

ACCESSIBLE, REALISTIC, AND FAIR EVALUATION OF POSITIVE-UNLABELED LEARNING ALGORITHMS

Wei Wang^{1,2*†} Dong-Dong Wu^{2,1*} Ming Li^{3,2}
 Jingxiong Zhang² Gang Niu¹ Masashi Sugiyama^{1,2}

¹ RIKEN, Tokyo 103-0027, Japan

² The University of Tokyo, Chiba 277-8561, Japan

³ University of Maryland, College Park, MD 20742, USA

wei.wang@riken.jp

{dongdongwu1230, lm1640362161, gang.niu.ml}@gmail.com

7656472907@edu.k.u-tokyo.ac.jp sugi@k.u-tokyo.ac.jp

ABSTRACT

Positive-unlabeled (PU) learning is a weakly supervised binary classification problem, in which the goal is to learn a binary classifier from only positive and unlabeled data, without access to negative data. In recent years, many PU learning algorithms have been developed to improve model performance. However, experimental settings are highly inconsistent, making it difficult to identify which algorithm performs better. In this paper, we propose the first PU learning benchmark to systematically compare PU learning algorithms. During our implementation, we identify subtle yet critical factors that affect the realistic and fair evaluation of PU learning algorithms. On the one hand, many PU learning algorithms rely on a validation set that includes negative data for model selection. This is unrealistic in traditional PU learning settings, where no negative data are available. To handle this problem, we systematically investigate model selection criteria for PU learning. On the other hand, PU learning involves different problem settings and corresponding solution families, i.e., the one-sample and two-sample settings. However, existing evaluation protocols are heavily biased towards the one-sample setting and neglect the significant difference between them. We identify the internal label shift problem of unlabeled training data for the one-sample setting and propose a simple yet effective calibration approach to ensure fair comparisons within and across families. We hope our framework will provide an accessible, realistic, and fair environment for evaluating PU learning algorithms in the future.

1 INTRODUCTION

In binary classification, both positive and negative data are usually necessary to train an effective classifier. However, in many real-world applications, collecting negative data can be more challenging than collecting positive data (Hsieh et al., 2015; Zhou et al., 2021). In positive-unlabeled (PU) learning, only positive and unlabeled data are needed. The objective is to train a binary classifier that assigns positive or negative labels to unseen instances. Therefore, PU learning is a promising weakly supervised binary classification approach for many real-world problems where negative data are difficult to obtain, including recommender systems (Yi et al., 2017; Chen et al., 2023), anomaly detection (Ju et al., 2020; Tian et al., 2024; Takahashi et al., 2025), knowledge graphs (Yin et al., 2024), and link prediction (Wu et al., 2024; Mao et al., 2025).

In recent years, there has been significant progress in PU learning algorithms. PU learning can be divided into three groups: cost-sensitive PU learning algorithms (du Plessis et al., 2014; Zhao et al., 2022), sample-selection PU learning algorithms (Chen et al., 2020b; Wang et al., 2023a), and biased PU learning algorithms (Teisseyre et al., 2025). Cost-sensitive algorithms assign different weights to positive and unlabeled data to approximate the classification risk. Sample-selection algorithms select high-confidence negative data from unlabeled data, which are then given to supervised learning

*Equal contributions. †Corresponding author: Wei Wang <wei.wang@riken.jp>.

algorithms. Biased PU learning algorithms model the biased generation process of positive data and exploit various correction approaches.

Although many PU learning algorithms have been developed to improve generalization performance, there is a lack of a unified experimental setup in the literature for fairly comparing different PU learning algorithms. The experimental settings of different papers are not consistent with each other, making it difficult to tell which algorithm is better. It has been observed that subtle differences in experimental settings can greatly affect the model performance of PU learning algorithms. Additionally, subtle algorithm details, including data augmentation, algorithm tricks, and warm-up strategies, can also greatly affect model performance (Zhu et al., 2023b; Wang et al., 2023a). Therefore, a unified experimental protocol is necessary to further promote the development of PU learning algorithms. In this paper, we propose the first PU learning benchmark to systematically and fairly compare state-of-the-art PU learning algorithms with unified experimental settings. We propose careful and unified implementations of the data generation, algorithm training, and evaluation processes for PU learning algorithms. This makes it easier for users to validate the effectiveness of their newly developed algorithms.

In our implementations, we observe that many PU learning algorithms rely on a validation set containing both positive and negative data for meta-learning, model selection, or early stopping (Chen et al., 2020b; Zhu et al., 2023b; Long et al., 2024). However, accessing negative data is unrealistic and contradicts the original motivation of PU learning (Elkan & Noto, 2008), which goes against the advantages of PU learning in not depending on negative data. Actually, if we can obtain some negative data, we can directly apply supervised learning techniques, which can greatly boost model performance (Sakai et al., 2017). Therefore, standardizing the composition and use of the validation set is vital to fairly and practically evaluating PU learning algorithms. In this paper, we systematically revisit the model selection criteria for PU learning by using only positive and unlabeled validation data, and validate their effectiveness with both theoretical and empirical analyses.

In addition, there are different families and corresponding solutions of PU learning algorithms, but existing evaluations fail to consider the differences between these families. From the perspective of data generation processes, there are two types of PU learning problems: the one-sample (OS) and two-sample (TS) settings. In the OS setting, the positive and unlabeled training sets are generated sequentially. An unlabeled dataset is first sampled from the marginal density. Then, if an instance in the unlabeled dataset is positive, its positive label is observed with a constant probability. If an instance in the unlabeled dataset is negative, its label is never observed, and the instance remains unlabeled. Finally, the observed positive data constitute the positive training set, while the remaining unlabeled data constitute the unlabeled training set. In the TS setting, the positive and unlabeled training sets are generated independently, meaning that the density of unlabeled training data is the same as the marginal density. This indicates that the density of unlabeled training data is different in these two settings. Figure 1 shows an example of the distribution of unlabeled data under the OS and TS settings. We can find that the class priors of the two settings are different. This causes an internal label shift (ILS) problem for the unlabeled training data when adopting the OS setting as the evaluation setting. Unfortunately, this problem has typically been overlooked. Existing evaluation protocols are heavily biased towards the OS setting and compare OS and TS algorithms together without specific manipulations. This can deteriorate the performance of TS PU learning algorithms and lead to unfair experimental comparisons. Therefore, we identify the ILS problem for the first time in the PU learning literature and propose a simple yet effective calibration approach to overcome it with theoretical guarantees.

We draw the following key takeaways from our benchmark results:

We draw the following key takeaways from our benchmark results:

- No single algorithm outperforms all others on every dataset or evaluation metric; some early, simple methods already achieve strong classification performance. Therefore, we should choose which PU learning algorithm to use on a case-by-case basis.

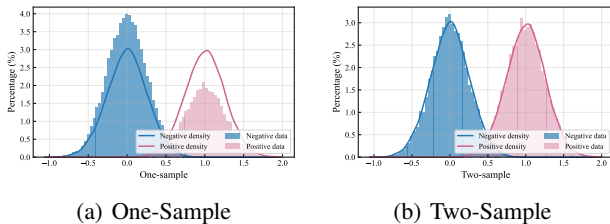


Figure 1: An example of the comparison of the distribution of unlabeled training data in different PU learning settings.

- The model-selection problem in PU learning must be addressed when designing new algorithms or conducting empirical comparisons, and different selection criteria should be used for different test metrics.
- The performance of TS PU learning algorithms degrades significantly when evaluated under the OS setting without proper adaptation, so OS protocols in the existing PU learning literature do not reflect the true performance of TS methods. Hence, differences between OS and TS settings must be considered to ensure fair cross-family comparisons.

2 PRELIMINARIES

In this section, we present the background of PU learning and existing state-of-the-art algorithms.

2.1 POSITIVE-UNLABELED LEARNING

Problem Setting. Let $\mathcal{X} \subseteq \mathbb{R}^d$ denote the d -dimensional feature space and $\mathcal{Y} = \{+1, -1\}$ denote the binary label space. Let $p(\mathbf{x}, y)$ denote the joint probability density over the random variables $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$. In PU learning, we are given a positive training set $D_P = \{(\mathbf{x}_i, +1)\}_{i=1}^{n_P}$ and an unlabeled training set $D_U = \{\mathbf{x}_i\}_{i=n_P+1}^{n_P+n_U}$. Let $\pi = p(y = +1)$ denote the class prior probability of the positive class. Let $p(\mathbf{x}|y = +1)$ and $p(\mathbf{x}|y = -1)$ denote the positive and negative class-conditional densities, respectively. Let $p(\mathbf{x})$ denote the marginal density. The goal of PU learning is to learn a binary classifier $f : \mathcal{X} \rightarrow \mathbb{R}$ from $D_P \cup D_U$ that maximizes the *expected accuracy*

$$\text{ACC}(f) = \mathbb{E}_{p(\mathbf{x}, y)} \mathbb{I}(yf(\mathbf{x}) \geq 0), \quad (1)$$

where \mathbb{E} denotes the expectation and \mathbb{I} denotes the indicator function. However, since the 0-1 loss function is difficult to optimize, we usually use a surrogate loss function ℓ , such as the logistic loss. Then, the *classification risk* to be minimized can be expressed as

$$R(f) = \mathbb{E}_{p(\mathbf{x}, y)} [\ell(f(\mathbf{x}), y)]. \quad (2)$$

Data Generation Assumption. There are mainly two data generation assumptions for PU learning, i.e., the TS setting (du Plessis et al., 2014; Niu et al., 2016; Chen et al., 2020a) and the OS setting (Elkan & Noto, 2008; Coudray et al., 2023). In the TS setting, we assume that D_P and D_U are generated *independently*, where D_P is sampled from the positive conditional density $p(\mathbf{x}|y = +1)$ and D_U is sampled from the marginal density $p(\mathbf{x})$. In the OS setting, D_U and D_P are generated *sequentially*. First, D_U is sampled from the marginal density $p(\mathbf{x})$. Second, for each example in D_U , if it is positive, its positive label is observed with a *constant probability* $c > 0$. If an example is negative, its negative label is never observed and the example remains unlabeled with probability 1. Finally, the observed positive data constitute D_P and all the unlabeled data left constitute D_U .

2.2 POSITIVE-UNLABELED LEARNING ALGORITHMS

From a methodology taxonomy perspective, PU learning algorithms can be divided into three groups: cost-sensitive algorithms, sample-selection algorithms, and biased PU learning algorithms. Cost-sensitive algorithms assign different weights to positive and unlabeled data to approximate the classification risk (du Plessis et al., 2015; Kiryo et al., 2017; Hsieh et al., 2019). Some algorithms are equipped with other regularization techniques to further improve performance, such as entropy minimization (Zhao et al., 2022; Jiang et al., 2023) and mixup technique (Chen et al., 2020a; Li et al., 2022). Sample-selection algorithms select reliable negative examples from the unlabeled dataset for supervised learning (Chen et al., 2020b; Garg et al., 2021; Wang et al., 2023a; Li et al., 2024). Biased PU learning algorithms consider the density of positive data to be biased and adopt different strategies to model the bias (Bekker et al., 2019; Gong et al., 2022; Coudray et al., 2023; Wang et al., 2023b; Teisseyre et al., 2025).

3 MODEL SELECTION FOR POSITIVE-UNLABELED LEARNING

In this section, we first explain our motivation for studying the model selection problem in PU learning. Next, we review the criteria used for model selection in PU learning, including the proxy accuracy, proxy area under the curve score, and oracle accuracy.

3.1 MOTIVATION

Although model selection is well established for supervised learning, it is non-trivial for PU learning because negative data are inaccessible. This problem is particularly important for deep learning algorithms because they have many hyperparameters, including universal hyperparameters (e.g., learning rates and weight decay) and algorithm-specific hyperparameters. Previous work has usually conducted model selection by assuming a validation set with labels (i.e., both positive and negative labels) is available. However, this assumption is inconsistent with the definition of PU learning, in which negative data are unavailable. Therefore, it is important to study the model selection problem systematically for PU learning. According to the original definition of PU learning (Bekker & Davis, 2020), we assume that the validation set consists of a positive validation set $D'_P = \{(\mathbf{x}'_i, +1)\}_{i=1}^{n'_P}$ and an unlabeled validation set $D'_U = \{\mathbf{x}'_i\}_{i=n'_P+1}^{n'_P+n'_U}$.

3.2 PROXY ACCURACY

Although the validation accuracy cannot be directly calculated because of the absence of negative data, it has been shown that the expected accuracy can be expressed using only positive and unlabeled data (du Plessis et al., 2014). This motivates us to apply it for model selection.

Definition 1 (Proxy accuracy (PA)). The proxy accuracy of a binary classifier f on the PU validation dataset is defined as

$$\text{PA}(f) = \begin{cases} \frac{2\pi}{n'_P} \sum_{i=1}^{n'_P} \mathbb{I}(f(\mathbf{x}'_i) \geq 0) + \frac{1}{n'_U} \sum_{i=n'_P+1}^{n'_P+n'_U} \mathbb{I}(f(\mathbf{x}'_i) < 0), & \text{if the setting is TS;} \\ \frac{2\pi}{n'_P} \sum_{i=1}^{n'_P} \mathbb{I}(f(\mathbf{x}'_i) \geq 0) + \frac{1}{n'_P+n'_U} \sum_{i=1}^{n'_P+n'_U} \mathbb{I}(f(\mathbf{x}'_i) < 0), & \text{if the setting is OS.} \end{cases} \quad (3)$$

PA can be calculated using only PU validation data when the class prior π is known or estimated (Ramaswamy et al., 2016; Yao et al., 2022; Zhu et al., 2023a). The following proposition then holds.

Proposition 1. For two classifiers f_1 and f_2 that satisfy $\mathbb{E}[\text{PA}(f_1)] < \mathbb{E}[\text{PA}(f_2)]$, we have $\text{ACC}(f_1) < \text{ACC}(f_2)$.

The proof can be found in Appendix B.1. According to Proposition 1, a classifier with a higher expected value of the proxy accuracy can achieve a higher expected accuracy even when the true labels are inaccessible. This means that when the number of validation data is large, the best model chosen using the PA metric will achieve the highest accuracy in expectation. One limitation of PA is that knowledge of the class prior is necessary. However, knowledge of π is an intrinsic and common issue in PU learning. Addressing this issue is beyond the scope of our paper. In practice, we can estimate it using off-the-shelf estimation methods (Ramaswamy et al., 2016; Garg et al., 2021; Yao et al., 2022), and we can even obtain this knowledge in some real-world applications (Sugiyama et al., 2022).

3.3 PROXY AUC SCORE

It has been shown that the area under the curve (AUC) score can be robust to corrupted labels for binary classification (Charoenphakdee et al., 2019; Wei et al., 2022). Therefore, it is promising to employ it for PU model selection. First, we introduce the expected AUC score as follows:

$$\text{AUC}(f) = \mathbb{E}_{p(\mathbf{x}|y=+1)} \mathbb{E}_{p(\mathbf{x}'|y'=-1)} \left[\mathbb{I}(f(\mathbf{x}) > f(\mathbf{x}')) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}) = f(\mathbf{x}')) \right]. \quad (4)$$

We then consider the unlabeled validation data to be corrupted negative data and calculate the AUC score as follows, which is suitable for both OS and TS settings.

Definition 2 (Proxy AUC score (PAUC)). The proxy AUC of a binary classifier f on the PU validation dataset is defined as

$$\text{PAUC}(f) = \frac{1}{n'_P n'_U} \sum_{i=1}^{n'_P} \sum_{j=n'_P+1}^{n'_P+n'_U} \left(\mathbb{I}(f(\mathbf{x}'_i) > f(\mathbf{x}'_j)) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}'_i) = f(\mathbf{x}'_j)) \right). \quad (5)$$

The following proposition then holds.

Proposition 2. *Under both the OS and TS settings, for two classifiers f_1 and f_2 that satisfy $\mathbb{E}[\text{PAUC}(f_1)] < \mathbb{E}[\text{PAUC}(f_2)]$, we have $\text{AUC}(f_1) < \text{AUC}(f_2)$.*

The proof can be found in Appendix B.2. Proposition 2 shows that a classifier with a higher expected value of the proxy AUC score will achieve a higher expected AUC score, regardless of whether the setting is OS or TS. Therefore, when the number of validation data is large, the model selected with the highest PAUC can also achieve the highest expected value of the AUC score. An advantage is that the class prior π is not necessary when calculating the PAUC.

3.4 ORACLE ACCURACY

Finally, we introduce the oracle accuracy metric if the true labels of unlabeled data are available.

Definition 3 (Oracle accuracy (OA)). The oracle accuracy of a binary classifier f on the PU validation dataset is defined as

$$\text{OA}(f) = \begin{cases} \frac{1}{n'_U} \sum_{i=n'_P+1}^{n'_P+n'_U} \mathbb{I}(y'_i f(\mathbf{x}'_i) \geq 0), & \text{if the setting is TS;} \\ \frac{1}{n'_P+n'_U} \sum_{i=1}^{n'_P+n'_U} \mathbb{I}(y'_i f(\mathbf{x}'_i) \geq 0), & \text{if the setting is OS.} \end{cases} \quad (6)$$

Here, y'_i is the true label of \mathbf{x}'_i .

Notably, the implementations for the OS and TS settings differ slightly, as it is important to ensure that the validation data have the same distribution as the test data. OA is a natural metric for supervised learning. However, due to the absence of negative data, it cannot be calculated in the traditional PU learning setting. Unfortunately, this metric has actually been widely used in the PU learning literature because of a lack of standardized benchmarking. Therefore, this paper only includes the results of OA for comparison. We recommend using PA and PAUC in future PU learning experiments, especially in real-world applications where negative data cannot be obtained.

4 INTERNAL LABEL SHIFT IN POSITIVE-UNLABELED LEARNING

In this section, we first introduce the ILS problem in PU learning. Then, we provide a calibration approach to solve it with both theoretical and empirical analysis.

4.1 PROBLEM STATEMENT

The difference between the OS and TS settings lies in the density of the unlabeled training data. Specifically, the density of the unlabeled training data equals the marginal density in the TS setting but differs from it in the OS setting. We formalize the ILS problem as follows.

Definition 4 (Internal label shift in OS PU learning). In the OS setting, the density of \mathcal{D}_U is $\bar{p}(\mathbf{x}) = \bar{\pi}p(\mathbf{x}|y = +1) + (1 - \bar{\pi})p(\mathbf{x}|y = -1)$, where $\bar{\pi}$ is the class prior under the OS setting. Here, the positive and negative class-conditional densities are the same as those of the test data; however, the class prior is $\bar{\pi} = (1 - c)\pi/(1 - c\pi)$, which differs from π , the class prior of the test data. This mismatch causes an internal label shift between the unlabeled training data and the test data.

Many cost-sensitive PU learning algorithms have been developed for the TS setting. In these algorithms, positive and unlabeled data are assigned different weights to approximate the classification risk (du Plessis et al., 2014; Chen et al., 2020a; Zhao et al., 2022). Because the weights are theoretically derived, small discrepancies in data assumptions can degrade performance. Conversely, sample-selection PU learning algorithms select reliable negative data from \mathcal{D}_U and need not rely strictly on the specific data generation process (Zhu et al., 2023b; Wang et al., 2023a; Li et al., 2024). However, many papers adopt only the OS setting and ignore the distribution mismatch, causing experimental datasets to violate the assumptions of TS approaches.

To demonstrate how ILS affects model performance, we use uPU (du Plessis et al., 2015) as an example in Section 4; it is a representative TS algorithm and underpins many subsequent cost-sensitive methods.¹ Under the TS assumption $\mathcal{D}_U \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})$, du Plessis et al. (2015) proposed the

¹Our analysis and calibration approach can be extended to other TS algorithms as well.

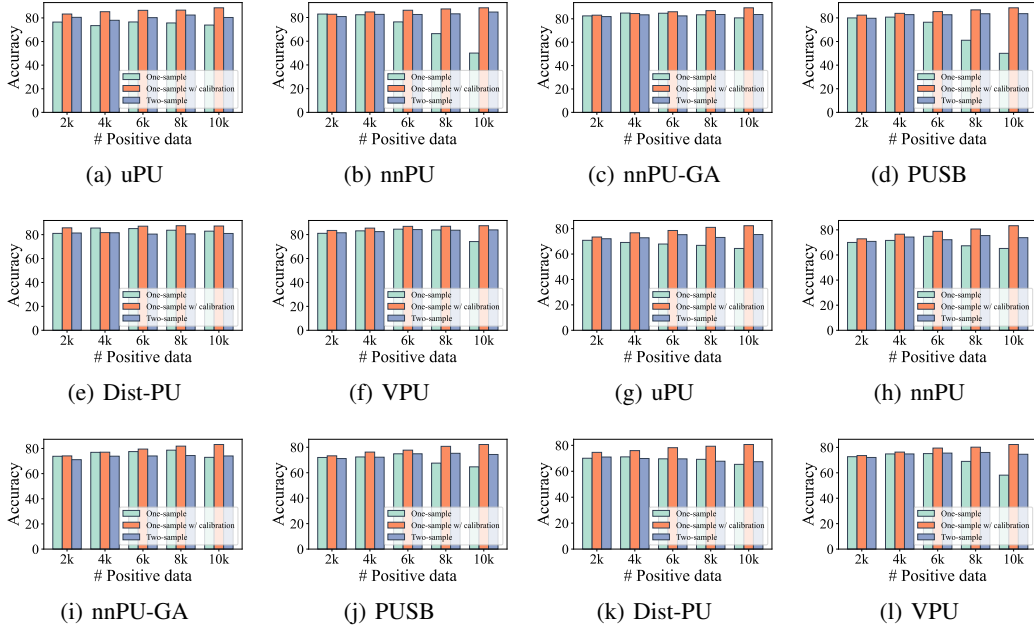


Figure 2: Positive classification accuracies of TS PU learning algorithms in OS and TS settings of a PU version of CIFAR-10 with varying amounts of positive data. Figures (a) to (f) are for Case 1, and Figures (g) to (l) are for Case 2.

unbiased risk estimator (URE)

$$\hat{R}(f) = \frac{\pi}{n_P} \sum_{i=1}^{n_P} (\ell(f(\mathbf{x}_i), +1) - \ell(f(\mathbf{x}_i), -1)) + \frac{1}{n_U} \sum_{i=n_P+1}^{n_P+n_U} \ell(f(\mathbf{x}_i), -1), \quad (7)$$

which enjoys risk consistency because $\mathbb{E}[\hat{R}(f)] = R(f)$. Let $\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}(f)$ and $f^* = \arg \min_{f \in \mathcal{F}} R(f)$ denote the classifiers that minimize the empirical risk in Eq. (7) and the risk in Eq. (2), respectively, where \mathcal{F} is the model class. It is known that $\hat{f} \rightarrow f^*$ as $n_P \rightarrow \infty$ and $n_U \rightarrow \infty$ under the TS setting (Niu et al., 2016). Under the OS setting, however, $\mathbb{E}[\hat{R}(f)] \neq R(f)$, so $\hat{f} \rightarrow f^*$ no longer holds (see Appendix B.3). Consequently, minimizing losses designed for the TS setting may not yield high-performing classifiers when datasets are generated under the OS setting, leading to unfair comparisons when all methods are evaluated in the OS setting. The bias stems from the ILS problem: under the OS setting, the class prior of \mathcal{D}_U differs from π , breaking the consistency of many TS algorithms and degrading their performance.

4.2 THE PROPOSED CALIBRATION APPROACH

To address the bias, we incorporate the true densities of \mathcal{D}_U for TS algorithms. The following theorem shows that the risk rewrite for the uPU approach differs under the OS setting.

Theorem 1. *Under the OS setting, the classification risk in Eq. (2) can be equivalently expressed as*

$$R(f) = \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} [\ell(f(\mathbf{x}), +1) + (c-1)\ell(f(\mathbf{x}), -1)] + (1-c\pi) \mathbb{E}_{\bar{p}(\mathbf{x})} [\ell(f(\mathbf{x}), -1)].$$

The proof is given in Appendix B.4. Theorem 1 shows that the classification risk can be equivalently expressed as expectations w.r.t. the densities of positive and unlabeled data under the OS setting. We then obtain a calibrated risk estimator using the positive and unlabeled datasets:

$$\bar{R}(f) = \frac{\pi}{n_P} \sum_{i=1}^{n_P} (\ell(f(\mathbf{x}_i), +1) + (c-1)\ell(f(\mathbf{x}_i), -1)) + \frac{1-c\pi}{n_U} \sum_{i=n_P+1}^{n_P+n_U} \ell(f(\mathbf{x}_i), -1). \quad (8)$$

When the class prior π is known or estimated, we obtain an unbiased estimate of c as $c = n_P/\pi(n_P + n_U)$. Let $\bar{f} = \arg \min_{f \in \mathcal{F}} \bar{R}(f)$ denote the optimal classifier that minimizes the calibrated risk

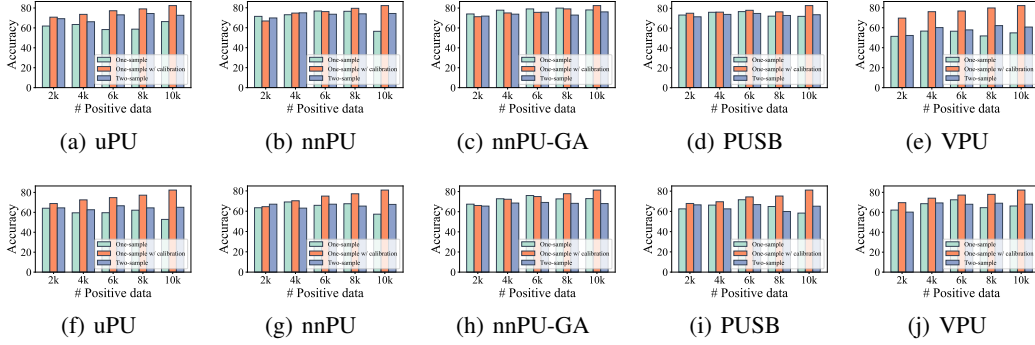


Figure 3: Classification accuracies of TS PU learning algorithms in OS and TS settings of a PU version of ImageNette with varying amounts of positive data. Figures (a) to (e) are for Case 1, and Figures (f) to (j) are for Case 2.

estimator in Eq. (8). Let $\mathfrak{R}_{n_P}(\mathcal{F})$ and $\mathfrak{R}'_{n_U}(\mathcal{F})$ denote the Rademacher complexities defined in Appendix B.5. Then, the following theorem holds.

Theorem 2. Assume that there exists a constant C_f such that $\sup_{f \in \mathcal{F}} \|f\|_\infty \leq C_f$ and a constant C_ℓ such that $\forall y, \sup_{|z| \leq C_f} \ell(z, y) \leq C_\ell$. We also assume that $\forall y$, the binary loss function $\ell(z, y)$ is Lipschitz continuous in z with a Lipschitz constant L_ℓ . For any $\delta > 0$, the following inequality holds with probability at least $1 - \delta$:

$$R(\bar{f}) - R(f^*) \leq (8 - 4c)\pi L_\ell \mathfrak{R}_{n_P}(\mathcal{F}) + (4 - 4c\pi)L_\ell \mathfrak{R}'_{n_U}(\mathcal{F}) + \left(\frac{(4 - 2c)\pi C_\ell}{\sqrt{n_P}} + \frac{(2 - 2c\pi)C_\ell}{\sqrt{n_U}} \right) \sqrt{\frac{\ln 2/\delta}{2}}. \quad (9)$$

The proof is given in Appendix B.5. Theorem 2 shows that $\bar{f} \rightarrow f^*$ as $n_P \rightarrow \infty$ and $n_U \rightarrow \infty$, because $\mathfrak{R}_{n_U, \bar{p}}(\mathcal{F}) \rightarrow 0$ and $\mathfrak{R}_{n_P, p_+}(\mathcal{F}) \rightarrow 0$ for all parametric models with a bounded norm, such as deep neural networks trained with weight decay (Golowich et al., 2018). Notably, Eq. (8) can be equivalently transformed into Eq. (7) if we incorporate \mathcal{D}_P into \mathcal{D}_U when computing the last loss term w.r.t. unlabeled data in Eq. (7) (see Appendix B.6). Thus, when \mathcal{D}_P is used in both loss terms, the ILS bias is eliminated, because the union of positive and unlabeled data is unbiased w.r.t. the marginal density. This motivates a simple yet effective calibration approach that adapts TS algorithms to the OS setting, summarized in Algorithm 1. We augment \mathcal{D}_U with \mathcal{D}_P when computing the loss on unlabeled data, so the replenished set is marginally unbiased and suitable for TS PU learners.

4.3 EMPIRICAL ANALYSIS

We validated the existence of the ILS problem and the effectiveness of the proposed calibration approach. We used uPU (du Plessis et al., 2015), nnPU (Kiryo et al., 2017), nnPU-GA (Kiryo et al., 2017), PUBS (Kato et al., 2019), VPU (Chen et al., 2020a), and Dist-PU (Zhao et al., 2022), six representative TS PU learning algorithms. We used

CIFAR-10 (Krizhevsky & Hinton, 2009) and ImageNette (Deng et al., 2009) as the datasets. We synthesized PU training datasets with different definitions of positive and negative labels, where the details are presented in Appendix C. We did not include the results of Dist-PU on ImageNette since Dist-PU did not work well on this dataset. We considered both the OS and TS cases using the same experimental settings, and the only difference lay in how positive data were generated. Figures 2 and 3 show the experimental results on CIFAR-10 and ImageNette with varying amounts of positive data, respectively. We can observe that using TS approaches directly in the OS setting yields inferior performance. Their performance consistently drops when the number of positive data increases,

Algorithm 1 Calibrated Two-Sample PU Learning

Require: Two-sample PU learning algorithm \mathcal{A} , positive training set \mathcal{D}_P , unlabeled training set \mathcal{D}_U , maximum epochs T_{\max} , maximum iterations I_{\max} .

Ensure: Classifier f produced by \mathcal{A} .

- 1: **for** $t = 1, 2, \dots, T_{\max}$ **do**
 - 2: **Shuffle** \mathcal{D}_P and \mathcal{D}_U ;
 - 3: **for** $k = 1, \dots, I_{\max}$ **do**
 - 4: **Fetch** mini-batch \mathcal{D}_k^P from \mathcal{D}_P and \mathcal{D}_k^U from \mathcal{D}_U ;
 - 5: **Call** $\mathcal{A}.$ TRAIN_ONE_BATCH($\mathcal{D}_k^P, \mathcal{D}_k^U \cup \mathcal{D}_k^P$)
 - 6: **end for**
 - 7: **end for**
-

Table 1: Test results (mean \pm std) of accuracy, AUC, and F1 score for each algorithm on CIFAR-10 (Case 1) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Accuracy			AUC			F1		
	PA	PAUC	OA	PA	PAUC	OA	PA	PAUC	OA
PubN	86.46 \pm 0.46	86.24 \pm 0.84	87.33 \pm 0.28	93.96 \pm 0.49	93.44 \pm 0.74	94.33 \pm 0.23	86.62 \pm 0.52	86.07 \pm 0.98	86.98 \pm 0.24
PAN	76.64 \pm 0.78	77.56 \pm 0.41	78.91 \pm 0.59	87.11 \pm 0.86	87.28 \pm 0.93	85.70 \pm 0.68	79.08 \pm 0.73	79.41 \pm 0.61	78.74 \pm 0.95
CVIR	85.45 \pm 1.03	83.32 \pm 0.44	86.47 \pm 0.48	93.74 \pm 0.73	93.67 \pm 0.62	93.73 \pm 0.31	86.19 \pm 0.88	84.71 \pm 0.33	86.51 \pm 0.40
P3MIX-E	72.68 \pm 6.26	50.00 \pm 0.00	73.96 \pm 5.63	88.80 \pm 2.65	92.62 \pm 0.67	89.56 \pm 2.18	77.65 \pm 3.65	66.67 \pm 0.00	67.45 \pm 12.03
P3MIX-C	86.36 \pm 0.58	85.75 \pm 0.76	86.65 \pm 0.57	92.70 \pm 0.71	93.09 \pm 0.65	93.16 \pm 0.43	86.44 \pm 0.51	85.93 \pm 0.70	86.72 \pm 0.58
LBE	82.71 \pm 0.73	73.60 \pm 1.29	85.03 \pm 0.38	92.09 \pm 0.15	93.21 \pm 0.04	92.26 \pm 0.31	83.79 \pm 0.49	78.72 \pm 0.76	84.31 \pm 0.31
Count Loss	80.89 \pm 0.32	79.86 \pm 0.88	82.39 \pm 0.37	90.63 \pm 0.69	90.40 \pm 0.45	89.20 \pm 1.27	82.60 \pm 0.28	81.83 \pm 0.39	83.11 \pm 0.39
Robust-PU	85.57 \pm 0.18	85.61 \pm 0.55	85.91 \pm 0.35	91.56 \pm 0.49	92.89 \pm 0.29	91.04 \pm 1.60	85.88 \pm 0.09	84.80 \pm 0.96	85.47 \pm 0.32
Holistic-PU	50.20 \pm 0.10	50.00 \pm 0.00	81.81 \pm 0.49	64.56 \pm 11.51	69.45 \pm 5.04	90.60 \pm 0.41	66.64 \pm 0.03	66.67 \pm 0.00	82.97 \pm 0.37
PUc	77.85 \pm 0.85	78.51 \pm 0.33	80.45 \pm 0.46	86.84 \pm 0.61	86.60 \pm 0.45	87.58 \pm 0.44	79.45 \pm 0.55	78.01 \pm 0.48	78.99 \pm 0.28
GLWS	84.46 \pm 0.45	79.83 \pm 2.30	85.66 \pm 0.44	93.55 \pm 0.07	93.54 \pm 0.14	93.48 \pm 0.16	85.65 \pm 0.36	82.69 \pm 1.46	86.26 \pm 0.32
uPU	80.24 \pm 1.25	76.07 \pm 2.83	82.04 \pm 0.49	88.72 \pm 0.40	89.05 \pm 0.17	87.36 \pm 0.73	81.05 \pm 0.90	77.01 \pm 1.41	80.34 \pm 0.56
uPU-c	85.89 \pm 0.44	84.20 \pm 0.49	86.48 \pm 0.21	92.65 \pm 0.38	93.03 \pm 0.22	93.22 \pm 0.15	85.96 \pm 0.43	83.04 \pm 0.92	86.12 \pm 0.10
nnPU	82.03 \pm 0.11	75.56 \pm 0.29	82.40 \pm 0.31	92.62 \pm 0.15	92.32 \pm 0.47	91.95 \pm 0.44	83.51 \pm 0.05	79.64 \pm 0.25	83.49 \pm 0.05
nnPU-c	85.52 \pm 0.20	86.03 \pm 0.68	86.35 \pm 0.26	92.19 \pm 0.33	93.07 \pm 0.55	92.95 \pm 0.38	85.90 \pm 0.28	85.71 \pm 0.70	86.29 \pm 0.30
nnPU-GA	84.26 \pm 0.80	84.18 \pm 0.40	84.93 \pm 0.70	92.79 \pm 0.47	92.26 \pm 0.36	92.25 \pm 0.46	84.87 \pm 0.62	84.63 \pm 0.42	84.58 \pm 0.53
nnPU-GA-c	85.80 \pm 0.29	86.28 \pm 0.31	86.13 \pm 0.25	92.81 \pm 0.42	92.96 \pm 0.47	93.00 \pm 0.42	85.90 \pm 0.31	85.66 \pm 0.27	85.57 \pm 0.19
PUSB	81.53 \pm 0.77	82.49 \pm 1.02	82.91 \pm 0.70	81.53 \pm 0.77	82.49 \pm 1.02	82.91 \pm 0.70	83.29 \pm 0.47	83.80 \pm 0.77	84.12 \pm 0.53
PUSB-c	86.15 \pm 0.37	84.76 \pm 0.17	86.49 \pm 0.17	86.15 \pm 0.37	84.76 \pm 0.17	86.49 \pm 0.17	86.09 \pm 0.44	83.89 \pm 0.19	86.23 \pm 0.18
VPU	84.93 \pm 0.52	65.71 \pm 7.32	85.80 \pm 0.40	91.89 \pm 0.08	92.89 \pm 0.54	92.86 \pm 0.20	84.15 \pm 0.59	42.73 \pm 17.09	84.91 \pm 0.49
VPU-c	86.41 \pm 0.75	82.85 \pm 1.68	87.65 \pm 0.25	92.30 \pm 0.31	93.51 \pm 0.53	91.79 \pm 1.62	86.73 \pm 0.55	84.56 \pm 1.15	87.41 \pm 0.29
Dist-PU	81.64 \pm 0.45	79.31 \pm 0.51	83.56 \pm 0.46	90.91 \pm 0.54	91.90 \pm 0.48	90.59 \pm 0.49	83.34 \pm 0.26	81.94 \pm 0.23	83.26 \pm 0.60
Dist-PU-c	87.06\pm0.45	87.38\pm0.23	88.47\pm0.25	94.93\pm0.31	94.55\pm0.21	94.90\pm0.32	87.63\pm0.33	87.28\pm0.29	88.18\pm0.25

even though we have more knowledge of the true labels of positive data in the unlabeled dataset. By using our proposed calibration approach, the performance can be improved greatly and can even sometimes surpass the performance in the TS setting. This shows the effectiveness of our calibration approach in improving TS approaches under the OS setting.

5 BENCHMARKING POSITIVE-UNLABELED LEARNING

In this section, we first introduce the benchmark settings, then we present the benchmark experimental results. The code package is available at <https://github.com/wu-dd/PUBench>.

5.1 BENCHMARK SETTINGS

We included seventeen representative PU learning algorithms: uPU (du Plessis et al., 2015), nnPU (Kiryo et al., 2017), nnPU-GA (Kiryo et al., 2017), PUSB (Kato et al., 2019), PubN (Hsieh et al., 2019), VPU (Chen et al., 2020a), PAN (Hu et al., 2021), CVIR (Garg et al., 2021), Dist-PU (Zhao et al., 2022), P³MIX-E (Li et al., 2022), P³MIX-C (Li et al., 2022), LBE (Gong et al., 2022), Count Loss (Shukla et al., 2023), Robust-PU (Zhu et al., 2023b), Holistic-PU (Wang et al., 2023a), PUc (Wang et al., 2023b), and GLWS (Chen et al., 2024). We evaluated our methods on two image datasets (CIFAR-10 (Krizhevsky & Hinton, 2009) and ImageNette (Deng et al., 2009)) and two UCI datasets (USPS and Letter) from the UCI Machine Learning Repository (Kelly et al., 2023). ImageNette is a curated subset of the larger ImageNet corpus, containing ten easily distinguishable categories: *tench*, *English springer*, *cassette player*, *chain saw*, *church*, *French horn*, *garbage truck*, *gas pump*, *golf ball*, and *parachute*. We synthesized PU versions of these datasets; detailed information can be found in Appendix C. We used ResNet-34 (He et al., 2016) for image datasets and a multilayer perceptron (MLP) with a hidden layer width of 500 equipped with the ReLU (Nair & Hinton, 2010) activation function for tabular datasets.

Following the widely used validation protocol (Raschka, 2018; Gulrajani & Lopez-Paz, 2021; Wang et al., 2025), we divided some training data from the positive and unlabeled datasets into the positive validation set D'_P and the unlabeled validation set D'_U , respectively. We used various test metrics, including accuracy, AUC score, F1 score, precision, and recall. We first trained a model with training sets D_P and D_U . Then, we evaluated its validation performance based on the metrics in Section 3 as well as its test performance on a test set with true labels. We randomly selected a set of hyperparameter configurations from a given pool. For each validation metric, we selected the checkpoint with the best validation performance on $D'_P \cup D'_U$, and recorded the corresponding test metrics. We recorded the mean test metrics and standard deviations obtained with different data splits.

Table 2: Test results (mean \pm std) of accuracy, AUC, and F1 score for each algorithm on CIFAR-10 (Case 2) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold.

Test metric	Test ACC			AUC			Test F1		
	PA	PAUC	OA	PA	PAUC	OA	PA	PAUC	OA
PubN	78.26 \pm 1.01	79.50\pm0.38	79.94 \pm 0.36	87.81 \pm 0.65	88.00 \pm 0.38	88.08 \pm 0.45	80.47 \pm 0.65	79.17 \pm 0.88	79.91 \pm 0.21
PAN	61.43 \pm 2.74	60.61 \pm 4.34	63.48 \pm 2.71	68.71 \pm 5.63	71.54 \pm 4.68	69.63 \pm 5.43	70.73 \pm 1.40	70.87 \pm 1.72	69.25 \pm 3.04
CVIR	78.49 \pm 1.49	79.50\pm1.46	80.44\pm0.68	88.10\pm0.87	87.98 \pm 1.33	88.68\pm0.81	80.86\pm0.97	80.69\pm1.33	81.44\pm0.58
P3MIX-E	59.04 \pm 4.54	50.00 \pm 0.00	59.13 \pm 4.62	74.26 \pm 4.26	84.52 \pm 0.84	74.11 \pm 4.16	70.45 \pm 2.00	44.44 \pm 18.14	70.45 \pm 2.00
P3MIX-C	78.05 \pm 0.95	77.42 \pm 1.40	78.70 \pm 0.50	85.87 \pm 1.02	84.92 \pm 1.40	86.13 \pm 0.79	79.82 \pm 0.56	79.06 \pm 0.92	79.90 \pm 0.49
LBE	72.47 \pm 1.50	63.54 \pm 2.86	75.96 \pm 0.88	84.02 \pm 0.40	84.26 \pm 0.78	83.47 \pm 0.97	77.13 \pm 0.72	72.96 \pm 1.42	76.04 \pm 0.83
Count Loss	74.44 \pm 0.68	74.75 \pm 0.45	76.87 \pm 0.75	82.88 \pm 1.02	82.99 \pm 1.03	84.44 \pm 0.75	77.41 \pm 0.54	76.70 \pm 0.55	78.27 \pm 0.99
Robust-PU	78.94 \pm 0.79	78.43 \pm 0.61	79.60 \pm 0.81	85.23 \pm 1.09	87.13 \pm 0.76	86.33 \pm 0.63	80.37 \pm 0.72	77.16 \pm 0.68	79.79 \pm 0.89
Holistic-PU	55.60 \pm 0.16	56.04 \pm 4.93	71.18 \pm 1.20	78.03 \pm 2.53	67.96 \pm 6.67	76.93 \pm 3.13	69.02 \pm 0.04	44.49 \pm 18.12	73.64 \pm 2.09
PUe	68.60 \pm 0.41	67.40 \pm 1.90	71.05 \pm 0.52	78.06 \pm 0.31	79.27 \pm 0.51	78.69 \pm 0.36	73.41 \pm 0.44	73.05 \pm 0.71	71.06 \pm 1.35
GLWS	77.71 \pm 0.71	76.22 \pm 1.33	79.58 \pm 0.61	87.86 \pm 0.33	88.08\pm0.43	87.44 \pm 0.51	80.40 \pm 0.37	79.75 \pm 0.81	80.47 \pm 0.47
uPU	66.21 \pm 1.40	69.03 \pm 1.04	70.46 \pm 0.70	76.46 \pm 1.65	78.80 \pm 0.74	77.97 \pm 0.90	71.52 \pm 0.73	72.78 \pm 0.47	70.89 \pm 1.53
uPU-c	77.22 \pm 0.26	79.29 \pm 0.37	79.02 \pm 0.99	85.19 \pm 0.46	87.76 \pm 0.38	87.11 \pm 0.83	79.48 \pm 0.22	77.19 \pm 0.45	78.60 \pm 1.22
nnPU	74.27 \pm 1.26	62.67 \pm 1.09	77.62 \pm 0.68	86.16 \pm 0.07	86.53 \pm 0.16	86.42 \pm 0.58	78.00 \pm 0.55	72.57 \pm 0.51	79.07 \pm 0.52
nnPU-c	77.74 \pm 0.53	78.49 \pm 0.35	79.37 \pm 0.30	84.84 \pm 0.44	86.63 \pm 0.31	86.16 \pm 0.22	79.79 \pm 0.18	77.25 \pm 0.61	79.20 \pm 0.39
nnPU-GA	76.59 \pm 1.15	76.73 \pm 0.88	78.38 \pm 0.74	86.41 \pm 1.24	86.09 \pm 1.23	86.58 \pm 0.84	79.14 \pm 0.95	78.76 \pm 1.11	78.22 \pm 0.53
nnPU-GA-c	78.00 \pm 0.52	78.32 \pm 0.71	79.12 \pm 0.91	83.75 \pm 1.30	85.82 \pm 1.04	85.63 \pm 1.27	79.26 \pm 0.81	77.78 \pm 0.48	79.03 \pm 0.92
PUSB	75.74 \pm 0.61	78.80 \pm 0.55	78.35 \pm 0.41	75.74 \pm 0.61	78.80 \pm 0.55	78.35 \pm 0.41	79.18 \pm 0.43	79.83 \pm 0.59	79.79 \pm 0.61
PUSB-c	79.06\pm0.45	77.98 \pm 0.54	79.19 \pm 0.32	79.06 \pm 0.45	77.98 \pm 0.54	79.19 \pm 0.32	80.06 \pm 0.36	77.43 \pm 0.40	79.29 \pm 0.40
VPU	76.99 \pm 1.00	63.22 \pm 5.30	77.31 \pm 0.86	85.47 \pm 0.98	87.08 \pm 0.43	86.07 \pm 0.67	75.15 \pm 1.31	39.92 \pm 15.74	73.45 \pm 1.33
VPU-c	77.70 \pm 0.41	78.20 \pm 0.90	79.81 \pm 0.66	86.90 \pm 0.39	87.50 \pm 0.46	86.32 \pm 0.28	80.12 \pm 0.27	80.52 \pm 0.53	80.56 \pm 0.71
Dist-PU	73.46 \pm 0.59	74.83 \pm 0.58	74.69 \pm 0.60	80.70 \pm 0.45	82.09 \pm 0.40	81.48 \pm 0.78	76.90 \pm 0.31	76.88 \pm 0.16	76.65 \pm 0.15
Dist-PU-c	72.57 \pm 3.47	74.41 \pm 2.67	74.30 \pm 2.73	80.34 \pm 3.48	82.49 \pm 2.68	81.94 \pm 2.90	75.50 \pm 2.34	75.27 \pm 2.67	73.68 \pm 3.25

Table 3: Test results (mean \pm std) of accuracy, AUC, and F1 score for each algorithm on ImageNet (Case 1) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold.

Test metric	Test ACC			AUC			Test F1		
	PA	PAUC	OA	PA	PAUC	OA	PA	PAUC	OA
PubN	75.69 \pm 0.02	77.07 \pm 0.47	78.99 \pm 0.57	84.82 \pm 0.28	86.25 \pm 0.83	87.45 \pm 0.37	77.63 \pm 0.16	76.30 \pm 1.26	79.25 \pm 0.76
PAN	50.74 \pm 1.17	51.52 \pm 0.46	56.93 \pm 1.83	53.71 \pm 3.39	55.48 \pm 1.03	55.73 \pm 1.79	65.24 \pm 0.19	32.31 \pm 14.81	45.43 \pm 2.65
CVIR	78.78 \pm 0.86	78.26 \pm 1.62	81.01\pm0.67	87.98\pm0.35	88.29\pm0.65	89.28\pm0.38	80.12 \pm 0.36	79.52 \pm 0.78	81.51\pm0.45
P3MIX-E	74.81 \pm 2.36	49.71 \pm 0.48	75.19 \pm 2.39	82.71 \pm 2.74	85.84 \pm 0.68	82.91 \pm 2.92	76.23 \pm 1.97	43.92 \pm 17.93	76.41 \pm 2.04
P3MIX-C	78.81\pm1.61	78.91\pm1.84	80.25 \pm 0.82	85.85 \pm 1.50	86.41 \pm 1.31	87.33 \pm 1.27	80.26\pm1.24	80.35\pm1.27	80.48 \pm 0.37
LBE	78.52 \pm 0.41	78.73 \pm 0.65	79.20 \pm 0.36	86.84 \pm 0.37	86.31 \pm 0.61	86.16 \pm 0.78	78.90 \pm 0.32	77.09 \pm 1.48	78.14 \pm 0.75
Count Loss	74.98 \pm 0.85	75.95 \pm 1.56	78.07 \pm 0.73	85.50 \pm 0.23	85.44 \pm 0.52	85.75 \pm 0.74	77.84 \pm 0.40	77.95 \pm 0.87	78.92 \pm 0.91
Robust-PU	77.67 \pm 0.27	75.53 \pm 2.04	78.73 \pm 0.43	83.93 \pm 0.64	85.22 \pm 0.11	84.46 \pm 0.93	77.86 \pm 0.47	71.78 \pm 4.44	78.06 \pm 0.59
Holistic-PU	51.16 \pm 0.47	54.42 \pm 3.66	53.62 \pm 0.24	58.85 \pm 1.01	56.45 \pm 6.07	55.25 \pm 0.43	65.18 \pm 0.31	64.23 \pm 1.18	51.58 \pm 1.27
PUe	67.47 \pm 1.88	71.46 \pm 1.27	70.90 \pm 1.28	75.35 \pm 1.52	77.29 \pm 1.49	77.47 \pm 1.55	70.39 \pm 0.48	70.97 \pm 1.81	71.46 \pm 1.56
GLWS	76.14 \pm 0.86	74.96 \pm 1.62	78.68 \pm 0.70	87.00 \pm 0.40	86.89 \pm 0.71	86.96 \pm 0.74	78.93 \pm 0.45	78.52 \pm 1.01	79.56 \pm 0.67
uPU	71.07 \pm 0.95	64.14 \pm 6.15	73.69 \pm 0.74	82.24 \pm 0.61	81.60 \pm 1.06	81.94 \pm 0.41	74.88 \pm 0.50	71.58 \pm 2.45	74.95 \pm 0.78
uPU-c	75.00 \pm 0.97	72.54 \pm 4.40	77.76 \pm 0.66	84.13 \pm 0.33	85.82 \pm 0.55	84.65 \pm 0.65	77.16 \pm 0.27	63.33 \pm 9.88	77.19 \pm 0.67
nnPU	75.63 \pm 1.34	66.81 \pm 1.09	77.80 \pm 0.74	86.56 \pm 0.38	86.12 \pm 0.71	86.72 \pm 0.26	78.52 \pm 0.77	73.97 \pm 0.57	78.19 \pm 0.46
nnPU-c	76.51 \pm 0.61	76.95 \pm 0.75	77.66 \pm 0.63	83.87 \pm 0.71	85.08 \pm 0.67	83.93 \pm 1.24	77.89 \pm 0.33	74.59 \pm 1.65	77.33 \pm 0.99
nnPU-GA	75.70 \pm 0.36	78.72 \pm 0.64	79.40 \pm 0.47	83.74 \pm 0.65	86.06 \pm 0.88	84.45 \pm 1.58	78.33 \pm 0.16	78.98 \pm 1.22	79.13 \pm 0.24
nnPU-GA-c	77.65 \pm 0.58	72.91 \pm 2.33	78.56 \pm 0.06	81.45 \pm 1.32	84.75 \pm 0.53	82.69 \pm 1.21	77.88 \pm 0.51	65.42 \pm 5.22	78.34 \pm 0.42
PUSB	72.73 \pm 0.54	77.03 \pm 0.74	76.73 \pm 0.35	73.10 \pm 0.53	77.19 \pm 0.68	76.91 \pm 0.33	77.26 \pm 0.33	78.65 \pm 0.09	78.66 \pm 0.16
PUSB-c	76.37 \pm 0.16	77.36 \pm 0.36	77.81 \pm 0.60	76.48 \pm 0.15	77.31 \pm 0.33	77.86 \pm 0.60	77.37 \pm 0.17	76.42 \pm 0.04	78.15 \pm 0.80
VPU	56.36 \pm 2.98	50.91 \pm 0.03	61.72 \pm 0.41	61.21 \pm 2.22	82.35 \pm 0.27	73.84 \pm 4.69	53.86 \pm 6.31	0.14 \pm 0.11	45.88 \pm 4.26
VPU-c	77.48 \pm 0.83	78.00 \pm 0.50	78.06 \pm 0.91	83.35 \pm 0.33	84.63 \pm 0.49	84.28 \pm 0.96	78.09 \pm 0.69	78.60 \pm 0.40	77.64 \pm 0.69
Dist-PU	70.40 \pm 2.37	71.86 \pm 2.34	74.68 \pm 0.79	83.97 \pm 1.16	83.18 \pm 1.51	83.92 \pm 0.73	75.58 \pm 1.50	75.76 \pm 1.61	77.00 \pm 0.75
Dist-PU-c	72.03 \pm 0.99	65.88 \pm 3.33	73.84 \pm 1.06	79.51 \pm 0.44	77.83 \pm 0.61	80.91 \pm 1.18	74.78 \pm 0.58	52.05 \pm 10.16	74.10 \pm 0.81

5.2 BENCHMARK RESULTS

Tables 1 to 4, and 7 to 18 in Appendix D report detailed experimental results in terms of different metrics on CIFAR-10, ImageNet, Letter, and USPS, and the hyperparameters were determined with PA, PAUC, and OA, respectively. In addition, Figures 4 to 7 show the overall performance of different algorithms. For ease of presentation in figures, we did not include the algorithms where the performance is obviously inferior. We can draw the following conclusions based on the experimental results: 1) The TS algorithms without calibration perform worse due to the ILS problem, indicating the existence of an evaluation pitfall in the literature. The proposed calibration technique consistently improves the classification performance for TS approaches, demonstrating the effectiveness of the proposed calibration technique. 2) There is no algorithm that can win in every case of the dataset and evaluation metric. Besides, some early algorithms can already achieve satisfactory classification performance. 3) Our proposed validation metrics are effective in hyperparameter selection. However, the effectiveness may also depend on the test metric. For example, we can observe from Table 1 that the model selected using PAUC can achieve better performance than using OA when the test metric is the AUC score.

Table 4: Test results (mean±std) of accuracy, AUC, and F1 score for each algorithm on ImageNette (Case 2) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold.

Test metric	Test ACC			AUC			Test F1		
	PA	PAUC	OA	PA	PAUC	OA	PA	PAUC	OA
PubN	75.30±0.58	75.97±0.61	77.39±0.45	83.97±0.64	84.03±0.50	85.20±0.66	76.89±0.51	75.98±1.24	76.44±0.60
PAN	53.31±0.36	64.73±1.69	64.37±2.18	65.28±1.32	70.38±1.92	69.19±2.61	66.49±0.27	58.68±5.06	63.03±2.55
CVIR	76.60±0.74	77.87±0.67	79.29±0.47	85.84±0.26	86.25±0.32	87.22±0.49	78.43±0.48	78.67±0.35	79.39±0.09
P3MIX-E	60.42±4.27	49.86±0.23	60.82±4.16	70.79±2.16	81.61±0.76	71.51±2.47	67.11±1.54	44.19±18.04	67.38±1.42
P3MIX-C	74.17±0.90	75.40±0.81	75.35±0.81	83.92±0.88	83.97±1.13	83.68±0.43	76.85±0.73	77.21±0.65	77.13±0.44
LBE	74.51±0.94	74.67±0.54	76.31±0.92	83.06±1.01	82.33±0.45	83.33±0.99	76.85±0.74	73.81±1.63	74.99±1.41
Count Loss	73.27±0.28	73.62±0.23	74.43±0.66	82.20±0.51	82.04±0.66	82.05±0.72	76.28±0.22	76.11±0.30	76.46±0.45
Robust-PU	72.58±1.19	72.78±0.43	75.52±0.68	80.19±0.81	80.57±0.71	81.69±0.38	73.75±0.62	69.94±1.57	74.06±1.05
Holistic-PU	56.12±0.93	54.70±2.16	59.19±0.53	61.46±0.21	59.83±1.01	62.22±0.69	64.75±1.38	60.81±1.16	58.83±0.57
PUc	64.65±0.59	65.89±1.36	67.63±0.53	72.74±1.48	72.62±1.26	74.27±0.79	69.33±1.07	68.66±0.71	69.42±0.90
GLWS	75.61±0.65	75.38±0.24	76.99±0.21	85.81±0.55	86.55±0.37	85.77±0.29	78.47±0.34	78.40±0.13	78.65±0.19
uPU	60.42±2.82	66.42±1.08	66.29±1.00	67.49±3.09	73.24±0.86	72.82±0.52	67.95±0.68	67.50±1.57	66.46±1.37
uPU-c	72.57±1.56	73.20±1.05	75.07±0.54	79.19±2.25	81.76±0.92	82.22±0.87	72.60±1.24	69.52±2.07	74.58±0.79
nnPU	69.83±0.52	55.99±3.35	72.76±0.55	82.83±1.26	80.53±0.83	82.65±0.76	74.92±0.35	69.09±1.48	75.65±0.45
nnPU-c	74.42±0.75	73.55±1.07	74.17±1.07	82.13±0.89	82.44±0.97	81.77±1.05	75.24±1.25	71.49±2.66	73.68±1.79
nnPU-GA	72.19±1.31	75.23±1.08	75.87±0.60	81.37±0.66	82.62±1.08	83.37±0.99	75.44±0.43	74.24±2.14	74.85±0.80
nnPU-GA-c	72.27±1.25	73.85±0.58	74.62±0.11	79.16±1.52	81.60±0.47	79.72±1.52	73.86±0.71	71.02±0.68	74.79±0.81
PUSB	71.29±1.80	72.57±0.76	75.30±0.56	71.44±1.78	72.65±0.76	75.35±0.54	75.76±0.85	74.73±0.74	76.77±0.29
PUSB-c	72.98±1.11	73.69±0.38	74.86±0.52	73.00±1.10	73.64±0.38	74.86±0.50	73.69±0.94	71.85±0.48	74.70±0.37
VPU	70.42±1.87	58.37±6.43	73.28±0.73	78.68±1.11	78.52±1.50	80.21±0.47	73.29±0.72	24.24±19.50	70.20±1.70
VPU-c	76.30±0.79	77.75±0.57	77.38±1.00	84.37±0.45	84.11±0.49	83.80±0.79	78.34±0.83	78.32±0.40	77.88±1.75
Dist-PU	63.97±1.03	67.74±0.50	68.58±1.04	72.64±0.26	75.26±0.57	74.62±0.96	69.88±0.13	71.92±0.37	70.92±0.75
Dist-PU-c	60.43±4.37	68.02±1.27	67.29±1.74	67.51±3.90	74.10±1.29	71.97±2.69	67.65±0.89	69.06±1.11	65.39±3.35

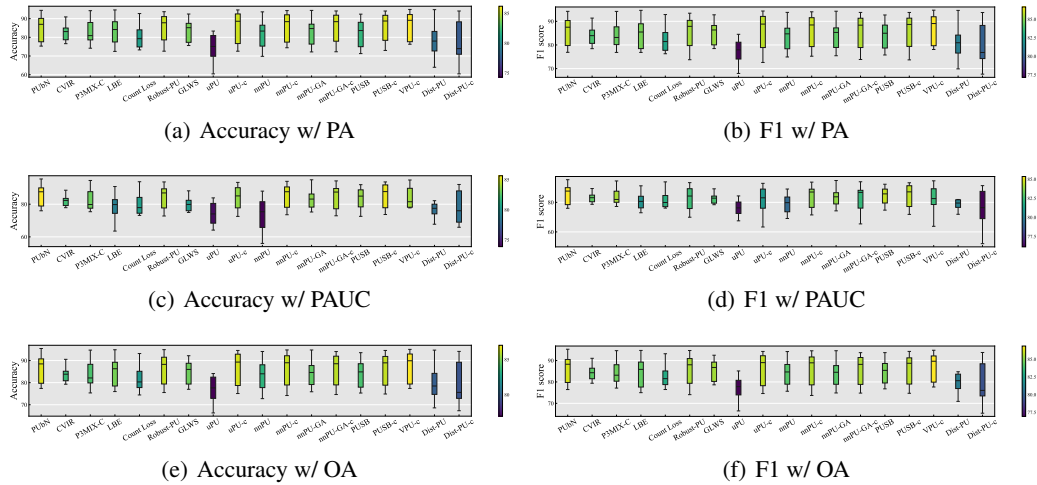


Figure 4: Overall performance w.r.t. accuracy and the F1 score across all datasets. Hyperparameters were tuned using PA, PAUC and OA, respectively; bar colors indicate means.

6 CONCLUSION

In this paper, we conducted a comprehensive empirical study of PU learning algorithms. To the best of our knowledge, we proposed the first PU learning benchmark to systematically compare different PU learning algorithms in a unified framework. We investigated model selection criteria to facilitate realistic evaluation of PU learning algorithms. We also identified the ILS problem for the one-sample setting of PU learning and proposed a calibration approach to ensure fair comparisons of different families of PU learning algorithms. We hope that our framework can facilitate accessible, realistic, and fair evaluation of PU learning algorithms in the future. A limitation of our work is that we use relatively small benchmark datasets following previous work. In the future, it is also promising to investigate the performance of different algorithms on collected large-scale PU benchmark datasets.

ACKNOWLEDGMENTS

WW and DDW were supported by the Junior Research Associate (JRA) program of RIKEN. MS was supported by JST ASPIRE Grant Number JPMJAP2405.

REFERENCES

- Jessa Bekker and Jesse Davis. Learning from positive and unlabeled data: A survey. *Machine Learning*, 109:719–760, 2020.
- Jessa Bekker, Pieter Robberechts, and Jesse Davis. Beyond the selected completely at random assumption for learning from positive and unlabeled data. In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2019, Proceedings, Part II*, volume 11907 of *Lecture Notes in Computer Science*, pp. 71–85, 2019.
- Nontawat Charoenphakdee, Jongyeong Lee, and Masashi Sugiyama. On symmetric losses for learning from corrupted labels. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 961–970, 2019.
- Hao Chen, Jindong Wang, Lei Feng, Xiang Li, Yidong Wang, Xing Xie, Masashi Sugiyama, Rita Singh, and Bhiksha Raj. A general framework for learning from weak supervision. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 7462–7485, 2024.
- Hui Chen, Fangqing Liu, Yin Wang, Liyue Zhao, and Hao Wu. A variational approach for learning from positive and unlabeled data. In *Advances in Neural Information Processing Systems 33*, pp. 14844–14854, 2020a.
- Jiawei Chen, Hande Dong, Xiang Wang, Fuli Feng, Meng Wang, and Xiangnan He. Bias and debias in recommender system: A survey and future directions. *ACM Transactions on Information Systems*, 41(3):1–39, 2023.
- Qiuyi Chen, Haiyang Zhang, Leqi Zhang, Changchun Li, Jia Wang, and Wei Wang. PU-Bench: A unified benchmark for rigorous and reproducible PU learning. In *Proceedings of the 14th International Conference on Learning Representations*, 2026.
- Xuxi Chen, Wuyang Chen, Tianlong Chen, Ye Yuan, Chen Gong, Kewei Chen, and Zhangyang Wang. Self-PU: Self boosted and calibrated positive-unlabeled training. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 1510–1519, 2020b.
- Olivier Coudray, Christine Keribin, Pascal Massart, and Patrick Pamphile. Risk bounds for positive-unlabeled learning under the selected at random assumption. *Journal of Machine Learning Research*, 24(107):1–31, 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009.
- Marthinus du Plessis, Gang Niu, and Masashi Sugiyama. Convex formulation for learning from positive and unlabeled data. In *Proceedings of the 32nd International Conference on Machine Learning*, pp. 1386–1394, 2015.
- Marthinus C. du Plessis, Gang Niu, and Masashi Sugiyama. Analysis of learning from positive and unlabeled data. In *Advances in Neural Information Processing Systems 27*, pp. 703–711, 2014.
- Charles Elkan and Keith Noto. Learning classifiers from only positive and unlabeled data. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 213–220, 2008.
- Saurabh Garg, Yifan Wu, Alexander J. Smola, Sivaraman Balakrishnan, and Zachary C. Lipton. Mixture proportion estimation and PU learning: A modern approach. In *Advances in Neural Information Processing Systems 34*, pp. 8532–8544, 2021.
- Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of neural networks. In *Proceedings of the 31st Conference On Learning Theory*, pp. 297–299, 2018.
- Chen Gong, Qizhou Wang, Tongliang Liu, Bo Han, Jane You, Jian Yang, and Dacheng Tao. Instance-dependent positive and unlabeled learning with labeling bias estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(8):4163–4177, 2022.

- Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *Proceedings of the 9th International Conference on Learning Representations*, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Cho-Jui Hsieh, Nagarajan Natarajan, and Inderjit Dhillon. PU learning for matrix completion. In *Proceedings of the 32nd International Conference on Machine Learning*, pp. 2445–2453, 2015.
- Yu-Guan Hsieh, Gang Niu, and Masashi Sugiyama. Classification from positive, unlabeled and biased negative data. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 2820–2829, 2019.
- Wenpeng Hu, Ran Le, Bing Liu, Feng Ji, Jinwen Ma, Dongyan Zhao, and Rui Yan. Predictive adversarial learning from positive and unlabeled data. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, pp. 7806–7814, 2021.
- Yangbangyan Jiang, Qianqian Xu, Yunrui Zhao, Zhiyong Yang, Peisong Wen, Xiaochun Cao, and Qingming Huang. Positive-unlabeled learning with label distribution alignment. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(12):15345–15363, 2023.
- Hyunjun Ju, Dongha Lee, Junyoung Hwang, Junghyun Namkung, and Hwanjo Yu. Pumad: Pu metric learning for anomaly detection. *Information Sciences*, 523:167–183, 2020.
- Masahiro Kato, Takeshi Teshima, and Junya Honda. Learning from positive and unlabeled data with a selection bias. In *Proceedings of the 7th International Conference on Learning Representations*, 2019.
- Markelle Kelly, Rachel Longjohn, and Kolby Nottingham. The UCI machine learning repository, 2023.
- Ryuichi Kiryo, Gang Niu, Marthinus C. du Plessis, and Masashi Sugiyama. Positive-unlabeled learning with non-negative risk estimator. In *Advances in Neural Information Processing Systems 30*, pp. 1674–1684, 2017.
- Alex Krizhevsky and Geoffrey E. Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- Changchun Li, Ximing Li, Lei Feng, and Jihong Ouyang. Who is your right mixup partner in positive and unlabeled learning. In *Proceedings of the 10th International Conference on Learning Representations*, 2022.
- Changchun Li, Yuanchao Dai, Lei Feng, Ximing Li, Bing Wang, and Jihong Ouyang. Positive and unlabeled learning with controlled probability boundary fence. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 27641–27652, 2024.
- Lin Long, Haobo Wang, Zhijie Jiang, Lei Feng, Chang Yao, Gang Chen, and Junbo Zhao. Positive-unlabeled learning by latent group-aware meta disambiguation. In *Proceedings of the 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23138–23147, 2024.
- Yuren Mao, Yu Hao, Xin Cao, Yunjun Gao, Chang Yao, and Xuemin Lin. Boosting GNN-based link prediction via PU-AUC optimization. *IEEE Transactions on Knowledge and Data Engineering*, 37(4):1635–1649, 2025.
- Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning*, pp. 807–814, 2010.
- Gang Niu, Marthinus C. du Plessis, Tomoya Sakai, Yao Ma, and Masashi Sugiyama. Theoretical comparisons of positive-unlabeled learning against positive-negative learning. In *Advances in Neural Information Processing Systems 29*, pp. 1199–1207, 2016.

- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, 2019.
- Harish G. Ramaswamy, Clayton Scott, and Ambuj Tewari. Mixture proportion estimation via kernel embeddings of distributions. In *Proceedings of the 33rd International Conference on Machine Learning*, pp. 2052–2060, 2016.
- Sebastian Raschka. Model evaluation, model selection, and algorithm selection in machine learning. *arXiv preprint arXiv:1811.12808*, 2018.
- Tomoya Sakai, Marthinus Christoffel du Plessis, Gang Niu, and Masashi Sugiyama. Semi-supervised classification based on classification from positive and unlabeled data. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 2998–3006, 2017.
- Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge University Press, Cambridge, UK, 2014.
- Vinay Shukla, Zhe Zeng, Kareem Ahmed, and Guy Van den Broeck. A unified approach to count-based weakly supervised learning. In *NeurIPS*, 2023.
- Masashi Sugiyama, Han Bao, Takashi Ishida, Nan Lu, Tomoya Sakai, and Gang Niu. *Machine learning from weak supervision: An empirical risk minimization approach*. MIT Press, 2022.
- Hiroshi Takahashi, Tomoharu Iwata, Atsutoshi Kumagai, Yuuki Yamanaka, and Tomoya Yamashita. Positive-unlabeled diffusion models for preventing sensitive data generation. In *Proceedings of the 13th International Conference on Learning Representations*, 2025.
- Paweł Teisseyre, Timo Martens, Jessa Bekker, and Jesse Davis. Learning from biased positive-unlabeled data via threshold calibration. In *Proceedings of the 28th International Conference on Artificial Intelligence and Statistics*, pp. 2314–2322, 2025.
- Yuchuan Tian, Hanting Chen, Xutao Wang, Zheyuan Bai, Qinghua ZHANG, Ruifeng Li, Chao Xu, and Yunhe Wang. Multiscale positive-unlabeled detection of AI-generated texts. In *Proceedings of the 12th International Conference on Learning Representations*, 2024.
- Vladimir Vapnik. Statistical learning theory. *John Wiley & Sons*, 2, 1998.
- Wei Wang, Dong-Dong Wu, Jindong Wang, Gang Niu, Min-Ling Zhang, and Masashi Sugiyama. Realistic evaluation of deep partial-label learning algorithms. In *Proceedings of the 13th International Conference on Learning Representations*, 2025.
- Xinrui Wang, Wenhai Wan, Chuanxing Geng, Shao-Yuan Li, and Songcan Chen. Beyond myopia: Learning from positive and unlabeled data through holistic predictive trends. In *Advances in Neural Information Processing Systems 36*, 2023a.
- Xutao Wang, Hanting Chen, Tianyu Guo, and Yunhe Wang. PUE: Biased positive-unlabeled learning enhancement by causal inference. In *Advances in Neural Information Processing Systems 36*, pp. 19783–19798, 2023b.
- Tong Wei, Hai Wang, Weiwei Tu, and Yufeng Li. Robust model selection for positive and unlabeled learning with constraints. *Science China Information Sciences*, 65(11):212101, 2022.
- Yuhao Wu, Jiangchao Yao, Bo Han, Lina Yao, and Tongliang Liu. Unraveling the impact of heterophilic structures on graph positive-unlabeled learning. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 53928–53943, 2024.
- Yu Yao, Tongliang Liu, Bo Han, Mingming Gong, Gang Niu, Masashi Sugiyama, and Dacheng Tao. Rethinking class-prior estimation for positive-unlabeled learning. In *Proceedings of the 10th International Conference on Learning Representations*, 2022.
- Jinfeng Yi, Cho-Jui Hsieh, Kush R. Varshney, Lijun Zhang, and Yao Li. Scalable demand-aware recommendation. In *Advances in Neural Information Processing Systems 30*, pp. 2412–2421, 2017.

- Hang Yin, Liyao Xiang, Dong Ding, Yuheng He, Yihan Wu, Pengzhi Chu, Xinbing Wang, and Chenghu Zhou. Lambda: Learning matchable prior for entity alignment with unlabeled dangling cases. In *Advances in Neural Information Processing Systems*, volume 37, pp. 78964–78995, 2024.
- Yunrui Zhao, Qianqian Xu, Yangbangyan Jiang, Peisong Wen, and Qingming Huang. Dist-PU: Positive-unlabeled learning from a label distribution perspective. In *Proceedings of the 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14441–14450, 2022.
- Yao Zhou, Jianpeng Xu, Jun Wu, Zeinab Taghavi, Evren Korpeoglu, Kannan Achan, and Jingrui He. Pure: Positive-unlabeled recommendation with generative adversarial network. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 2409–2419, 2021.
- Yilun Zhu, Aaron Fjeldsted, Darren Holland, George Landon, Azaree Lintereur, and Clayton Scott. Mixture proportion estimation beyond irreducibility. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 42962–42982, 2023a.
- Zhangchi Zhu, Lu Wang, Pu Zhao, Chao Du, Wei Zhang, Hang Dong, Bo Qiao, Qingwei Lin, Saravan Rajmohan, and Dongmei Zhang. Robust positive-unlabeled learning via noise negative sample self-correction. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 3663–3673, 2023b.

THE USE OF LARGE LANGUAGE MODELS (LLMs)

LLMs were used solely for correcting grammatical and spelling errors.

A COMPARISON WITH A CONCURRENT WORK

A concurrent work (Chen et al., 2026) also proposed a PU learning benchmark. The authors primarily focused on establishing a fair comparison and evaluation pipeline for PU learning algorithms. They noted the differences between one-sample and two-sample settings, yet they did not analyze the ILS problem or propose solutions. Furthermore, they assumed that the validation set had true labels, which is unrealistic, as discussed in the paper.

B PROOFS

B.1 PROOF OF PROPOSITION 1

$$\begin{aligned}
\text{ACC}(f) &= \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} [\mathbb{I}(f(\mathbf{x}) \geq 0)] + (1 - \pi) \mathbb{E}_{p(\mathbf{x}|y=-1)} [\mathbb{I}(f(\mathbf{x}) < 0)] \\
&= \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} [\mathbb{I}(f(\mathbf{x}) \geq 0)] + \mathbb{E}_{p(\mathbf{x})} [\mathbb{I}(f(\mathbf{x}) < 0)] - \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} [\mathbb{I}(f(\mathbf{x}) < 0)] \\
&= \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} [\mathbb{I}(f(\mathbf{x}) \geq 0)] + \mathbb{E}_{p(\mathbf{x})} [\mathbb{I}(f(\mathbf{x}) < 0)] - \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} [1 - \mathbb{I}(f(\mathbf{x}) \geq 0)] \\
&= 2\pi \mathbb{E}_{p(\mathbf{x}|y=+1)} [\mathbb{I}(f(\mathbf{x}) \geq 0)] + \mathbb{E}_{p(\mathbf{x})} [\mathbb{I}(f(\mathbf{x}) < 0)] - \pi \\
&= \mathbb{E}[\text{PA}(f)] - \pi.
\end{aligned}$$

Here, the last equation is obtained since $\mathcal{D}_U \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})$ for the TS setting and $\mathcal{D}_P \cup \mathcal{D}_U \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})$ for the OS setting. Therefore, for two classifiers f_1 and f_2 that satisfy $\mathbb{E}[\text{PA}(f_1)] < \mathbb{E}[\text{PA}(f_2)]$, we have $\text{ACC}(f_1) < \text{ACC}(f_2)$. The proof is complete. \square

B.2 PROOF OF PROPOSITION 2

For the TS setting,

$$\begin{aligned}
\text{AUC}(f) &= \mathbb{E}_{p(\mathbf{x}|y=+1)} \mathbb{E}_{p(\mathbf{x}'|y'=-1)} \left[\mathbb{I}(f(\mathbf{x}) > f(\mathbf{x}')) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}) = f(\mathbf{x}')) \right] \\
&= \frac{1}{1 - \pi} \mathbb{E}_{p(\mathbf{x}|y=+1)} \mathbb{E}_{p(\mathbf{x}')} \left[\mathbb{I}(f(\mathbf{x}) > f(\mathbf{x}')) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}) = f(\mathbf{x}')) \right] \\
&\quad - \frac{\pi}{1 - \pi} \mathbb{E}_{p(\mathbf{x}|y=+1)} \mathbb{E}_{p(\mathbf{x}'|y'=+1)} \left[\mathbb{I}(f(\mathbf{x}) > f(\mathbf{x}')) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}) = f(\mathbf{x}')) \right] \\
&= \frac{1}{1 - \pi} \mathbb{E}_{p(\mathbf{x}|y=+1)} \mathbb{E}_{p(\mathbf{x}')} \left[\mathbb{I}(f(\mathbf{x}) > f(\mathbf{x}')) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}) = f(\mathbf{x}')) \right] - \frac{\pi}{2 - 2\pi} \\
&= \frac{1}{1 - \pi} \mathbb{E}[\text{PAUC}(f)] - \frac{\pi}{2 - 2\pi}.
\end{aligned}$$

For the OS setting,

$$\begin{aligned}
\text{AUC}(f) &= \mathbb{E}_{p(\mathbf{x}|y=+1)} \mathbb{E}_{p(\mathbf{x}'|y'=-1)} \left[\mathbb{I}(f(\mathbf{x}) > f(\mathbf{x}')) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}) = f(\mathbf{x}')) \right] \\
&= \frac{1}{1 - \bar{\pi}} \mathbb{E}_{p(\mathbf{x}|y=+1)} \mathbb{E}_{p(\mathbf{x}')} \left[\mathbb{I}(f(\mathbf{x}) > f(\mathbf{x}')) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}) = f(\mathbf{x}')) \right] \\
&\quad - \frac{\bar{\pi}}{1 - \bar{\pi}} \mathbb{E}_{p(\mathbf{x}|y=+1)} \mathbb{E}_{p(\mathbf{x}'|y'=+1)} \left[\mathbb{I}(f(\mathbf{x}) > f(\mathbf{x}')) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}) = f(\mathbf{x}')) \right] \\
&= \frac{1}{1 - \bar{\pi}} \mathbb{E}_{p(\mathbf{x}|y=+1)} \mathbb{E}_{p(\mathbf{x}')} \left[\mathbb{I}(f(\mathbf{x}) > f(\mathbf{x}')) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}) = f(\mathbf{x}')) \right] - \frac{\bar{\pi}}{2 - 2\bar{\pi}} \\
&= \frac{1}{1 - \bar{\pi}} \mathbb{E}[\text{PAUC}(f)] - \frac{\bar{\pi}}{2 - 2\bar{\pi}}.
\end{aligned}$$

Therefore, under both OS and TS settings, for two classifiers f_1 and f_2 that satisfy $\mathbb{E}[\text{PAUC}(f_1)] < \mathbb{E}[\text{PAUC}(f_2)]$, we have $\text{AUC}(f_1) < \text{AUC}(f_2)$. \square

B.3 BIAS OF THE RISK ESTIMATOR

Under the OS setting, we have

$$\begin{aligned} \mathbb{E}[\widehat{R}(f)] - R(f) &= \mathbb{E}_{\bar{p}(\mathbf{x})}[\ell(f(\mathbf{x}), -1)] - \mathbb{E}_{p(\mathbf{x})}[\ell(f(\mathbf{x}), -1)] \\ &= (\bar{\pi} - \pi) \left(\mathbb{E}_{p(\mathbf{x}|y=+1)}[\ell(f(\mathbf{x}), -1)] - \mathbb{E}_{p(\mathbf{x}|y=-1)}[\ell(f(\mathbf{x}), -1)] \right), \end{aligned}$$

which is not equal to 0. Therefore, it means that the bias of the risk estimator always exists. Then, the minimizers of $\mathbb{E}[\widehat{R}(f)]$ and $R(f)$ are not the same.

B.4 PROOF OF THEOREM 1

First, we have

$$\begin{aligned} \bar{p}(\mathbf{x}) &= \bar{\pi}p(\mathbf{x}|y=+1) + (1 - \bar{\pi})p(\mathbf{x}|y=-1) \\ &= \frac{(1-c)\pi}{1-c\pi}p(\mathbf{x}|y=+1) + \frac{1-\pi}{1-c\pi}p(\mathbf{x}|y=-1). \end{aligned}$$

Therefore, we have

$$p(\mathbf{x}|y=-1) = \frac{1-c\pi}{1-\pi}\bar{p}(\mathbf{x}) - \frac{(1-c)\pi}{1-\pi}p(\mathbf{x}|y=+1).$$

Then,

$$\begin{aligned} R(f) &= \pi\mathbb{E}_{p(\mathbf{x}|y=+1)}[\ell(f(\mathbf{x}), +1)] + (1 - \pi)\mathbb{E}_{p(\mathbf{x}|y=-1)}[\ell(f(\mathbf{x}), -1)] \\ &= \pi\mathbb{E}_{p(\mathbf{x}|y=+1)}[\ell(f(\mathbf{x}), +1)] + (1 - c\pi)\mathbb{E}_{\bar{p}(\mathbf{x})}[\ell(f(\mathbf{x}), -1)] - (1 - c)\pi\mathbb{E}_{p(\mathbf{x}|y=+1)}[\ell(f(\mathbf{x}), -1)] \\ &= \pi\mathbb{E}_{p(\mathbf{x}|y=+1)}[\ell(f(\mathbf{x}), +1) + (c-1)\ell(f(\mathbf{x}), -1)] + (1 - c\pi)\mathbb{E}_{\bar{p}(\mathbf{x})}[\ell(f(\mathbf{x}), -1)], \end{aligned}$$

which concludes the proof. \square

B.5 PROOF OF THEOREM 2

Definition 5 (Rademacher complexity). Let $\mathcal{X}_{n_P}^P = \{\mathbf{x}_1, \dots, \mathbf{x}_{n_P}\}$ denote n_P i.i.d. random variables drawn from density $p(\mathbf{x}|y=+1)$. Let $\mathcal{X}_{n_U}^U = \{\mathbf{x}_{n_P+1}, \dots, \mathbf{x}_{n_P+n_U}\}$ denote n_U i.i.d. random variables drawn from density $\bar{p}(\mathbf{x})$. Let $\mathcal{F} = \{f : \mathcal{X} \mapsto \mathbb{R}\}$ denote a class of measurable functions, $\boldsymbol{\sigma}_P = (\sigma_1, \sigma_2, \dots, \sigma_{n_P})$, and $\boldsymbol{\sigma}_U = (\sigma_{n_P+1}, \sigma_{n_P+2}, \dots, \sigma_{n_P+n_U})$ denote Rademacher variables taking values from $\{+1, -1\}$ uniformly. Then, the (expected) Rademacher complexities of \mathcal{F} are defined as

$$\begin{aligned} \mathfrak{R}_{n_P}(\mathcal{F}) &= \mathbb{E}_{\mathcal{X}_{n_P}^P} \mathbb{E}_{\boldsymbol{\sigma}_P} \left[\sup_{f \in \mathcal{F}} \frac{1}{n_P} \sum_{i=1}^{n_P} \sigma_i f(\mathbf{x}_i) \right], \\ \mathfrak{R}'_{n_U}(\mathcal{F}) &= \mathbb{E}_{\mathcal{X}_{n_U}^U} \mathbb{E}_{\boldsymbol{\sigma}_U} \left[\sup_{f \in \mathcal{F}} \frac{1}{n_U} \sum_{i=n_P+1}^{n_P+n_U} \sigma_i f(\mathbf{x}_i) \right]. \end{aligned}$$

Lemma 1. For any $\delta > 0$, we have the following inequality with probability at least $1 - \delta$:

$$\begin{aligned} \sup_{f \in \mathcal{F}} |\bar{R}(f) - R(f)| &\leq 2(2-c)\pi L_\ell \mathfrak{R}_{n_P}(\mathcal{F}) + 2(1-c\pi)L_\ell \mathfrak{R}'_{n_U}(\mathcal{F}) \\ &\quad + \left(\frac{\pi(2-c)C_\ell}{\sqrt{n_P}} + \frac{(1-c\pi)C_\ell}{\sqrt{n_U}} \right) \sqrt{\frac{\ln 2/\delta}{2}}. \end{aligned}$$

Proof. First, we give the upper bound for the one-side uniform deviation $\sup_{f \in \mathcal{F}} (\bar{R}(f) - R(f))$. When an instance in $\mathcal{X}_{n_P}^P$ is replaced by another instance, the value of $\sup_{f \in \mathcal{F}} (\bar{R}(f) - R(f))$

changes at most $\pi(2-c)C_\ell/n_P$; when an instance in $\mathcal{X}_{n_U}^U$ is replaced by another instance, the value of $\sup_{f \in \mathcal{F}} (\bar{R}(f) - R(f))$ changes at most $(1-c\pi)C_\ell/n_U$. Therefore, according to McDiarmid's inequality, we have the following inequality with probability at least $1 - \delta/2$:

$$\begin{aligned} \sup_{f \in \mathcal{F}} (\bar{R}(f) - R(f)) &\leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} (\bar{R}(f) - R(f)) \right] + \sqrt{\frac{\pi^2(2-c)^2 C_\ell^2}{n_P} + \frac{(1-c\pi)^2 C_\ell^2}{n_U}} \sqrt{\frac{\ln 2/\delta}{2}} \\ &\leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} (\bar{R}(f) - R(f)) \right] + \left(\frac{\pi(2-c)C_\ell}{\sqrt{n_P}} + \frac{(1-c\pi)C_\ell}{\sqrt{n_U}} \right) \sqrt{\frac{\ln 2/\delta}{2}}. \end{aligned}$$

Then, by symmetrization (Vapnik, 1998), it is a routine work to have

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} (\bar{R}(f) - R(f)) \right] \leq 2(2-c)\pi \mathfrak{R}_{n_P}(\ell \circ \mathcal{F}) + 2(1-c\pi) \mathfrak{R}'_{n_U}(\ell \circ \mathcal{F}).$$

According to Talagrand's contraction lemma (Shalev-Shwartz & Ben-David, 2014), we have

$$\mathfrak{R}_{n_P}(\ell \circ \mathcal{F}) \leq L_\ell \mathfrak{R}_{n_P}(\mathcal{F}), \quad \mathfrak{R}'_{n_U}(\ell \circ \mathcal{F}) \leq L_\ell \mathfrak{R}'_{n_U}(\mathcal{F}).$$

By combining the above inequalities, we have the following inequality with probability at least $1 - \delta/2$:

$$\begin{aligned} \sup_{f \in \mathcal{F}} (\bar{R}(f) - R(f)) &\leq 2(2-c)\pi L_\ell \mathfrak{R}_{n_P}(\mathcal{F}) + 2(1-c\pi) L_\ell \mathfrak{R}'_{n_U}(\mathcal{F}) \\ &\quad + \left(\frac{\pi(2-c)C_\ell}{\sqrt{n_P}} + \frac{(1-c\pi)C_\ell}{\sqrt{n_U}} \right) \sqrt{\frac{\ln 2/\delta}{2}}. \end{aligned}$$

In a similar way, we have the following inequality with probability at least $1 - \delta/2$:

$$\begin{aligned} \sup_{f \in \mathcal{F}} (R(f) - \bar{R}(f)) &\leq 2(2-c)\pi L_\ell \mathfrak{R}_{n_P}(\mathcal{F}) + 2(1-c\pi) L_\ell \mathfrak{R}'_{n_U}(\mathcal{F}) \\ &\quad + \left(\frac{\pi(2-c)C_\ell}{\sqrt{n_P}} + \frac{(1-c\pi)C_\ell}{\sqrt{n_U}} \right) \sqrt{\frac{\ln 2/\delta}{2}}. \end{aligned}$$

Therefore, we have the following inequality with probability at least $1 - \delta$:

$$\begin{aligned} \sup_{f \in \mathcal{F}} |\bar{R}(f) - R(f)| &\leq 2(2-c)\pi L_\ell \mathfrak{R}_{n_P}(\mathcal{F}) + 2(1-c\pi) L_\ell \mathfrak{R}'_{n_U}(\mathcal{F}) \\ &\quad + \left(\frac{\pi(2-c)C_\ell}{\sqrt{n_P}} + \frac{(1-c\pi)C_\ell}{\sqrt{n_U}} \right) \sqrt{\frac{\ln 2/\delta}{2}}. \end{aligned}$$

The proof is complete. \square

Then, we give the proof of Theorem 2.

Proof of Theorem 2.

$$\begin{aligned} R(\bar{f}) - R(f^*) &= R(\bar{f}) - \bar{R}(\bar{f}) + \bar{R}(\bar{f}) - \bar{R}(f^*) + \bar{R}(f^*) - R(f^*) \\ &\leq R(\bar{f}) - \bar{R}(\bar{f}) + \bar{R}(\bar{f}) - \bar{R}(f^*) + \bar{R}(f^*) - R(f^*) \\ &\leq 2 \sup_{f \in \mathcal{F}} |\bar{R}(f) - R(f)|. \end{aligned}$$

By Lemma 1, the proof is complete. \square

B.6 DERIVATION OF EQUIVALENCE OF RISK ESTIMATORS

$$\begin{aligned}
& \bar{R}(f) \\
&= \frac{\pi}{n_P} \sum_{i=1}^{n_P} (\ell(f(\mathbf{x}_i), +1) + (c-1)\ell(f(\mathbf{x}_i), -1)) + \frac{1-c\pi}{n_U} \sum_{i=n_P+1}^{n_P+n_U} \ell(f(\mathbf{x}_i), -1) \\
&= \sum_{i=1}^{n_P} \left(\frac{\pi}{n_P} \ell(f(\mathbf{x}_i), +1) + \left(\frac{1}{n_P+n_U} - \frac{\pi}{n_P} \right) \ell(f(\mathbf{x}_i), -1) \right) + \frac{1}{n_P+n_U} \sum_{i=n_P+1}^{n_P+n_U} \ell(f(\mathbf{x}_i), -1) \\
&= \frac{\pi}{n_P} \sum_{i=1}^{n_P} (\ell(f(\mathbf{x}_i), +1) - \ell(f(\mathbf{x}_i), -1)) + \frac{1}{n_U} \sum_{i=1}^{n_P+n_U} \ell(f(\mathbf{x}_i), -1), \tag{10}
\end{aligned}$$

where the second equation uses the estimation $c = n_P/\pi(n_P + n_U)$.

C MORE EXPERIMENTAL DETAILS

C.1 MORE DETAILS OF BENCHMARK DATASETS

Table 5 summarizes their key characteristics, including the number of examples, feature dimensionality, positive class configurations, and task domains. For all datasets, we vary the positive rate in $\{10\%, 20\%, 30\%, 40\%, 50\%\}$. For the benchmark experiments in Section 5, we used the positive rate 30%.

Table 5: Summary of datasets used in this PU learning benchmark.

Dataset	# Examples	# Features	Positive Classes (Case 1)	Positive Classes (Case 2)	Task Domain
CIFAR-10	20,000	3,072	{0,1,2,8,9}	{2,3,5,7,9}	Image classification
ImageNette	6,000	12,288	{0,1,2,8,9}	{2,3,5,7,9}	Image classification
USPS	4,000	256	{4,7,9,5,8}	{1,6,4,9,8}	Digit recognition
Letter	13,000	16	{B,V,L,R,I,O,W,S,J,K,C,H,Z}	{D,T,A,Y,Q,G,B,L,I,W,J,C,Z}	Character recognition

C.2 DESCRIPTIONS OF ALGORITHMS

- **uPU** (du Plessis et al., 2015): An unbiased risk estimator that is convex when the loss function satisfies certain linear-odd conditions.
- **nnPU** (Kiryo et al., 2017): A non-negative risk estimator that alleviates the overfitting issue in PU learning.
- **nnPU-GA** (Kiryo et al., 2017): A gradient ascent formulation of nnPU.
- **PUSB** (Kato et al., 2019): A method that accounts for selection bias in the labeling process.
- **PubN** (Hsieh et al., 2019): A framework that incorporates biased negative data into empirical risk minimization.
- **VPU** (Chen et al., 2020a): A variational approach that directly evaluates the modeling error of a Bayesian classifier from data.
- **PAN** (Hu et al., 2021): A predictive adversarial network built upon the generative adversarial network framework.
- **CVIR** (Garg et al., 2021): A mixture-proportion estimation method combining best bin estimation and conditional Value Ignoring Risk.
- **Dist-PU** (Zhao et al., 2022): A method that enforces consistency between predicted and ground-truth label distributions.
- **P³MIX-E** (Li et al., 2022): A mixup-based method that pairs marginal pseudo-negative instances with boundary-near positive instances, with early-learning regularization.
- **P³MIX-C** (Li et al., 2022): A mixup-based method that pairs marginal pseudo-negative instances with boundary-near positive instances, with pseudo-negative correction.
- **LBE** (Gong et al., 2022): An instance-dependent PU algorithm that jointly estimates labeling bias and learns the classifier.
- **Count Loss** (Shukla et al., 2023): A unified approach introducing a count-based loss penalizing deviations from arithmetic label-count constraints.

- Robust-PU (Zhu et al., 2023b): A reweighted learning framework that dynamically adjusts sample weights based on training progress and sample hardness.
- Holistic-PU (Wang et al., 2023a): A holistic method interpreting prediction scores as a temporal point process.
- PUE (Wang et al., 2023b): A causality-based method that reconstructs the loss via normalized propensity scores and inverse probability weighting.
- GLWS (Chen et al., 2024): A general weak-supervision framework formulated as Expectation-Maximization, accommodating PU data as one supervision source.

C.3 IMPLEMENTATION DETAILS

All algorithms were implemented in PyTorch (Paszke et al., 2019), and all experiments were conducted on a single NVIDIA Tesla V100 GPU. We used the SGD optimizer and trained for 20,000 iterations across all datasets. Model performance on the validation and test sets was recorded every 100 iterations. For each dataset, we generated three random data splits. For each split, 10 random hyperparameter configurations were sampled from a predefined pool. Table 6 provides the details of the hyperparameter configurations used for all algorithms.

Table 6: Hyperparameters, their default values, and distributions for random search.

Condition	Parameter	Default Value	Random Distribution
ResNet	learning rate	0.001	$10^{\text{Uniform}(-4.5, -2.5)}$
	batch size	64	$2^{\text{Uniform}(5, 8)}$
	momentum	0.9	0.9
MLP	learning rate	0.001	$10^{\text{Uniform}(-4.5, -2.5)}$
	batch size	128	$2^{\text{Uniform}(4, 7)}$
	momentum	0.9	0.9
nnPU	tolerance threshold	0.0	0.0
PUBN	importance of unlabeled data	0.5	RandomChoice([0.5, 0.7, 0.9])
PAN	balance factor of the KL-divergences	0.0001	0.0001
P ³ MIX-E	predictive score threshold	0.85	0.85
	size of the candidate mixup pool	96	96
	weight of the positive loss	1	1
	weight of the unlabeled loss	1	1
	weight of the entropy loss	0.5	0.5
P ³ MIX-C	weight of the early-learning regularization	5	5
	predictive score threshold	0.8	0.8
	size of the candidate mixup pool	96	96
	mixup coefficient	1.0	1.0
	weight of the positive loss	1	1
	weight of the unlabeled loss	1	1
LBE	weight of the entropy loss	0.1	0.1
	warm up iteration	2000	2000
Robust-PU	warm up iteration	2000	2000
	training scheduler	linear	linear
	temperature in the logistic loss	1	RandomChoice([1, 1.3])
	initial threshold	0.1	RandomChoice([0.1, 0.11])
	final threshold	2	RandomChoice([1, 2])
Holistic-PU	growing step	10	RandomChoice([5, 10])
	warm up iteration	2000	2000

D DETAILS OF EXPERIMENTAL RESULTS

Tables 7 to 18 report detailed experimental results in terms of different metrics on CIFAR-10, ImageNette, Letter, and USPS, and the hyperparameters are determined with PA, PAUC, and OA, respectively.

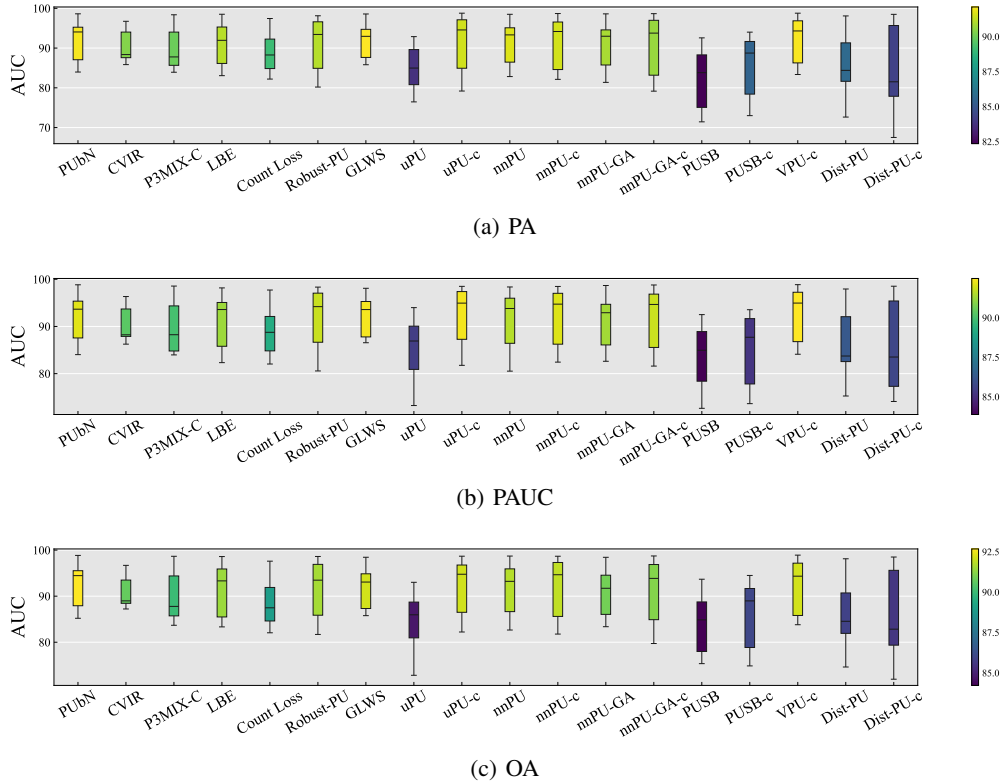


Figure 5: Overall performance w.r.t. the AUC score of different algorithms across all datasets. Hyperparameters were tuned using PA, PAUC and OA, respectively; bar colors indicate means.

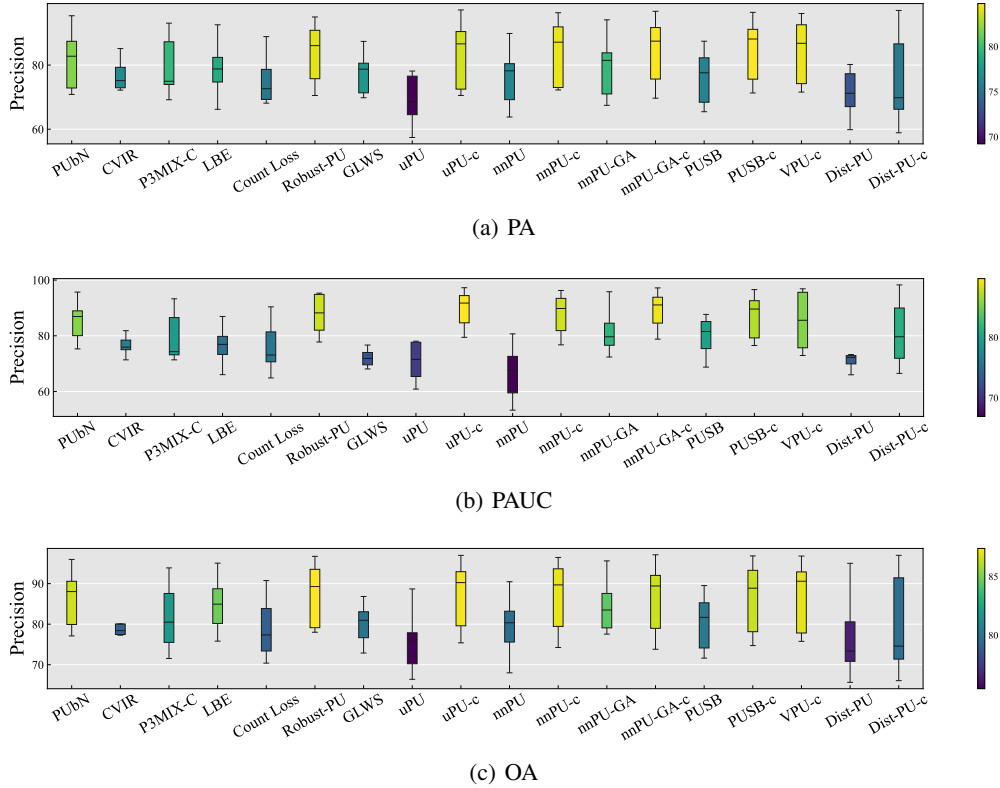


Figure 6: Overall performance w.r.t. precision of different algorithms across all datasets. Hyperparameters were tuned using PA, PAUC and OA, respectively; bar colors indicate means.

Table 7: Test results (mean \pm std) of precision and recall for each algorithm on CIFAR-10 (Case 1) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Precision			Recall		
Val metric	PA	PAUC	OA	PA	PAUC	OA
PUBN	85.58 \pm 0.37	86.97 \pm 0.18	89.46 \pm 0.65	87.71 \pm 1.00	85.25 \pm 1.87	84.64 \pm 0.24
PAN	71.61 \pm 0.86	73.33 \pm 0.14	79.51 \pm 1.74	88.39 \pm 1.65	86.65 \pm 1.61	78.41 \pm 3.09
CVIR	82.12 \pm 1.38	78.27 \pm 1.03	86.30 \pm 0.93	90.72 \pm 0.42	92.42 \pm 1.35	86.74 \pm 0.26
P3MIX-E	68.93 \pm 6.62	50.00 \pm 0.00	82.77 \pm 5.96	90.99 \pm 1.97	100.00\pm0.00	67.19 \pm 17.74
P3MIX-C	86.03 \pm 0.91	84.91 \pm 1.00	86.27 \pm 0.66	86.86 \pm 0.24	86.97 \pm 0.45	87.19 \pm 0.81
LBE	79.00 \pm 1.37	66.06 \pm 1.22	88.64 \pm 0.96	89.31 \pm 0.95	97.45 \pm 0.33	80.41 \pm 0.36
Count Loss	75.81 \pm 0.30	74.78 \pm 1.63	79.88 \pm 0.66	90.73 \pm 0.31	90.57 \pm 1.51	86.65 \pm 1.07
Robust-PU	84.19 \pm 1.05	89.77 \pm 1.85	88.23 \pm 0.69	87.73 \pm 1.21	80.72 \pm 2.97	82.89 \pm 0.26
Holistic-PU	50.10 \pm 0.05	50.00 \pm 0.00	78.03 \pm 0.77	99.49\pm0.22	100.00\pm0.00	88.60 \pm 0.41
PUe	74.23 \pm 1.27	80.12 \pm 1.97	85.70 \pm 2.17	85.52 \pm 0.48	76.39 \pm 2.72	73.49 \pm 1.77
GLWS	79.61 \pm 0.65	73.12 \pm 2.81	82.86 \pm 0.85	92.70 \pm 0.34	95.53 \pm 1.03	89.99 \pm 0.32
uPU	78.14 \pm 1.90	78.06 \pm 7.35	88.69 \pm 0.56	84.29 \pm 0.41	80.13 \pm 7.43	73.44 \pm 0.66
uPU-c	85.58 \pm 0.88	89.53 \pm 1.39	88.50 \pm 0.93	86.39 \pm 0.98	77.69 \pm 2.63	83.91 \pm 0.66
nnPU	77.17 \pm 0.27	68.25 \pm 0.27	78.73 \pm 1.19	91.00 \pm 0.32	95.62 \pm 0.62	88.99 \pm 1.44
nnPU-c	83.72 \pm 0.48	87.75 \pm 0.73	86.66 \pm 0.46	88.22 \pm 0.98	83.76 \pm 0.68	85.95 \pm 0.77
nnPU-GA	81.92 \pm 1.66	82.29 \pm 0.38	86.81 \pm 1.62	88.18 \pm 1.20	87.12 \pm 0.72	82.54 \pm 0.65
nnPU-GA-c	85.33 \pm 0.90	89.78 \pm 1.03	89.25 \pm 0.78	86.55 \pm 1.15	81.95 \pm 0.76	82.19 \pm 0.41
PUSB	76.20 \pm 1.34	78.13 \pm 1.39	78.64 \pm 0.98	91.93 \pm 0.81	90.41 \pm 0.19	90.44\pm0.20
PUSB-c	86.43 \pm 0.03	89.07 \pm 1.32	87.87 \pm 0.17	85.76 \pm 0.84	79.38 \pm 1.23	84.66 \pm 0.25
VPU	84.97 \pm 1.65	77.43 \pm 2.41	89.08 \pm 0.15	80.05 \pm 0.84	33.15 \pm 15.89	79.93 \pm 1.08
VPU-c	84.97 \pm 1.65	77.43 \pm 2.41	89.08 \pm 0.15	88.67 \pm 0.83	93.37 \pm 0.83	85.82 \pm 0.61
Dist-PU	76.34 \pm 0.77	72.78 \pm 0.89	84.75 \pm 0.15	91.79 \pm 0.56	93.81 \pm 0.86	81.86 \pm 1.26
Dist-PU-c	84.07 \pm 1.13	88.22 \pm 2.06	90.49 \pm 0.84	91.58 \pm 0.99	86.67 \pm 2.28	86.02 \pm 0.83

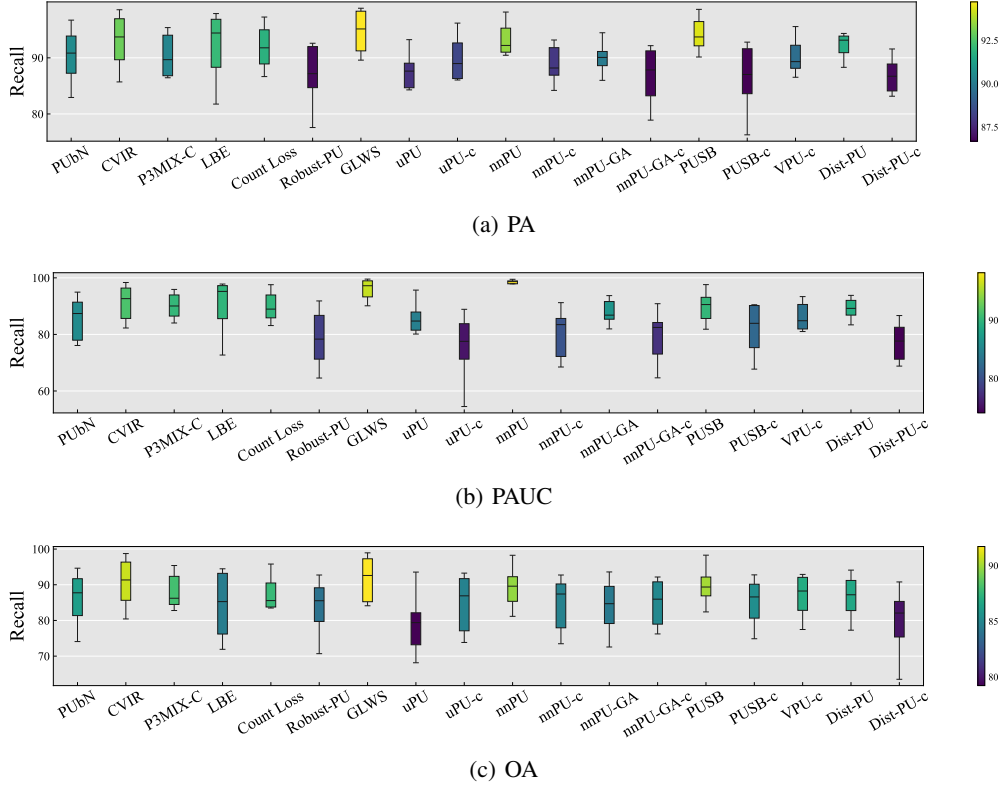


Figure 7: Overall performance w.r.t. recall of different algorithms across all datasets. Hyperparameters were tuned using PA, PAUC and OA, respectively; bar colors indicate means.

Table 8: Test results (mean \pm std) of precision and recall for each algorithm on CIFAR-10 (Case 2) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Precision			Recall		
Val metric	PA	PAUC	OA	PA	PAUC	OA
PUBN	73.21 \pm 1.41	80.55 \pm 1.72	80.14 \pm 1.05	89.42 \pm 0.70	78.25 \pm 3.04	79.75 \pm 0.94
PAN	57.23 \pm 1.87	56.98 \pm 2.87	59.50 \pm 1.67	92.77 \pm 0.38	94.70 \pm 2.27	83.05 \pm 5.42
CVIR	73.13 \pm 1.91	76.32 \pm 1.46	77.51 \pm 0.90	90.57 \pm 0.80	85.61 \pm 1.27	85.81 \pm 0.63
P3MIX-E	55.91 \pm 3.28	33.33 \pm 13.61	56.04 \pm 3.39	96.32 \pm 1.97	66.67 \pm 27.22	96.07\pm2.17
P3MIX-C	74.10 \pm 1.68	74.09 \pm 2.13	75.63 \pm 0.62	86.67 \pm 1.25	85.03 \pm 1.78	84.71 \pm 0.96
LBE	66.21 \pm 1.66	58.31 \pm 2.11	75.81 \pm 1.05	92.61 \pm 1.25	97.81 \pm 0.77	76.30 \pm 0.92
Count Loss	69.39 \pm 0.70	71.20 \pm 0.18	73.73 \pm 0.09	87.54 \pm 0.76	83.13 \pm 1.10	83.48 \pm 2.12
Robust-PU	75.32 \pm 1.11	82.16 \pm 1.82	79.06 \pm 0.85	86.25 \pm 1.51	72.97 \pm 1.95	80.59 \pm 1.55
Holistic-PU	53.00 \pm 0.10	59.93 \pm 4.67	67.62 \pm 0.43	98.93\pm0.21	54.64 \pm 23.85	81.48 \pm 5.21
PUe	63.65 \pm 0.22	62.69 \pm 2.20	71.16 \pm 1.47	86.71 \pm 0.82	88.03 \pm 2.08	71.53 \pm 3.84
GLWS	71.86 \pm 1.10	69.66 \pm 1.62	77.16 \pm 0.93	91.35 \pm 1.11	93.40 \pm 0.68	84.11 \pm 0.63
uPU	62.03 \pm 1.56	65.14 \pm 1.52	69.80 \pm 0.66	84.79 \pm 2.37	82.83 \pm 2.57	72.40 \pm 3.72
uPU-c	72.31 \pm 0.26	82.55 \pm 0.25	80.15 \pm 0.99	88.23 \pm 0.26	74.27 \pm 0.66	77.23 \pm 2.25
nnPU	68.39 \pm 1.63	57.41 \pm 0.78	74.03 \pm 0.92	91.01 \pm 1.51	98.65\pm0.40	85.19 \pm 0.73
nnPU-c	73.19 \pm 1.09	81.98 \pm 0.60	80.25 \pm 0.74	87.81 \pm 1.17	73.10 \pm 1.44	77.99 \pm 1.18
nnPU-GA	71.42 \pm 1.13	72.38 \pm 0.59	78.98 \pm 1.66	88.75 \pm 1.03	86.52 \pm 2.59	77.65 \pm 1.40
nnPU-GA-c	74.94 \pm 0.69	80.00 \pm 1.79	79.37 \pm 0.90	84.28 \pm 2.37	75.86 \pm 1.44	78.70 \pm 0.96
PUSB	69.38 \pm 0.61	76.30 \pm 1.64	74.79 \pm 0.62	92.21 \pm 0.15	84.05 \pm 2.62	85.62 \pm 1.96
PUSB-c	76.45 \pm 0.77	79.54 \pm 1.37	78.93 \pm 0.69	84.07 \pm 0.83	75.54 \pm 1.16	79.71 \pm 1.19
VPU	81.54\pm0.63	92.66\pm1.97	82.10\pm0.62	69.74 \pm 1.96	29.59 \pm 11.86	69.91 \pm 2.43
VPU-c	72.30 \pm 0.58	72.93 \pm 1.34	77.65 \pm 0.43	89.87 \pm 0.48	90.00 \pm 0.94	83.71 \pm 1.10
Dist-PU	68.12 \pm 0.72	71.23 \pm 1.19	71.27 \pm 1.24	88.31 \pm 0.40	83.64 \pm 1.35	83.07 \pm 1.45
Dist-PU-c	69.46 \pm 4.49	72.97 \pm 2.81	75.24 \pm 2.54	83.36 \pm 0.70	78.05 \pm 3.74	72.82 \pm 5.27

Table 9: Test results (mean \pm std) of precision and recall for each algorithm on ImageNette (Case 1) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Precision			Recall		
Val metric	PA	PAUC	OA	PA	PAUC	OA
PUBN	70.84 \pm 0.32	78.57 \pm 4.30	77.10 \pm 1.61	85.89 \pm 0.85	76.09 \pm 5.82	81.90 \pm 2.84
PAN	50.00 \pm 0.66	53.43 \pm 1.40	60.22 \pm 3.11	94.12 \pm 2.52	38.73 \pm 23.64	36.64 \pm 2.60
CVIR	74.49 \pm 1.81	75.11 \pm 3.74	78.28 \pm 1.36	86.93 \pm 1.69	85.67 \pm 3.79	85.13 \pm 1.14
P3MIX-E	71.37 \pm 2.70	32.75 \pm 13.37	71.98 \pm 2.64	81.92 \pm 1.42	66.67 \pm 27.22	81.48 \pm 1.38
P3MIX-C	74.23 \pm 1.92	74.56 \pm 2.46	78.59 \pm 2.18	87.47 \pm 0.80	87.34 \pm 0.62	82.78 \pm 1.94
LBE	76.29\pm0.82	81.74 \pm 1.60	80.73\pm0.77	81.76 \pm 0.93	73.48 \pm 3.85	75.86 \pm 2.11
Count Loss	69.02 \pm 1.15	71.41 \pm 2.49	74.81 \pm 1.18	89.37 \pm 0.93	86.22 \pm 1.82	83.77 \pm 2.66
Robust-PU	75.89 \pm 0.38	81.60 \pm 2.37	79.16 \pm 0.86	80.00 \pm 1.31	66.11 \pm 7.91	77.07 \pm 1.53
Holistic-PU	50.17 \pm 0.26	54.20 \pm 3.82	52.95 \pm 0.29	93.12 \pm 1.81	84.66 \pm 10.32	50.48 \pm 2.51
PUe	64.33 \pm 2.68	70.74 \pm 0.59	68.85 \pm 0.95	78.49 \pm 2.87	71.33 \pm 3.07	74.38 \pm 2.68
GLWS	69.81 \pm 1.18	68.10 \pm 1.82	75.17 \pm 0.67	90.92 \pm 1.13	92.88 \pm 0.54	84.51 \pm 0.74
uPU	65.37 \pm 1.04	60.87 \pm 4.80	70.38 \pm 0.55	87.71 \pm 0.58	89.47 \pm 4.66	80.15 \pm 1.07
uPU-c	70.53 \pm 2.49	85.36 \pm 2.00	77.92 \pm 1.59	86.05 \pm 3.74	54.50 \pm 12.20	76.69 \pm 1.98
nnPU	69.49 \pm 1.81	60.22 \pm 0.88	76.11 \pm 2.64	90.46 \pm 1.00	95.92\pm0.48	81.15 \pm 3.55
nnPU-c	72.51 \pm 1.01	81.53 \pm 2.56	77.04 \pm 0.57	84.20 \pm 0.62	69.52 \pm 4.15	77.77 \pm 2.26
nnPU-GA	69.74 \pm 0.60	76.68 \pm 1.43	79.13 \pm 2.26	89.37 \pm 0.77	81.95 \pm 3.81	79.63 \pm 2.88
nnPU-GA-c	75.88 \pm 1.16	86.08\pm2.79	77.78 \pm 0.87	80.10 \pm 1.46	54.93 \pm 8.53	79.05 \pm 1.77
PUSB	65.46 \pm 0.55	72.72 \pm 2.12	71.63 \pm 0.81	94.26\pm0.58	86.17 \pm 2.78	87.31\pm1.24
PUSB-c	73.08 \pm 0.57	78.35 \pm 1.28	75.79 \pm 1.41	82.24 \pm 0.99	74.69 \pm 1.16	80.96 \pm 2.74
VPU	61.32 \pm 5.55	33.33 \pm 27.22	80.12 \pm 7.64	59.47 \pm 17.51	0.07 \pm 0.06	34.73 \pm 6.80
VPU-c	74.86 \pm 1.16	75.32 \pm 0.83	77.88 \pm 1.56	81.67 \pm 0.95	82.23 \pm 0.80	77.46 \pm 0.23
Dist-PU	63.85 \pm 2.20	66.01 \pm 2.40	69.52 \pm 0.61	92.79 \pm 0.44	89.14 \pm 1.38	86.31 \pm 1.05
Dist-PU-c	67.29 \pm 1.38	81.49 \pm 4.56	72.42 \pm 1.95	84.35 \pm 1.74	43.02 \pm 11.82	76.21 \pm 2.39

Table 10: Test results (mean \pm std) of precision and recall for each algorithm on ImageNette (Case 2) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Precision			Recall		
Val metric	PA	PAUC	OA	PA	PAUC	OA
PUBN	71.79 \pm 1.10	75.30 \pm 1.17	79.28\pm1.74	82.94 \pm 1.86	77.08 \pm 3.45	74.08 \pm 2.39
PAN	51.63 \pm 0.24	70.75 \pm 3.83	64.81 \pm 2.19	93.47\pm1.86	53.01 \pm 8.18	61.43 \pm 3.03
CVIR	72.31 \pm 1.00	75.51 \pm 1.45	78.53 \pm 1.48	85.71 \pm 0.51	82.27 \pm 1.44	80.42 \pm 1.40
P3MIX-E	59.73 \pm 5.02	33.05 \pm 13.49	60.00 \pm 5.03	82.05 \pm 9.32	66.67 \pm 27.22	82.13 \pm 9.07
P3MIX-C	69.19 \pm 0.86	71.42 \pm 0.87	71.54 \pm 1.30	86.45 \pm 0.93	84.04 \pm 0.35	83.78 \pm 0.78
LBE	70.02 \pm 1.26	75.72 \pm 1.63	78.52 \pm 0.55	85.34 \pm 1.94	72.71 \pm 4.32	71.91 \pm 2.57
Count Loss	68.13 \pm 0.40	69.07 \pm 0.31	70.39 \pm 0.88	86.66 \pm 0.67	84.79 \pm 0.94	83.73 \pm 0.75
Robust-PU	70.52 \pm 2.00	77.77 \pm 2.80	78.01 \pm 1.20	77.58 \pm 1.81	64.58 \pm 4.76	70.69 \pm 2.37
Holistic-PU	54.06 \pm 1.09	53.98 \pm 1.86	58.96 \pm 0.87	82.10 \pm 6.17	71.70 \pm 7.18	58.89 \pm 1.87
PUe	60.79 \pm 0.12	63.24 \pm 1.65	65.25 \pm 0.23	80.82 \pm 2.80	75.25 \pm 0.71	74.24 \pm 2.00
GLWS	69.84 \pm 0.87	69.39 \pm 0.41	72.88 \pm 0.68	89.59 \pm 0.58	90.12 \pm 0.64	85.49\pm1.26
uPU	57.42 \pm 2.33	65.47 \pm 2.34	66.40 \pm 2.65	84.31 \pm 3.94	71.10 \pm 5.64	68.14 \pm 6.00
uPU-c	72.56\pm2.89	79.46\pm1.02	75.39 \pm 0.26	73.42 \pm 3.51	62.21 \pm 3.91	73.83 \pm 1.50
nnPU	63.81 \pm 0.91	53.31 \pm 2.09	68.00 \pm 0.93	91.01 \pm 2.53	98.56\pm0.94	85.41 \pm 1.95
nnPU-c	72.24 \pm 0.48	76.75 \pm 2.56	74.25 \pm 0.69	78.73 \pm 3.00	68.48 \pm 6.78	73.47 \pm 3.67
nnPU-GA	67.45 \pm 1.90	76.34 \pm 0.93	77.55 \pm 1.48	85.99 \pm 2.22	72.88 \pm 4.68	72.56 \pm 2.15
nnPU-GA-c	69.67 \pm 2.08	78.80 \pm 0.62	73.81 \pm 1.46	78.90 \pm 1.89	64.65 \pm 0.71	76.22 \pm 3.06
PUSB	65.52 \pm 2.07	68.76 \pm 0.60	72.14 \pm 1.68	90.15 \pm 1.50	81.84 \pm 0.98	82.39 \pm 2.42
PUSB-c	71.30 \pm 1.31	76.51 \pm 0.32	74.72 \pm 1.46	76.28 \pm 0.92	67.73 \pm 0.61	74.87 \pm 1.76
VPU	66.96 \pm 2.61	55.05 \pm 22.59	78.42 \pm 1.24	81.55 \pm 2.30	22.49 \pm 18.22	63.98 \pm 3.62
VPU-c	71.58 \pm 0.54	75.82 \pm 0.96	75.76 \pm 1.51	86.54 \pm 1.40	81.02 \pm 0.43	80.18 \pm 0.64
Dist-PU	59.86 \pm 1.20	63.30 \pm 0.69	65.65 \pm 1.35	84.36 \pm 2.94	83.38 \pm 1.52	77.27 \pm 1.64
Dist-PU-c	58.89 \pm 4.09	66.51 \pm 1.59	68.25 \pm 0.62	83.16 \pm 7.58	72.05 \pm 2.14	63.51 \pm 5.76

Table 11: Test results (mean±std) of accuracy, AUC, and F1 score for each algorithm on Letter (Case 1) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Test ACC			AUC			Test F1		
	PA	PAUC	OA	PA	PAUC	OA	PA	PAUC	OA
PubN	88.92±1.90	89.05±2.11	89.70±1.38	94.24±1.48	94.48±1.35	94.61±1.20	89.59±1.58	89.43±1.70	89.57±1.46
PAN	49.28±0.27	48.20±0.54	52.18±1.24	47.05±2.18	55.92±0.51	46.69±2.30	65.40±0.26	65.04±0.49	42.19±17.25
CVIR	83.35±0.56	82.60±0.75	84.67±0.58	86.40±0.65	87.63±0.99	87.78±0.90	85.16±0.50	84.33±0.62	85.86±0.35
P3MIX-E	51.80±1.39	49.62±0.87	61.42±4.12	60.70±5.00	81.42±0.26	67.00±7.82	67.12±0.64	43.85±17.57	42.69±17.49
P3MIX-C	80.03±1.13	77.58±2.53	80.92±1.14	85.08±1.23	84.43±1.62	84.50±0.68	82.46±0.83	80.56±1.73	82.83±0.96
LBE	85.63±1.13	81.37±2.19	87.55±0.28	91.81±1.52	93.96±0.29	94.38±0.23	87.17±0.85	83.32±1.15	87.44±0.32
Count Loss	77.67±0.86	73.15±1.96	78.27±1.01	86.31±1.48	87.17±1.55	84.67±0.78	80.27±0.67	77.19±1.62	79.98±0.84
Robust-PU	90.02±0.67	89.17±0.33	90.63±0.31	95.30±0.29	95.51±0.32	95.91±0.31	90.20±0.61	89.09±0.66	90.58±0.32
Holistic-PU	85.80±0.99	75.22±9.45	87.32±1.27	94.12±1.36	95.72±1.49	94.74±1.64	87.14±0.83	80.97±5.61	88.17±1.02
PUe	79.50±0.24	81.83±1.08	82.00±0.78	89.77±1.07	91.42±0.98	90.88±0.50	81.54±0.21	82.32±1.64	81.95±1.08
GLWS	85.87±0.95	80.93±1.54	86.32±0.58	92.91±0.63	93.62±0.45	92.65±0.83	87.03±0.75	83.53±1.17	87.28±0.54
uPU	74.98±1.19	79.75±0.63	77.72±0.79	85.87±0.59	88.34±0.29	86.19±0.71	78.05±1.00	79.62±0.40	77.65±1.10
uPU-c	92.23±0.26	85.97±4.01	92.73±0.15	96.84±0.15	97.26±0.05	96.40±0.18	92.18±0.14	83.09±6.02	92.60±0.14
nnPU	85.13±0.46	79.53±1.62	85.60±0.31	94.16±0.51	95.44±0.44	94.49±0.66	86.19±0.37	82.46±1.10	85.85±0.40
nnPU-c	91.87±0.34	89.25±1.14	91.82±0.14	96.15±0.30	96.39±0.69	96.36±0.38	91.85±0.25	88.24±1.71	91.58±0.21
nnPU-GA	85.12±0.13	82.85±0.68	84.27±0.58	93.17±0.44	93.56±0.61	91.18±0.41	85.74±0.25	82.86±1.69	84.46±0.64
nnPU-GA-c	90.97±0.30	88.60±0.57	90.97±0.30	94.72±0.23	96.37±1.16	94.72±0.23	90.86±0.25	87.75±0.25	90.86±0.25
PUSB	85.73±0.70	87.43±0.21	86.82±0.54	86.09±0.63	87.42±0.25	86.81±0.56	86.63±0.67	87.66±0.50	86.70±0.78
PUSB-c	91.42±0.86	90.68±0.58	91.43±0.92	91.45±0.87	90.66±0.57	91.46±0.92	91.35±1.02	90.30±0.67	91.29±1.04
VPU	89.85±1.07	67.88±8.64	90.13±0.77	95.67±0.40	96.03±0.77	95.44±0.57	89.69±0.98	44.13±20.00	89.86±0.67
VPU-c	91.83±0.54	90.28±0.98	92.15±0.52	96.32±0.38	97.06±0.26	96.96±0.30	91.95±0.42	89.38±1.17	91.93±0.48
Dist-PU	77.07±0.77	77.45±0.78	77.55±0.78	81.95±1.07	82.71±1.23	82.07±1.66	80.15±0.45	79.68±0.14	80.07±0.40
Dist-PU-c	67.65±2.41	69.33±2.52	70.03±2.28	72.96±2.78	75.75±2.56	74.72±2.66	72.61±1.10	68.78±2.64	72.81±2.05

Table 12: Test results (mean±std) of precision and recall for each algorithm on Letter (Case 1) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Precision			Recall		
	PA	PAUC	OA	PA	PAUC	OA
PubN	84.12±2.75	86.86±4.01	88.33±1.64	96.00±0.44	92.79±2.16	90.85±1.27
PAN	49.11±0.23	48.20±0.54	34.01±13.90	97.89±1.15	100.00±0.00	56.40±23.56
CVIR	75.82±0.92	74.75±0.73	77.28±0.52	97.17±0.70	96.73±0.42	96.58±0.14
P3MIX-E	50.79±0.93	65.70±14.00	45.23±18.89	99.04±0.74	66.80±27.10	42.64±18.48
P3MIX-C	73.49±1.22	71.37±2.78	75.15±0.90	94.00±1.10	92.77±0.92	92.26±1.02
LBE	78.58±1.20	75.84±4.43	85.17±1.79	97.89±0.29	93.97±3.87	90.14±2.36
Count Loss	69.56±1.08	64.88±1.99	72.33±0.83	94.95±0.69	95.42±1.62	89.44±0.97
Robust-PU	87.94±0.82	86.68±1.06	90.32±0.77	92.60±0.67	91.84±2.47	90.86±0.32
Holistic-PU	79.36±1.07	71.38±8.69	82.39±2.08	96.62±0.46	96.74±1.50	94.99±1.21
PUe	73.33±0.21	78.59±0.74	80.47±0.23	91.82±0.32	86.90±4.26	83.62±2.52
GLWS	78.56±1.32	72.32±1.87	79.44±0.61	97.60±0.27	98.98±0.30	96.84±0.48
uPU	67.24±1.28	77.65±2.09	74.98±0.74	93.05±0.71	81.96±1.61	80.56±1.73
uPU-c	89.10±0.60	93.95±2.97	92.00±0.92	95.49±0.48	77.42±10.67	93.26±0.91
nnPU	79.27±0.28	70.89±1.82	81.97±1.57	94.44±0.61	98.70±0.39	90.42±2.48
nnPU-c	90.57±0.67	93.06±1.83	93.26±1.26	93.18±0.49	84.51±4.36	90.08±1.46
nnPU-GA	81.01±0.13	81.22±3.50	80.81±1.19	91.07±0.63	86.51±6.63	88.48±0.53
nnPU-GA-c	89.60±0.51	92.37±3.02	89.60±0.51	92.17±0.32	84.05±2.57	92.17±0.32
PUSB	79.01±1.09	84.95±0.84	84.96±1.33	95.94±0.96	90.70±2.01	88.78±2.76
PUSB-c	90.00±1.00	90.13±0.90	89.87±0.76	92.79±1.59	90.49±0.97	92.77±1.40
VPU	89.11±1.55	65.54±26.76	91.24±1.24	90.42±1.99	35.30±17.59	88.61±1.31
VPU-c	88.60±0.48	95.21±0.75	92.09±0.68	95.57±0.65	84.40±2.68	91.84±1.42
Dist-PU	70.13±0.47	72.05±1.16	71.50±1.04	93.51±0.55	89.28±1.49	91.14±1.75
Dist-PU-c	63.03±3.53	68.85±3.51	66.07±3.28	87.01±3.71	68.80±1.90	81.54±2.35

Table 13: Test results (mean \pm std) of accuracy, AUC, and F1 score for each algorithm on Letter (Case 2) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Test ACC			AUC			Test F1		
	PA	PAUC	OA	PA	PAUC	OA	PA	PAUC	OA
PubN	87.47 \pm 0.58	88.98 \pm 1.45	89.63 \pm 0.98	94.15 \pm 1.14	93.88 \pm 1.45	94.59 \pm 1.09	88.40 \pm 0.61	89.15 \pm 1.44	89.74 \pm 1.02
PAN	50.02 \pm 0.48	49.88 \pm 0.85	51.73 \pm 1.43	45.39 \pm 4.63	57.60 \pm 2.16	51.85 \pm 4.33	66.64 \pm 0.40	44.18 \pm 18.05	21.43 \pm 17.50
CVIR	84.83 \pm 0.73	84.22 \pm 0.89	84.72 \pm 0.76	88.63 \pm 1.49	88.18 \pm 0.65	88.67 \pm 1.62	86.57 \pm 0.61	85.38 \pm 0.87	86.38 \pm 0.63
P3MIX-E	55.70 \pm 2.92	55.57 \pm 2.96	65.08 \pm 3.09	71.43 \pm 4.39	81.48 \pm 2.05	71.19 \pm 3.66	68.60 \pm 0.97	52.06 \pm 13.86	64.22 \pm 0.99
P3MIX-C	81.80 \pm 2.04	80.70 \pm 2.16	83.32 \pm 2.22	89.68 \pm 2.56	90.09 \pm 2.58	88.23 \pm 3.66	83.89 \pm 1.46	83.35 \pm 1.46	83.46 \pm 2.56
LBE	87.32 \pm 0.50	80.82 \pm 3.97	88.18 \pm 0.96	94.51 \pm 0.18	94.43 \pm 0.11	95.34 \pm 0.57	88.44 \pm 0.39	82.61 \pm 2.51	88.65 \pm 1.08
Count Loss	81.35 \pm 0.64	82.20 \pm 1.16	82.93 \pm 0.85	90.22 \pm 0.71	90.35 \pm 0.73	90.08 \pm 1.03	83.36 \pm 0.36	83.19 \pm 0.35	83.30 \pm 0.37
Robust-PU	90.88 \pm 0.52	90.18 \pm 0.84	91.07 \pm 0.39	96.31 \pm 0.65	96.92\pm0.53	96.63 \pm 0.57	91.00 \pm 0.44	89.64 \pm 1.06	90.83 \pm 0.42
Holistic-PU	87.88 \pm 1.37	86.12 \pm 1.82	88.65 \pm 1.12	95.09 \pm 0.62	95.36 \pm 0.69	95.40 \pm 0.83	88.79 \pm 0.90	87.58 \pm 1.22	89.49 \pm 0.85
PUc	79.50 \pm 0.70	78.03 \pm 1.17	82.53 \pm 0.04	88.18 \pm 1.99	91.92 \pm 0.19	90.65 \pm 0.40	80.97 \pm 0.48	80.73 \pm 0.63	81.94 \pm 0.20
GLWS	86.27 \pm 0.43	79.88 \pm 1.42	88.18 \pm 0.67	93.00 \pm 0.61	94.46 \pm 0.33	93.75 \pm 0.14	87.50 \pm 0.43	82.97 \pm 0.76	89.07 \pm 0.51
uPU	75.22 \pm 1.07	72.03 \pm 2.17	77.52 \pm 0.34	84.07 \pm 1.04	85.49 \pm 0.25	85.72 \pm 0.35	77.80 \pm 0.45	75.49 \pm 0.59	77.93 \pm 0.57
uPU-c	91.32 \pm 0.57	89.72 \pm 0.59	92.13 \pm 0.18	96.48 \pm 0.25	96.86 \pm 0.22	96.30 \pm 0.53	91.71 \pm 0.35	89.05 \pm 0.91	92.14 \pm 0.27
nnPU	84.68 \pm 0.34	75.22 \pm 2.11	87.38 \pm 0.39	94.02 \pm 0.74	95.30 \pm 0.51	95.34 \pm 0.24	85.90 \pm 0.44	79.94 \pm 1.42	87.73 \pm 0.39
nnPU-c	91.27 \pm 0.43	90.50 \pm 0.17	91.65 \pm 0.19	96.21 \pm 0.40	96.92\pm0.12	97.09\pm0.22	91.44 \pm 0.44	90.42 \pm 0.19	91.65 \pm 0.22
nnPU-GA	85.63 \pm 0.60	83.43 \pm 1.40	86.15 \pm 0.13	93.84 \pm 0.34	93.79 \pm 0.09	93.68 \pm 0.02	86.37 \pm 0.67	85.00 \pm 1.00	86.42 \pm 0.41
nnPU-GA-c	91.55 \pm 0.33	89.28 \pm 1.89	91.70 \pm 0.39	96.79\pm0.42	96.65 \pm 0.43	96.61 \pm 0.56	91.58 \pm 0.37	88.44 \pm 2.57	91.69 \pm 0.41
PUSB	87.42 \pm 0.31	87.83 \pm 0.13	87.63 \pm 0.24	87.39 \pm 0.34	87.85 \pm 0.13	87.61 \pm 0.23	88.15 \pm 0.18	88.30 \pm 0.24	87.98 \pm 0.44
PUSB-c	91.33 \pm 0.77	91.48\pm0.40	91.53 \pm 0.71	91.34 \pm 0.76	91.46 \pm 0.41	91.47 \pm 0.76	91.29 \pm 0.84	91.23\pm0.51	91.22 \pm 0.95
VPU	90.85 \pm 0.28	74.93 \pm 6.54	91.18 \pm 0.08	96.26 \pm 0.24	95.91 \pm 0.10	96.23 \pm 0.26	90.60 \pm 0.36	64.86 \pm 10.97	90.98 \pm 0.10
VPU-c	91.95\pm0.38	89.55 \pm 0.05	92.85\pm0.29	96.63 \pm 0.26	96.40 \pm 0.48	96.89 \pm 0.12	91.94\pm0.33	89.13 \pm 0.38	92.74\pm0.27
Dist-PU	78.92 \pm 0.89	77.52 \pm 0.51	79.42 \pm 0.71	84.82 \pm 0.29	84.29 \pm 0.73	85.17 \pm 0.34	81.73 \pm 0.94	79.36 \pm 0.51	81.10 \pm 0.82
Dist-PU-c	75.33 \pm 1.22	77.58 \pm 0.65	76.87 \pm 0.77	82.69 \pm 0.74	84.55 \pm 0.23	83.73 \pm 0.43	78.14 \pm 1.03	77.51 \pm 1.23	78.00 \pm 1.36

Table 14: Test results (mean \pm std) of precision and recall for each algorithm on Letter (Case 2) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Precision			Recall		
	PA	PAUC	OA	PA	PAUC	OA
PubN	81.40 \pm 0.45	87.45 \pm 1.99	87.75 \pm 1.02	96.72 \pm 0.84	90.94 \pm 0.83	91.83 \pm 1.24
PAN	50.04 \pm 0.51	33.05 \pm 13.52	18.13 \pm 14.80	99.74\pm0.21	66.67 \pm 27.22	26.22 \pm 21.41
CVIR	78.35 \pm 0.93	79.00 \pm 0.54	78.58 \pm 1.09	96.73 \pm 0.14	92.88 \pm 1.34	95.95\pm0.56
P3MIX-E	53.00 \pm 1.72	69.22 \pm 12.66	66.29 \pm 4.29	97.68 \pm 1.77	67.78 \pm 23.64	63.34 \pm 2.72
P3MIX-C	75.69 \pm 1.88	73.75 \pm 1.92	82.40 \pm 1.72	94.13 \pm 0.75	95.91 \pm 1.19	85.25 \pm 5.36
LBE	81.36 \pm 0.74	79.16 \pm 6.57	84.70 \pm 1.08	96.89 \pm 0.35	89.61 \pm 6.60	93.03 \pm 1.51
Count Loss	75.71 \pm 1.09	80.17 \pm 3.40	82.45 \pm 2.34	92.84 \pm 1.32	87.33 \pm 3.35	84.47 \pm 1.68
Robust-PU	89.96 \pm 1.62	94.75 \pm 0.38	93.22\pm0.86	92.17 \pm 0.98	85.16 \pm 2.21	88.58 \pm 0.79
Holistic-PU	84.34 \pm 2.63	81.09 \pm 3.15	85.42 \pm 2.12	94.05 \pm 1.39	95.65 \pm 1.66	94.18 \pm 1.43
PUc	74.35 \pm 1.10	70.71 \pm 1.81	82.69 \pm 0.51	89.04 \pm 1.83	94.32 \pm 1.43	81.25 \pm 0.86
GLWS	78.89 \pm 0.68	71.48 \pm 1.33	83.68 \pm 0.94	98.23 \pm 0.11	98.94 \pm 0.40	95.21 \pm 0.38
uPU	70.01 \pm 1.43	67.60 \pm 4.04	77.28 \pm 0.22	87.70 \pm 1.18	87.40 \pm 4.78	78.61 \pm 1.09
uPU-c	87.67 \pm 0.95	94.08 \pm 1.70	92.56 \pm 0.46	96.19 \pm 0.92	84.81 \pm 2.70	91.72 \pm 0.28
nnPU	79.59 \pm 0.97	66.88 \pm 1.98	85.51 \pm 1.48	93.41 \pm 1.66	99.49\pm0.08	90.21 \pm 1.41
nnPU-c	91.19 \pm 1.12	91.85 \pm 0.63	92.71 \pm 0.80	91.80 \pm 1.68	89.05 \pm 0.60	90.66 \pm 1.05
nnPU-GA	82.09 \pm 1.25	78.07 \pm 1.98	86.18 \pm 1.35	91.33 \pm 2.21	93.44 \pm 0.65	86.84 \pm 1.84
nnPU-GA-c	91.10 \pm 0.34	93.73 \pm 1.91	91.70 \pm 0.36	92.09 \pm 0.96	84.69 \pm 5.82	91.69 \pm 0.83
PUSB	83.69 \pm 1.06	85.27 \pm 0.95	86.21 \pm 0.86	93.19 \pm 1.01	91.66 \pm 1.39	89.95 \pm 1.78
PUSB-c	89.84 \pm 0.67	92.04 \pm 0.36	92.96 \pm 0.25	92.80 \pm 1.20	90.44 \pm 0.72	89.58 \pm 1.67
VPU	92.35\pm0.79	98.33\pm0.85	92.56 \pm 0.96	89.00 \pm 1.39	51.86 \pm 13.08	89.52 \pm 0.97
VPU-c	92.02 \pm 0.61	93.71 \pm 1.97	92.76 \pm 0.53	91.85 \pm 0.15	85.28 \pm 2.42	92.75 \pm 0.68
Dist-PU	72.28 \pm 1.26	72.46 \pm 1.03	75.23 \pm 1.44	94.04 \pm 0.38	87.89 \pm 1.96	88.08 \pm 1.38
Dist-PU-c	70.19 \pm 2.12	77.82 \pm 0.81	73.93 \pm 0.49	88.32 \pm 0.66	77.27 \pm 1.93	82.61 \pm 2.50

Table 15: Test results (mean \pm std) of accuracy, AUC, and F1 score for each algorithm on USPS (Case 1) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Test ACC			AUC			Test F1		
	PA	PAUC	OA	PA	PAUC	OA	PA	PAUC	OA
PubN	93.76\pm0.23	92.89 \pm 0.24	93.95\pm0.13	98.28\pm0.03	98.01 \pm 0.08	98.29\pm0.04	92.60\pm0.28	91.42 \pm 0.36	92.77\pm0.17
PAN	84.52 \pm 0.50	84.52 \pm 0.65	84.97 \pm 0.48	89.98 \pm 0.20	90.89 \pm 0.48	90.06 \pm 0.46	81.23 \pm 0.46	80.56 \pm 1.19	80.61 \pm 0.58
CVIR	82.79 \pm 1.48	82.01 \pm 0.96	82.98 \pm 1.39	94.88 \pm 0.36	93.80 \pm 0.16	93.42 \pm 0.69	82.72 \pm 1.31	81.95 \pm 0.79	82.76 \pm 1.26
P3MIX-E	88.99 \pm 1.40	89.49 \pm 1.29	89.84 \pm 1.17	96.18 \pm 0.44	96.33 \pm 0.43	96.23 \pm 0.47	87.54 \pm 1.30	88.02 \pm 1.23	87.90 \pm 1.30
P3MIX-C	92.69 \pm 0.66	93.47\pm0.49	93.22 \pm 0.31	97.98 \pm 0.22	98.16\pm0.14	98.09 \pm 0.11	91.41 \pm 0.78	92.38\pm0.57	92.05 \pm 0.36
LBE	91.45 \pm 0.62	87.10 \pm 1.25	92.29 \pm 0.33	97.67 \pm 0.12	97.04 \pm 0.46	97.60 \pm 0.18	90.52 \pm 0.55	86.49 \pm 1.17	91.16 \pm 0.18
Count Loss	91.99 \pm 0.34	90.08 \pm 0.84	91.76 \pm 0.81	97.44 \pm 0.27	97.27 \pm 0.09	97.60 \pm 0.09	90.97 \pm 0.31	88.91 \pm 0.66	90.64 \pm 0.69
Robust-PU	91.73 \pm 0.27	88.19 \pm 3.28	92.79 \pm 0.12	97.51 \pm 0.20	97.48 \pm 0.22	97.73 \pm 0.15	89.88 \pm 0.20	83.74 \pm 5.36	91.20 \pm 0.14
Holistic-PU	91.94 \pm 0.82	92.56 \pm 0.11	93.46 \pm 0.36	97.22 \pm 0.34	97.47 \pm 0.17	97.76 \pm 0.16	90.88 \pm 0.84	91.12 \pm 0.02	92.27 \pm 0.40
PUe	84.82 \pm 1.01	84.22 \pm 0.30	86.93 \pm 0.27	95.41 \pm 0.12	95.25 \pm 0.13	94.40 \pm 1.03	84.23 \pm 0.76	83.60 \pm 0.19	85.24 \pm 0.59
GLWS	91.13 \pm 0.37	86.78 \pm 0.60	90.52 \pm 0.47	98.21 \pm 0.02	97.78 \pm 0.17	98.18 \pm 0.09	90.40 \pm 0.36	86.38 \pm 0.55	89.81 \pm 0.45
uPU	83.14 \pm 0.93	83.87 \pm 0.11	83.86 \pm 0.83	92.88 \pm 0.15	93.10 \pm 0.18	93.01 \pm 0.05	81.51 \pm 0.81	81.98 \pm 0.19	82.04 \pm 0.70
uPU-c	93.44 \pm 0.26	91.30 \pm 1.16	93.32 \pm 0.10	97.95 \pm 0.12	97.79 \pm 0.11	97.85 \pm 0.09	92.05 \pm 0.34	88.94 \pm 1.72	91.94 \pm 0.16
nnPU	90.60 \pm 0.28	87.49 \pm 0.81	90.22 \pm 0.42	97.94 \pm 0.09	97.63 \pm 0.06	97.70 \pm 0.15	89.82 \pm 0.27	87.02 \pm 0.73	89.44 \pm 0.42
nnPU-c	92.64 \pm 0.08	90.82 \pm 0.94	93.24 \pm 0.18	97.60 \pm 0.05	97.34 \pm 0.17	97.99 \pm 0.03	91.03 \pm 0.12	88.41 \pm 1.37	91.76 \pm 0.23
nnPU-GA	91.28 \pm 0.16	92.46 \pm 0.11	92.51 \pm 0.31	96.79 \pm 0.10	97.41 \pm 0.11	97.17 \pm 0.27	89.80 \pm 0.36	91.09 \pm 0.07	91.30 \pm 0.35
nnPU-GA-c	92.76 \pm 0.38	90.60 \pm 1.44	92.79 \pm 0.22	97.58 \pm 0.05	97.46 \pm 0.16	97.66 \pm 0.11	91.23 \pm 0.55	88.05 \pm 2.15	91.33 \pm 0.32
PUSB	89.90 \pm 0.73	91.73 \pm 0.26	91.38 \pm 0.83	90.98 \pm 0.65	92.51 \pm 0.26	92.17 \pm 0.70	89.17 \pm 0.71	90.91 \pm 0.29	90.56 \pm 0.80
PUSB-c	92.91 \pm 0.30	92.84 \pm 0.24	92.83 \pm 0.18	92.30 \pm 0.29	92.26 \pm 0.27	92.25 \pm 0.29	91.34 \pm 0.35	91.28 \pm 0.31	91.26 \pm 0.28
VPU	88.14 \pm 2.21	57.71 \pm 0.04	89.89 \pm 1.71	92.98 \pm 3.98	97.31 \pm 0.13	97.76 \pm 0.19	84.36 \pm 3.09	0.31 \pm 0.17	86.58 \pm 2.62
VPU-c	92.92 \pm 0.07	80.17 \pm 7.36	93.29 \pm 0.32	97.55 \pm 0.13	97.82 \pm 0.24	97.79 \pm 0.18	91.40 \pm 0.08	63.80 \pm 17.92	91.97 \pm 0.38
Dist-PU	87.73 \pm 0.55	82.15 \pm 2.23	86.10 \pm 0.14	92.52 \pm 0.85	92.58 \pm 0.43	91.03 \pm 0.77	86.69 \pm 0.55	81.64 \pm 1.82	84.77 \pm 0.17
Dist-PU-c	92.01 \pm 0.19	90.47 \pm 0.77	91.50 \pm 0.34	97.92 \pm 0.16	97.95 \pm 0.21	97.74 \pm 0.21	90.16 \pm 0.22	87.84 \pm 1.11	89.44 \pm 0.44

Table 16: Test results (mean \pm std) of precision and recall for each algorithm on USPS (Case 1) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Precision			Recall		
	PA	PAUC	OA	PA	PAUC	OA
PubN	92.93 \pm 0.28	93.46 \pm 0.72	93.93 \pm 0.10	92.27 \pm 0.31	89.53 \pm 1.26	91.65 \pm 0.35
PAN	84.90 \pm 5.07	87.64 \pm 5.36	90.59 \pm 5.35	79.49 \pm 4.77	76.59 \pm 6.00	74.27 \pm 5.22
CVIR	72.19 \pm 1.80	71.37 \pm 1.25	72.62 \pm 1.66	96.90 \pm 0.42	96.27 \pm 0.23	96.27 \pm 0.86
P3MIX-E	85.21 \pm 3.71	86.04 \pm 3.34	89.31 \pm 3.14	90.47 \pm 1.33	90.43 \pm 1.08	87.02 \pm 2.80
P3MIX-C	90.96 \pm 0.73	91.40 \pm 0.59	91.48 \pm 0.43	91.88 \pm 0.91	93.37 \pm 0.55	92.63 \pm 0.31
LBE	85.53 \pm 1.52	78.00 \pm 1.86	89.02 \pm 2.12	96.24 \pm 0.81	97.18 \pm 1.13	93.69 \pm 1.91
Count Loss	87.21 \pm 1.11	85.10 \pm 2.53	88.05 \pm 2.45	95.14 \pm 0.74	93.45 \pm 1.97	93.65 \pm 1.43
Robust-PU	93.49 \pm 1.50	95.16 \pm 0.47	94.46 \pm 0.25	86.63 \pm 0.95	75.96 \pm 8.06	88.16 \pm 0.16
Holistic-PU	87.61 \pm 1.66	92.29 \pm 1.53	92.36 \pm 0.77	94.47 \pm 0.28	90.12 \pm 1.45	92.20 \pm 0.26
PUe	75.50 \pm 1.85	74.71 \pm 0.58	81.82 \pm 1.80	95.45 \pm 1.08	94.90 \pm 0.45	89.41 \pm 3.20
GLWS	83.50 \pm 0.64	76.68 \pm 0.78	82.47 \pm 0.75	98.55\pm0.12	98.90\pm0.19	98.59\pm0.06
uPU	76.29 \pm 1.56	77.81 \pm 0.18	77.77 \pm 1.65	87.57 \pm 0.17	86.63 \pm 0.61	86.90 \pm 0.47
uPU-c	94.51 \pm 0.21	95.48 \pm 0.55	94.10 \pm 0.32	89.73 \pm 0.68	83.45 \pm 3.19	89.88 \pm 0.60
nnPU	83.01 \pm 0.49	77.78 \pm 1.21	82.43 \pm 0.66	97.84 \pm 0.22	98.78 \pm 0.14	97.76 \pm 0.36
nnPU-c	94.11 \pm 0.36	94.52 \pm 0.71	94.85 \pm 0.20	88.16 \pm 0.50	83.18 \pm 2.52	88.86 \pm 0.44
nnPU-GA	89.02 \pm 1.35	91.22 \pm 0.96	89.89 \pm 0.46	90.78 \pm 2.06	91.02 \pm 0.92	92.75 \pm 0.37
nnPU-GA-c	93.49 \pm 0.40	94.18 \pm 0.39	93.02 \pm 0.36	89.14 \pm 1.38	82.98 \pm 3.91	89.73 \pm 0.94
PUSB	81.80 \pm 1.12	85.07 \pm 0.28	84.72 \pm 1.55	98.04 \pm 0.28	97.61 \pm 0.28	97.33 \pm 0.42
PUSB-c	94.59\pm0.58	94.29 \pm 0.46	94.23 \pm 0.58	88.31 \pm 0.42	88.47 \pm 0.62	88.51 \pm 1.03
VPU	94.38 \pm 2.09	66.67 \pm 27.22	97.19\pm0.39	76.55 \pm 4.48	0.16 \pm 0.08	78.43 \pm 4.44
VPU-c	94.22 \pm 0.79	96.82\pm0.89	93.30 \pm 0.50	88.78 \pm 0.76	55.53 \pm 18.31	90.67 \pm 0.28
Dist-PU	80.19 \pm 0.78	73.27 \pm 3.53	79.18 \pm 1.24	94.35 \pm 0.22	92.71 \pm 1.44	91.41 \pm 1.93
Dist-PU-c	94.24 \pm 0.58	95.22 \pm 0.37	94.27 \pm 0.38	86.43 \pm 0.43	81.61 \pm 2.06	85.10 \pm 0.71

Table 17: Test results (mean \pm std) of accuracy, AUC, and F1 score for each algorithm on USPS (Case 2) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Test ACC			AUC			Test F1		
	PA	PAUC	OA	PA	PAUC	OA	PA	PAUC	OA
PUBN	94.45 \pm 0.26	95.45\pm0.14	95.45\pm0.20	98.62 \pm 0.18	98.83 \pm 0.08	98.85 \pm 0.10	94.23 \pm 0.29	95.31\pm0.13	95.29\pm0.20
PAN	80.70 \pm 3.24	83.54 \pm 1.27	84.22 \pm 1.12	88.16 \pm 3.61	92.49 \pm 0.23	92.89 \pm 0.18	78.89 \pm 3.82	81.47 \pm 2.11	82.45 \pm 1.90
CVIR	90.93 \pm 0.26	88.57 \pm 0.29	90.55 \pm 0.20	96.74 \pm 0.23	96.34 \pm 0.22	96.69 \pm 0.21	91.37 \pm 0.23	89.34 \pm 0.24	91.05 \pm 0.16
P3MIX-E	94.04 \pm 0.43	93.90 \pm 0.43	93.90 \pm 0.39	98.26 \pm 0.27	98.24 \pm 0.26	98.14 \pm 0.19	93.92 \pm 0.39	93.79 \pm 0.38	93.81 \pm 0.36
P3MIX-C	94.27 \pm 0.52	94.54 \pm 0.51	94.72 \pm 0.35	98.38 \pm 0.24	98.56 \pm 0.17	98.66 \pm 0.13	94.20 \pm 0.48	94.47 \pm 0.48	94.62 \pm 0.34
LBE	94.67 \pm 0.20	90.82 \pm 1.28	94.88 \pm 0.05	98.51 \pm 0.16	98.17 \pm 0.08	98.60 \pm 0.07	94.65 \pm 0.16	91.18 \pm 0.98	94.73 \pm 0.10
Count Loss	92.73 \pm 0.22	93.76 \pm 0.45	93.17 \pm 0.14	97.15 \pm 0.27	97.72 \pm 0.20	97.33 \pm 0.26	92.87 \pm 0.16	93.84 \pm 0.40	93.18 \pm 0.14
Robust-PU	93.72 \pm 0.41	93.64 \pm 0.31	94.93 \pm 0.19	98.13 \pm 0.11	98.34 \pm 0.15	98.62 \pm 0.23	93.44 \pm 0.45	93.33 \pm 0.31	94.69 \pm 0.20
Holistic-PU	95.15\pm0.28	94.83 \pm 0.24	95.02 \pm 0.53	98.76 \pm 0.11	98.73 \pm 0.17	98.49 \pm 0.20	94.99\pm0.29	94.65 \pm 0.24	94.84 \pm 0.57
PUe	85.27 \pm 1.11	85.00 \pm 0.63	86.05 \pm 0.33	93.95 \pm 0.48	95.26 \pm 0.23	93.48 \pm 0.92	86.55 \pm 0.96	86.38 \pm 0.50	87.08 \pm 0.25
GLWS	92.48 \pm 0.50	88.19 \pm 0.44	92.18 \pm 0.26	98.58 \pm 0.05	98.09 \pm 0.29	98.45 \pm 0.06	92.76 \pm 0.44	89.15 \pm 0.35	92.49 \pm 0.23
uPU	83.36 \pm 0.48	82.68 \pm 0.81	84.12 \pm 0.06	92.33 \pm 0.30	93.99 \pm 0.25	92.79 \pm 0.79	84.51 \pm 0.42	84.34 \pm 0.56	85.14 \pm 0.27
uPU-c	94.67 \pm 0.10	93.36 \pm 0.76	94.57 \pm 0.28	98.78\pm0.10	98.50 \pm 0.30	98.68 \pm 0.11	94.36 \pm 0.11	92.85 \pm 0.88	94.26 \pm 0.32
nnPU	93.64 \pm 1.13	88.01 \pm 1.30	94.10 \pm 0.37	98.50 \pm 0.15	98.37 \pm 0.09	98.71 \pm 0.08	93.79 \pm 1.03	89.01 \pm 1.07	94.20 \pm 0.33
nnPU-c	94.32 \pm 0.20	94.00 \pm 0.24	94.80\pm0.12	98.69 \pm 0.02	98.48 \pm 0.04	98.67 \pm 0.04	94.03 \pm 0.23	93.67 \pm 0.27	94.56 \pm 0.10
nnPU-GA	94.42 \pm 0.26	94.95 \pm 0.13	94.78 \pm 0.07	98.61 \pm 0.10	98.68 \pm 0.10	98.44 \pm 0.12	94.28 \pm 0.24	94.76 \pm 0.14	94.59 \pm 0.08
nnPU-GA-c	94.12 \pm 0.04	94.27 \pm 0.17	94.07 \pm 0.36	98.66 \pm 0.06	98.79 \pm 0.07	98.73 \pm 0.07	93.77 \pm 0.04	93.92 \pm 0.17	93.69 \pm 0.40
PUSB	92.41 \pm 0.67	91.96 \pm 1.50	93.56 \pm 0.34	92.57 \pm 0.65	92.10 \pm 1.45	93.68 \pm 0.33	92.69 \pm 0.59	92.25 \pm 1.30	93.70 \pm 0.30
PUSB-c	94.09 \pm 0.18	93.64 \pm 0.21	94.57 \pm 0.25	94.01 \pm 0.18	93.55 \pm 0.22	94.50 \pm 0.25	93.76 \pm 0.19	93.24 \pm 0.24	94.27 \pm 0.26
VPU	89.82 \pm 2.61	76.63 \pm 10.29	89.54 \pm 1.88	97.91 \pm 0.35	98.58 \pm 0.10	97.99 \pm 0.57	88.33 \pm 3.43	58.65 \pm 23.78	88.11 \pm 2.41
VPU-c	94.93 \pm 0.13	94.82 \pm 0.12	95.05 \pm 0.14	98.78\pm0.06	98.86\pm0.09	98.91\pm0.11	94.72 \pm 0.17	94.55 \pm 0.16	94.81 \pm 0.14
Dist-PU	94.82 \pm 0.12	94.02 \pm 0.18	94.72 \pm 0.38	98.09 \pm 0.17	97.94 \pm 0.28	98.12 \pm 0.19	94.63 \pm 0.13	93.72 \pm 0.25	94.55 \pm 0.39
Dist-PU-c	94.10 \pm 0.44	92.09 \pm 0.50	94.14 \pm 0.37	98.49 \pm 0.16	98.53 \pm 0.20	98.50 \pm 0.03	93.73 \pm 0.49	91.30 \pm 0.59	93.77 \pm 0.42

Table 18: Test results (mean \pm std) of precision and recall for each algorithm on USPS (Case 2) under different model selection criteria. The best performance w.r.t. each validation metric is shown in bold. Here, “-c” indicates using the proposed calibration technique in Algorithm 1.

Test metric	Precision			Recall		
	PA	PAUC	OA	PA	PAUC	OA
PUBN	95.38 \pm 0.94	95.68 \pm 0.48	95.97 \pm 0.59	93.18 \pm 1.09	94.95 \pm 0.27	94.64 \pm 0.42
PAN	83.76 \pm 3.59	89.42 \pm 1.97	89.52 \pm 2.59	75.13 \pm 5.70	75.67 \pm 5.30	77.35 \pm 5.13
CVIR	85.15 \pm 0.42	81.84 \pm 0.41	84.48 \pm 0.37	98.57 \pm 0.05	98.36 \pm 0.22	98.74 \pm 0.18
P3MIX-E	93.52 \pm 1.27	93.30 \pm 1.32	92.82 \pm 0.87	94.37 \pm 0.51	94.34 \pm 0.58	94.85 \pm 0.29
P3MIX-C	93.07 \pm 1.19	93.28 \pm 0.94	93.88 \pm 0.75	95.39 \pm 0.26	95.70 \pm 0.08	95.39 \pm 0.10
LBE	92.57 \pm 0.82	86.94 \pm 3.54	95.05 \pm 0.89	96.86 \pm 0.56	96.45 \pm 2.16	94.47 \pm 1.09
Count Loss	88.88 \pm 0.75	90.40 \pm 0.96	90.75 \pm 0.87	97.27 \pm 0.57	97.58 \pm 0.41	95.80 \pm 1.05
Robust-PU	95.00 \pm 0.46	95.31 \pm 0.48	96.73 \pm 0.26	91.95 \pm 0.83	91.44 \pm 0.17	92.73 \pm 0.39
Holistic-PU	95.52 \pm 0.49	95.53 \pm 0.63	95.60 \pm 0.47	94.47 \pm 0.42	93.79 \pm 0.31	94.10 \pm 0.92
PUe	77.96 \pm 1.20	77.46 \pm 0.78	79.31 \pm 0.72	97.30 \pm 0.88	97.65 \pm 0.21	96.59 \pm 0.91
GLWS	87.39 \pm 0.88	80.72 \pm 0.66	86.83 \pm 0.38	98.84\pm0.12	99.56\pm0.20	98.94\pm0.06
uPU	77.29 \pm 0.61	75.45 \pm 1.20	78.20 \pm 0.75	93.24 \pm 0.79	95.67 \pm 0.58	93.55 \pm 1.69
uPU-c	97.19\pm0.48	97.26 \pm 0.65	96.97 \pm 0.16	91.71 \pm 0.50	88.88 \pm 1.65	91.71 \pm 0.70
nnPU	89.86 \pm 1.86	80.69 \pm 1.82	90.47 \pm 0.71	98.16 \pm 0.35	99.32 \pm 0.07	98.26 \pm 0.13
nnPU-c	96.32 \pm 0.13	96.26 \pm 0.18	96.46 \pm 0.51	91.85 \pm 0.53	91.23 \pm 0.63	92.73 \pm 0.29
nnPU-GA	94.11 \pm 0.67	95.82 \pm 0.28	95.61 \pm 0.20	94.47 \pm 0.26	93.72 \pm 0.44	93.59 \pm 0.28
nnPU-GA-c	96.77 \pm 0.20	97.19 \pm 0.29	97.11 \pm 0.18	90.96 \pm 0.16	90.86 \pm 0.11	90.52 \pm 0.59
PUSB	87.44 \pm 1.11	87.68 \pm 2.47	89.52 \pm 0.75	98.64 \pm 0.10	97.48 \pm 0.32	98.29 \pm 0.27
PUSB-c	96.44 \pm 0.41	96.58 \pm 0.49	96.84 \pm 0.41	91.23 \pm 0.44	90.14 \pm 0.63	91.85 \pm 0.27
VPU	96.67 \pm 0.70	98.82\pm0.56	97.26\pm0.20	82.02 \pm 5.96	52.95 \pm 21.53	80.76 \pm 3.82
VPU-c	96.09 \pm 0.48	96.87 \pm 0.51	96.80 \pm 0.21	93.42 \pm 0.78	92.36 \pm 0.79	92.90 \pm 0.07
Dist-PU	95.45 \pm 0.21	95.77 \pm 0.70	95.01 \pm 0.40	93.82 \pm 0.18	91.81 \pm 1.11	94.10 \pm 0.50
Dist-PU-c	97.04 \pm 0.15	98.27 \pm 0.26	96.98 \pm 0.06	90.65 \pm 0.83	85.26 \pm 0.88	90.79 \pm 0.80