tag{37}$$

Recall that we assume that the full label information $Y = \{y, \tilde{Y}\}$. Then the second term of Eq. (37) can be developed as follows.

$$\begin{aligned}
\mathbb{E}_{(\mathcal{X}, \mathcal{Y})} \left\{ \log \frac{P(y|\mathbf{x})}{P(Y|\mathbf{x})} \right\} &= -\mathbb{E}_{(\mathcal{X}, \mathcal{Y})} \left\{ \log \frac{P(Y|\mathbf{x})}{P(y|\mathbf{x})} \right\}, \\
&= -\sum_{(\mathbf{x}, Y)} P(\mathbf{x}, Y) \log \frac{P(y, \tilde{Y}|\mathbf{x})}{P(y|\mathbf{x})}, \\
&= -\sum_{(\mathbf{x}, Y)} P(\mathbf{x}, Y) \log \frac{P(y, \tilde{Y}, \mathbf{x})}{P(y, \mathbf{x})}, \\
&= -\sum_{(\mathbf{x}, Y)} P(\mathbf{x}, Y) \log \frac{P(\tilde{Y}, \mathbf{x}|y)}{P(\mathbf{x}|y)}, \\
&= -\sum_{(\mathbf{x}, Y)} P(\mathbf{x}, Y) \log \frac{P(\tilde{Y}, \mathbf{x}|y) P(\tilde{Y}|y)}{P(\mathbf{x}|y) P(\tilde{Y}|y)}, \\
&= -\sum_{(\mathbf{x}, Y)} P(\mathbf{x}, Y) \log \frac{P(\tilde{Y}, \mathbf{x}|y)}{P(\mathbf{x}|y) P(\tilde{Y}|y)} - \sum_{(\mathbf{x}, Y)} P(\mathbf{x}, Y) \log P(\tilde{Y}|y), \\
&= -\sum_{(\mathbf{x}, Y)} P(\mathbf{x}, Y) \log \frac{P(\tilde{Y}, \mathbf{x}|y)}{P(\mathbf{x}|y) P(\tilde{Y}|y)} - \sum_{(\tilde{Y}, y)} P(\tilde{Y}, y) \log P(\tilde{Y}|y), \\
&= -I(\tilde{Y}, X|y) + H(\tilde{Y}|y), \\
&= -I(\tilde{Y}, X|y). \tag{38}
\end{aligned}$$

The last equality holds because the conditional entropy $H(\tilde{Y}|y) = 0$. This is because in CLL setting, once y is known, then the full label information is of course known in advance, meanwhile, thus the uncertainty remaining in \tilde{Y} is zero, i.e., $H(\tilde{Y}|y) = 0$. By combining Eq. (38) and Eq. (37), we have the objective function as follows.

$$f(\theta|\tilde{Y}) = -\mathbb{E}_{(\mathcal{X}, \mathcal{Y})} \left\{ \log \frac{P(y|\mathbf{x})}{Q(y|\mathbf{x})} \right\} - I(\tilde{Y}, X|y) - H(Y|X). \tag{39}$$

The first term is a log likelihood ratio between the true and the predicted ground-truth label distributions given features. The value of this term depends on how well the model Q can approximate P . The second term is the conditional mutual information between the complementary labels and the features, given the ground-truth label. The last term is a constant independent of parameters.

Then, we discuss the mild assumption under which PLG method is effective for CLL.

Assumption 2. Let $y \in Y'$ and $y' \in Y'$ denote any ground-truth label and unidentified negative label respectively. Let X denote the random variables of \mathbf{x} . Let $I(y, X)$ denote the mutual information between ground-truth label y and feature X . Let $I(y', X)$ denote the mutual information between any unidentified negative label y' and feature X . Then, with probability no more than ψ , we have

$$I(y, X) \leq I(y', X), \tag{40}$$

where $\psi < \frac{1}{K-1-s}$.

Remark. This assumption ensures that the feature tends to have more mutual information with positive labels than negative labels.

Motivated by the simplification for identification method in (Gong et al., 2022). In PLG, the training objective actually is the second term of Eq. (39): $f_y(\mathbf{x}) = I(y, X|\tilde{Y})$. As a result, the error rate of PLG can be calculated as follows. Suppose \hat{y} is the guessed positive label, and y' is any negative

918 label.

$$\begin{aligned}
919 \quad \mathbb{P}(y \neq j, \hat{y} = j) &= \sum_{y' \notin \tilde{Y}} P(f_{\hat{y}}(\mathbf{x}) - f_{y'}(\mathbf{x}) \leq 0), \\
920 &= \sum_{y' \notin \tilde{Y}} P(I(\hat{y}, X|\tilde{Y}) \leq I(y', X|\tilde{Y})), \\
921 &\leq (K - 1 - s)\psi, \\
922 & \\
923 & \\
924 & \tag{41}
\end{aligned}$$

925 where the second line is based on replacing \tilde{Y} with any y in 39, the fourth line is based on assumption
926 2. This concludes the proof.

927 D PROOF OF THEOREM 3

928 It is evident that reliable representation information is of great importance to the performance of k -NN
929 based NLE. Motivated by the label distinguishability setting in (He et al., 2024), a mild assumption
930 for CLL datasets are discussed to ensure the reliability of the representation information.

931 According to assumption 1, we have the following lemma.

932 **Lemma 3.** $\forall (\mathbf{x}_i, \bar{Y}_i) \in \mathcal{D}$, let p denote the probability of the true label $y_i \in Y_i$ appearing in its
933 k -NN instance's complementary label set. Let q denote the probability of each non-complementary
934 negative label $y \in Y_i \setminus \{y_i\}$ appearing in its k -NN instance's complementary label set. Then, we have

$$935 \quad p \leq \alpha_k, q \geq \beta_k. \tag{42}$$

936 Then, We derive the error bound of NLE in a step by step manner as follows.

$$\begin{aligned}
937 \quad \mathbb{P}(\mathbf{F}_i^{(\tau)} < \mathbf{F}_{iy}) &= \sum_{j=0}^k P(\mathbf{F}_i^{(\tau)} < \mathbf{F}_{iy} | \mathbf{F}_{iy} = j) P(\mathbf{F}_{iy} = j) \\
938 &= \sum_{j=1}^k P(\mathbf{F}_i \leq j - 1) P(\mathbf{F}_i = j), \\
939 &= \sum_{j=1}^k \left(\underbrace{P(\mathbf{F}_i^{(1)} > j) \dots P(\mathbf{F}_i^{(\tau)} < j)}_{(\tau-1)\text{items}} \dots \underbrace{P(\mathbf{F}_i^{(|Y_i|-1)} < j)}_{(|Y_i|-\tau)\text{items}} \right) P(\mathbf{F}_i = j), \\
940 &\leq \sum_{j=1}^k \binom{|Y_i| - 1}{|Y_i| - \tau} F_{\beta_k}(k - j + 1, j)^{(|Y_i| - \tau)} b_{\alpha_k}(k, j), \\
941 & \tag{43}
\end{aligned}$$

942 where $F_{\beta_k}(k, j) = \int_0^{\beta_k} p^{k-1} (1-p)^{j-1} dt$ denotes the regularized incomplete beta function.
943 $b_{\alpha_k}(k, j) = \binom{k}{j} \alpha_k^j (1 - \alpha_k)^{k-j}$ is simply the probability mass function of a binomial distribution
944 $B(k, \alpha_k)$ which is used to describe the counting process of the label frequency \mathbf{F}_{iy} .

945 E DETAILS OF MCLL SETTINGS.

946 Ishida *et al.* (Ishida et al., 2017) assumed that $\bar{p}(\mathbf{x}, \bar{y})$ is expressed as:

$$947 \quad \bar{p}(\mathbf{x}, \bar{y}) = \frac{1}{c-1} \sum_{y \neq \bar{y}} p(\mathbf{x}, y). \tag{44}$$

948 This assumption implies that all other labels except the correct label are picked as the complementary
949 label with uniform probabilities. Later, Feng *et al.* (Feng et al., 2020) considered a more general
950 setting where each instance is associated with multiple complementary labels (Multiple CLL).
951 Suppose a Multiple CLL dataset is represented as $\{(\mathbf{x}_i, \bar{Y}_i)\}_{i=1}^N$, where \bar{Y}_i is the complementary
952 label set for instance \mathbf{x}_i . Let us denote the number of the complementary labels by a random variable
953 s_i , which is sampled from a distribution $p(s_i)$. Then, we assume that each sample is drawn from the
954 following distribution:

$$955 \quad \bar{p}(\mathbf{x}_i, \bar{Y}_i) = \sum_{j=1}^{c-1} \bar{p}(\mathbf{x}_i, \bar{Y}_i | s_i = j) p(s_i = j), \tag{45}$$

972 where

$$973 \quad \bar{p}(\mathbf{x}_i, \bar{Y}_i | s_i = j) p(s_i = j) := \left\{ \begin{array}{ll} \frac{1}{\binom{c-1}{j}} \sum_{y \notin \bar{Y}_i} p(\mathbf{x}, y), & \text{if } s_i = j, \\ 0, & \text{otherwise.} \end{array} \right\} \quad (46)$$

976 F MORE EXPERIMENT DETAILS

977 Details of Compared Methods.

- 980 • UB-EXP (Feng et al., 2020), an unbiased risk estimator with an estimation error bound, which is derived for Multiple CLL specially.
- 981 • UB-LOG (Feng et al., 2020), another unbiased risk estimator with an estimation error bound but with a different multi-class classification loss function.
- 982 • SCL-EXP (Chou et al., 2020), a surrogate complementary loss with the use of exponential loss function.
- 983 • SCL-LOG (Chou et al., 2020), a surrogate complementary loss with the use of negative log loss function.
- 984 • POCR (Wang et al., 2021), an algorithm which combines the SCL-LOG loss and the consistency regularization technique.
- 985 • SELF-CL (Liu et al., 2022), a self-supervised learning algorithm which integrates self-distillation to CLL.
- 986 • ComCo (Jiang et al., 2024), a contrastive learning framework which leverages the contrastive learning technique on CLL.

987 Details of Implementation.

988 **Implementation.** Values of hyperparameters in PLNL are set as follows. The queue size t is selected from $\{2, 3, 4, 5\}$, k -NN parameter k is selected from $\{100, 250, 500\}$. The α in the instance-aware self-adaptive threshold is selected from $\{0.1, 0.3, 0.5, 0.7, 0.9, 0.99\}$. For each method, we train the commonly used PreAct-ResNet18 (He et al., 2016) with 200 epochs (initial 20 epochs for warm-up), and use SGD as the optimizer with a momentum of 0.9, a weight decay of $1e-4$. We set the batch size from $\{64, 128\}$, the initial learning rate from $\{10^{-1}, 10^{-2}\}$, and we use cosine learning rate scheduling with final learning rate 10^{-3} . We employed *faiss* (Johnson et al., 2019) to compute k -NN instances in the output space, which is a library for efficient similarity search and clustering of dense vectors. For weak augmentations, we employ normalization, horizontal flipping and random cropping. For strong augmentations, we use RandAugment strategy for all, which selects the type and magnitude of augmentation based on uniform probability.

1000 For implementation of SSL methods, we firstly pre-train the network with complete CLL dataset for 200 epochs. Next, we perform pseudo-labeling iteratively and train the model for another 200 epochs.

1001 All of the experiments are implemented based on PyTorch (Paszke et al., 2019) and all of our experiments are conducted with 8 NVIDIA 4090 GPUs.

1002 G MORE EXPERIMENTAL RESULTS.

1003 G.1 QUANTATIVE RESULTS: PERFORMANCE OF PLG AND NLE.

1004 As shown in Table 5, PLG achieves high selected ratio in all datasets, even reaching 98.86% on FMNIST (MCLL), and 76.08% on the difficult CIFAR-100. And the precision remains high even at a high selected ratio. NLE also has high precision, reaching 80.84% even on CIFAR-100.

1005 G.2 PARAMETER SENSITIVITY ANALYSIS.

1006 **Influence of memory bank size t .** As is shown in Table 6, it is obvious that t has little effect on the experimental results on the three datasets. We chose $t = 5$, which has a slightly better effect. It is important to utilize historical confidence information to alleviate confirmation bias, but the size of the memory banks does not matter.

Table 5: The performance of PLG and NLE on three settings. The results (mean \pm std) are reported over 3 random trials. η denotes PLG selected ratio, $1 - \epsilon_1$ denotes PLG precision, $1 - \epsilon_2$ denotes NLE precision.

Dataset	Case	η	Performance	
			$1 - \epsilon_1$	$1 - \epsilon_2$
CIFAR-10	SCLL	93.40 \pm 0.14%	96.79 \pm 0.08%	97.79 \pm 0.09%
	MCLL	97.49 \pm 0.12%	99.30 \pm 0.05%	98.91 \pm 0.07%
FMNIST	SCLL	97.25 \pm 0.08%	95.24 \pm 0.11%	97.56 \pm 0.08%
	MCLL	98.86 \pm 0.06%	97.89 \pm 0.11%	98.97 \pm 0.03%
CIFAR-100	MCLL	76.08 \pm 0.22%	79.84 \pm 0.41%	80.84 \pm 0.04%

Table 6: Classification accuracy of PLNL with different memory bank size t on three benchmark datasets. The best results are highlighted in bold and the second best are underlined (The same applies hereinafter).

t	STL-10 SCLL	CIFAR-10 SCLL	CIFAR-100 MCLL
2	54.45%	94.12%	<u>64.21%</u>
3	<u>54.95%</u>	93.98%	63.95%
4	54.88%	<u>94.56%</u>	63.81%
5	55.25%	94.78%	64.33%

Influence of confidence threshold λ . From Table 8, we observe that there is a trade-off between η and $1 - \epsilon_1$, i.e., a higher threshold will lead the precision to increase but result in less reliable instances selected while a lower threshold will decrease the precision but select more reliable instances. Instance-aware self-adaptive threshold (IST) shows obvious performance gain compared with a fixed global threshold, showcasing its effectiveness.

Influence of k -NN parameter k . As is shown in Table 7, the selection of k depends on the specific situation of the dataset. For example, STL-10 has only 500 labeled instances for each category. If k is set too large, there will be instances that are not in the category near the decision boundary, which will induce more noise in NLE and cause performance degradation.

Table 7: Classification accuracy of PLNL with different k -NN parameter k on three benchmark datasets.

k	STL-10 SCLL	CIFAR-10 SCLL	CIFAR-100 MCLL
100	<u>54.45%</u>	94.25%	<u>63.27%</u>
250	55.25%	<u>94.73%</u>	64.33%
500	47.23%	94.78%	59.65%

Table 8: Classification accuracy and PLG performance of PLNL with different confidence threshold λ on two benchmark datasets. η denotes PLG selected ratio, $1 - \epsilon_1$ denotes PLG precision, $1 - \epsilon_2$ denotes NLE precision. IST denotes instance-aware self-adaptive threshold.

λ	CIFAR-10, SCLL			CIFAR-100, MCLL		
	Accuracy	η	$1 - \epsilon_1$	Accuracy	η	$1 - \epsilon_1$
0.85	92.89%	99.86%	90.79%	62.12%	81.52%	72.88%
0.90	93.85%	95.78%	92.43%	63.95%	77.25%	77.34%
0.95	<u>94.74%</u>	87.28%	98.68%	<u>64.12%</u>	73.48%	83.61%
IST ($\alpha = 0.5$)	94.78%	93.40%	96.79%	64.33%	76.08%	79.84%

G.3 EXPERIMENTAL RESULTS ON TINY IMAGENET

Tiny-ImageNet (Le and Yang, 2015) contains 100000 images of 200 classes. Each class has 500 training images, 50 validation images and 50 test images. Due to its huge number of categories, it is an extremely difficult dataset for complementary label learning. Most existing CLL methods have only tested their performance on 10-class small datasets. Most of their backbones are ResNet and

Table 9: Comparison of classification accuracies between different methods on Tiny-ImageNet with multiple complementary labels per instance.

Method	Tiny-ImageNet
UB-EXP	3.89%
UB-LOG	7.17%
SCL-EXP	3.36%
SCL-LOG	8.96%
POCR	4.29%
SELF-CL	7.87%
ComCo	8.52%
Ours	11.87%

a single complementary label is a rather difficult setting for large datasets, these will lead to poor performance of traditional CLL methods.

However, we have tested the performance of our method on Tiny-ImageNet with MCLL settings. Though most of the methods perform poorly, our method still outperforms traditional CLL methods obtrusively.

G.4 EXTRACT FEATURES BASED ON VISUAL LANGUAGE MODELS (VLMS) AND SELF-SUPERVISED LEARNING (SSL) TECHNIQUES.

In the paper, we compute k -NN instances based on the model output space information, namely the feature extracted by the model itself. It strikes us that different features may have much to do with the NLE performance. To confirm this, we employ different feature extractors for computing k -NN, including PreActResNet-18-MoCo, BLIP-2 (Li et al., 2023a). Note that these results are just for performance comparison which has nothing to do with the results presented in the main body of the paper. For MoCo, we train a PreActResNet by self-supervised learning method MoCo (He et al., 2020) without any supervision. The weak and strong data augmentations used in MoCo follow the original configurations mentioned in the main body. Then we compute k -NN on the 512-dimensional feature output of the PreActResNet. For BLIP-2, we first employ the visual encoder to extract 768-dimensional high-quality representations and then leverage *faiss* (Johnson et al., 2019) to compute k -NN instances in this feature space. We compute the average precision of NLE $1 - \epsilon_2$ and accuracy on CIFAR-10 (SCLL) and CIFAR-100 (MCLL). As shown in Table 10, the feature extracted from BLIP-2 outperforms MoCo and ResNet itself significantly. This shows the powerful visual representation ability of VLMS, which has a great potential for facilitating innovation in weakly-supervised learning in the future.

Table 10: Comparison of classification accuracies and NLE precision between different methods on CIFAR-100 with multiple complementary labels per instance. PreActResNet-18 denotes leveraging the model output space information for k -NN calculation

Feature Extractor	CIFAR-10, SCLL		CIFAR-100, MCLL	
	Accuracy	$1 - \epsilon_2$	Accuracy	$1 - \epsilon_2$
MoCo	93.12%	95.64%	61.82%	75.34%
BLIP-2	95.84%	99.91%	69.85%	93.34%
PreActResNet-18	94.78%	97.79%	64.33%	80.84%

G.5 COMPARISON WITH SEMIREWARD

We further compare the performance of PLNL with one of the recently published SSL method SemiReward Li et al. (2023b). It can be seen that PLNL still outperforms SemiReward significantly in selection ratio and average size of NLS.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

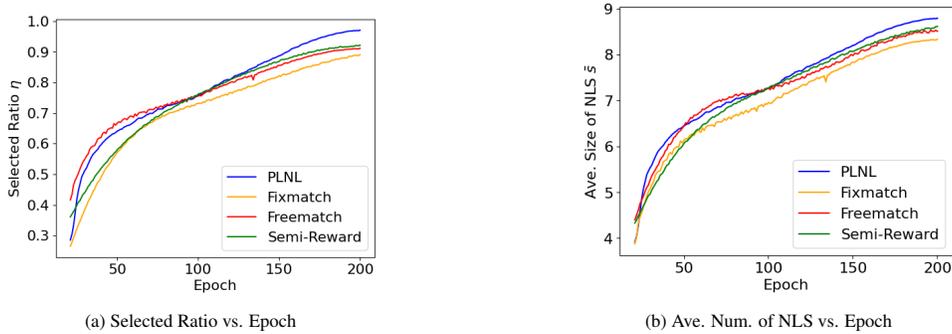


Figure 4: The experiments is conducted on CIFAR-10 with single complementary labels (SCLL). (a) shows that selected ratio of PLNL transcends Fixmatch, Freematch and SemiReward significantly. (b) shows that average size of NLS of PLNL is significantly larger due to specially designed technique NLE for enhancing the untrustworthy negative labels. Nearly all negative labels are revealed at the end of training, almost reaching 9 negative labels for each instances in CIFAR-10.

H DEVIATION REPORTS OF TABLE 3 AND TABLE 4

Table 11: The performance of PLNL with single network on two settings.

Method	CIFAR-10 SCLL		CIFAR-100 MCLL	
	η	$1-\epsilon_1$	η	$1-\epsilon_1$
Single	$84.65 \pm 0.12\%$	$90.21 \pm 0.09\%$	$67.43 \pm 0.24\%$	$70.62 \pm 0.31\%$
Two-view	$93.40 \pm 0.14\%$	$96.79 \pm 0.08\%$	$76.08 \pm 0.22\%$	$79.84 \pm 0.41\%$

Table 12: Classification accuracy of degenerated methods on three settings.

Method	STL-10 SCLL	CIFAR-10 SCLL	CIFAR-100 MCLL
PLNL	$55.25 \pm 0.35\%$	$94.78 \pm 0.12\%$	$64.33 \pm 0.43\%$
PLNL v1	$49.25 \pm 0.41\%$	$93.75 \pm 0.08\%$	$63.09 \pm 0.27\%$
PLNL v2	$49.82 \pm 0.39\%$	$92.01 \pm 0.12\%$	$58.94 \pm 0.31\%$
PLNL v3	$53.22 \pm 0.35\%$	$94.28 \pm 0.10\%$	$63.14 \pm 0.32\%$
POCR	$34.96 \pm 0.32\%$	$94.15 \pm 0.09\%$	$53.16 \pm 0.11\%$

I DETAILED DERIVATION OF EQUATION 30

$$\begin{aligned}\widehat{R}'(f) &= \frac{1}{N} \sum_{i=1}^N \widehat{\mathcal{L}}_{CLL}(f(\mathbf{x}_i), \widehat{Y}_i) \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{y \notin \widehat{Y}_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), y).\end{aligned}\quad (47)$$

This step is due to the definition of CLL loss function:

$$\mathcal{L}_{CLL}(f(\mathbf{x}), \widehat{Y}_i) = \sum_{y \notin \widehat{Y}_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}), y).$$

Then we have:

$$\begin{aligned}\widehat{R}'(f) &= \frac{1}{N} \sum_{i=1}^N \ell(f(\mathbf{x}_i), y_i) + \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \notin \widehat{Y}_i) \left[\sum_{c \notin \widehat{Y}_i, c \neq y_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), c) - \frac{K - |\widehat{Y}_i| - 1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), y_i) \right] \\ &\quad + \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \in \widehat{Y}_i) \left[\sum_{c \notin \widehat{Y}_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), c) - \ell(f(\mathbf{x}_i), y_i) \right].\end{aligned}\quad (48)$$

The first term is the loss on the ground-truth label, summing up to the empirical risk $\widehat{R}(f)$. The second and third term is the difference between empirical risk $\widehat{R}(f)$ and practical empirical risk with pseudo-labeling error $\widehat{R}'(f)$.

Since pseudo-labeling error may occur during PLG and NLE processes, the ground-truth label of \mathbf{x}_i may be mistakenly included or correctly excluded in the enhanced negative label set \widehat{Y}_i . We discuss two situations separately (i.e. $\mathbb{I}(y_i \in \widehat{Y}_i)$ and $\mathbb{I}(y_i \notin \widehat{Y}_i)$). We extract the empirical risk term $\widehat{R}(f)$ by subtracting the equivalent value in the second and third terms, which is shown between square brackets above.

As $\ell(f(\mathbf{x}), y)$ is bounded by a positive value B . We can scale up the second and third terms by directly removing the terms after minus sign, which is double underlined above, which will yield the following line.

$$\begin{aligned}\widehat{R}'(f) &\leq \widehat{R}(f) + \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \notin \widehat{Y}_i) \sum_{c \notin \widehat{Y}_i, c \neq y_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), c) \\ &\quad + \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \in \widehat{Y}_i) \sum_{c \notin \widehat{Y}_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), c).\end{aligned}\quad (49)$$

Then we utilize Jensen's Inequality for further scaling up.

Jensen's Inequality for concave function φ :

$$\varphi \left(\sum_{i=1}^n g(x_i) \lambda_i \right) \geq \sum_{i=1}^n \varphi(g(x_i)) \lambda_i,$$

where $\lambda_1 + \lambda_2 + \dots + \lambda_n = 1$, $\lambda_i \geq 0$.

Here, the concave function is $\varphi(|\widehat{Y}_i|) = \frac{1}{K - |\widehat{Y}_i|}$, where $\lambda_i = \frac{1}{N}$, since $\sum_{c \notin \widehat{Y}_i, c \neq y_i} \ell(f(\mathbf{x}_i), c)$ indicates the sum of binary losses on non-complementary labels $c \notin \widehat{Y}_i$, excluding the ground-truth

label $c \neq y_i$. The number of non-complementary labels is $K - |\widehat{Y}_i|$, excluding the ground-truth label will get $K - |\widehat{Y}_i| - 1$, which means computing the binary loss $\ell(f(\mathbf{x}_i), c)$ for $K - |\widehat{Y}_i| - 1$ times. Finally, since ℓ is bounded by B , we can make such a scale:

$$\begin{aligned}
& \frac{1}{N} \sum_{i=1}^N \sum_{c \notin \widehat{Y}_i, c \neq y_i} \frac{1}{K - |\widehat{Y}_i|} \ell(f(\mathbf{x}_i), c) \\
& \leq \frac{1}{N} \sum_{i=1}^N \frac{K - |\widehat{Y}_i| - 1}{K - |\widehat{Y}_i|} B \\
& \leq \frac{K - \frac{1}{N} \sum_{i=1}^N |\widehat{Y}_i| - 1}{K - \frac{1}{N} \sum_{i=1}^N |\widehat{Y}_i|} B.
\end{aligned} \tag{50}$$

Since we have defined that $\bar{s} = \frac{1}{N} \sum_{i=1}^N |\widehat{Y}_i|$ and $\mathbb{I}(y_i \notin \widehat{Y}_i) = 1 - \epsilon$, we have

$$\begin{aligned}
& \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \notin \widehat{Y}_i) \frac{K - \frac{1}{N} \sum_{i=1}^N |\widehat{Y}_i| - 1}{K - \frac{1}{N} \sum_{i=1}^N |\widehat{Y}_i|} B \\
& \leq (1 - \epsilon) \frac{K - \bar{s} - 1}{K - \bar{s}} B.
\end{aligned} \tag{51}$$

The third term in Eq.(48) follows similar scaling procedures.