

Article

Robust Positive-Unlabeled Learning via Bounded Loss Functions under Label Noise

Lalit Awasthi^{1,*}, Eric Danso²¹ School of Artificial Intelligence, Nanjing University of Information Science and Technology, Nanjing, 210044, China;
202452200021@nuist.edu.cn; lawasthi12@gmail.com² School of Computer Science, Nanjing University of Information Science and Technology, Nanjing, 210044, China;
202452200035@nuist.edu.cn; aericdanso98@gmail.com

* Correspondence

The author(s) received no financial support for the research, authorship, and/or publication of this article.

ABSTRACT: Positive-Unlabeled (PU) learning has become a pivotal tool in scenarios where only positive samples are labeled, and negative labels are unavailable. However, in practical applications, the labeled positive data often contains noise such as mislabeled or outlier instances that can severely degrade model performance. This issue is exacerbated using traditional surrogate loss functions, many of which are unbounded and overly sensitive to mislabeled examples. To address this limitation, we propose a robust PU learning framework that integrates bounded loss functions, including ramp loss and truncated logistic loss, into the non-negative risk estimation paradigm. Unlike conventional loss formulations that allow noisy samples to disproportionately influence training, our approach caps each instance's contribution, thereby reducing the sensitivity to label noise. We mathematically reformulate the PU risk estimator using bounded surrogates and demonstrate that this formulation maintains risk consistency while offering improved noise tolerance. A detailed framework diagram and algorithmic description are provided, along with theoretical analysis that bounds the influence of corrupted labels. Extensive experiments are conducted on both synthetic and real-world datasets under varying noise levels. Our method consistently outperforms baseline models such as unbiased PU (uPU) and non-negative PU (nnPU) in terms of classification accuracy, area under the receiver operating characteristic curve (ROC AUC), and precision-recall area under the curve (PR AUC). The ramp loss variant exhibits particularly strong robustness without sacrificing optimization efficiency. These results demonstrate that incorporating bounded losses is a principled and effective strategy for enhancing the reliability of PU learning in noisy environments.

Keywords: Positive-Unlabeled Learning, Label Noise Robustness, Bounded Loss Functions, Weak Supervision, Risk Estimation

Copyright: © 2025 by the authors. This is an open-access article under the CC-BY-SA license.



1. Introduction

Positive-Unlabeled (PU) learning is a variant of binary classification in which only a subset of positive examples is labeled, while the rest of the dataset remains unlabeled, possibly containing both positive and negative instances [1]. This learning paradigm arises naturally in many real-world scenarios, such as medical diagnosis, fraud detection, and text classification, where acquiring negative labels is often impractical or costly. In such cases, traditional supervised learning techniques become inadequate, and PU learning offers a principled alternative by

leveraging the information present in the unlabeled data to infer decision boundaries [2].

Formally, PU learning aims to estimate a binary classifier by assuming access to two disjoint subsets of data: a labeled set containing positive examples, and an unlabeled set drawn from the entire population, which includes both positive and negative samples [3]. A key assumption is that the labeled positives are randomly sampled from the true positive distribution [4]. However, in practice, this assumption may not hold due to the presence of label noise. That is, the labeled positive set may be contaminated with samples that do not actually belong to the positive class.

This contamination introduces a mismatch between the empirical and true data distributions, leading to biased risk estimates and degraded model performance.

A persistent challenge in PU learning is the vulnerability of existing algorithms to label noise, particularly when the labeled positive data contains incorrect annotations [5]. Unlike conventional supervised learning settings, where both classes may suffer from mislabeling, PU learning is more sensitive because the model heavily relies on a small, trusted set of labeled positives [6], [7]. Noise in this set can severely distort the risk estimation and misguide the classifier, leading to suboptimal performance. Notably, widely used approaches such as unbiased PU (uPU) and non-negative PU (nnPU) learning rely on loss functions that are often unbounded, making them particularly susceptible to outliers and noisy labels [8].

In this work, we propose a robust approach to PU learning by incorporating bounded loss functions into the risk estimation framework. Bounded loss functions, such as the ramp loss or modified logistic loss, inherently limit the influence of incorrectly labeled examples by capping the contribution of each instance to the overall loss [9], [10], [11]. This property is especially valuable in the PU setting, where robustness to mislabeled positive examples can significantly improve model stability and generalization. We further provide a theoretical foundation demonstrating the robustness properties of our bounded loss formulation and show, both mathematically and empirically, that our method consistently outperforms existing techniques under varying levels of label noise.

To aid in understanding the pipeline of our method, we provide a visual representation of the full framework. We also implement performance visualizations to generate clean and interpretable plots. To ensure reproducibility and clarity, the training process is formalized into an algorithmic table, and all essential components of our method are rigorously expressed using mathematical equations. This comprehensive treatment not only enhances robustness in PU learning but also sets a practical guideline for noise-aware model design.

The key contributions of this paper are as follows:

- 1) We introduce a novel PU learning framework that integrates bounded loss functions to mitigate the impact of noisy labels in positive samples.
- 2) We present a theoretical analysis to support the robustness of bounded losses under label noise.
- 3) We provide a structured algorithmic procedure and use visual diagrams to illustrate the effectiveness of the proposed approach.
- 4) We conduct extensive experiments on synthetic and real-world datasets, demonstrating superior performance compared to baseline PU learning models under varying noise conditions.

The rest of this paper is organized as follows: Section 2 reviews related work in PU learning, label noise robustness, and bounded loss functions. Section 3 introduces the necessary preliminaries and risk formulations. Section 4 describes the proposed method with theoretical justifications and algorithmic details. Section 5 presents our experimental setup, visualizations, and results. Section 6 offers a discussion on the implications and limitations of the method, and Section 7 concludes the paper with future research directions.

2. Related Work

The related work section situates our study within the broader landscape of PU learning, robustness to label noise, and loss function design. It reviews prior research that informs and contrasts with our approach, identifying their strengths and limitations in handling noisy labels, particularly within the PU learning setting. This section is structured into three main areas: existing PU learning frameworks, methods for handling label noise in PU scenarios, and the development and use of bounded loss functions in robust machine learning.

2.1. Positive-Unlabeled (PU) Learning

PU learning has received significant attention over the past decade due to its applicability in domains where only a subset of positive examples can be reliably labeled. Early efforts such as the method of Elkan and Noto introduced a probabilistic framework to estimate the true class posterior by assuming that labeled positives are selected randomly from the true positive distribution [1]. This foundational idea was later expanded in various directions.

One major branch of PU learning research involves risk estimators that approximate the true classification risk using only positive and unlabeled data. Jonathan et al. [12] proposed an unbiased PU (uPU) learning formulation based on empirical risk minimization by reweighting losses over positive and unlabeled samples. However, the uPU method suffers from high variance due to the cancellation of large loss terms, especially when the classifier is highly flexible.

To address this, Wang et al. [13] introduced the non-negative PU (nnPU) learning approach, which corrects for the high variance of uPU by imposing a non-negativity constraint on the risk estimator. This method ensures more stable training but still relies on standard surrogate loss functions like sigmoid or logistic that can be sensitive to noisy labels. These methods assume that labeled positives are clean and do not explicitly address the issue of noise within the labeled set, a gap our work aims to fill.

2.2. Label Noise in PU Learning

While extensive research exists on label noise in fully supervised settings, its treatment in PU learning remains

limited [14], [15], [16]. In real-world datasets, the assumption that all labeled positives are correctly annotated often does not hold. The impact of noisy positives is particularly severe in PU learning since the learning algorithm implicitly trusts the positive labels to calibrate the learning signal [17].

Although recent studies [18], [19] have investigated techniques like partial label correction and noise-tolerant classifiers, most PU learning methods either overlook the problem of label noise or operate under the impractical assumption that clean validation data is available. Other approaches, like PUBN (Positive-Unlabeled biased Negative), attempt to model bias in the unlabeled set but do not directly tackle noise in the positive set [20].

To the best of our knowledge, no prior PU learning method has systematically addressed robustness to noisy positive labels through loss function design. This absence motivates our introduction of bounded loss functions, which inherently dampen the influence of corrupted examples, providing a noise-robust alternative without requiring additional data or assumptions.

2.3. Bounded and Robust Loss Functions

Loss function design plays a central role in ensuring robustness in the presence of label noise. Unbounded losses such as logistic, exponential, or squared losses can assign disproportionately high penalties to misclassified or noisy points, leading to unstable gradients and overfitting.

To mitigate this, robust alternatives have been proposed. Ramp loss, Savage loss, and the symmetric loss family are notable examples that limit the influence of individual examples. Ghosh et al. [21] showed that symmetric loss functions, those satisfying $\ell(y, f(x)) = \ell(-y, -f(x))$ exhibit inherent noise tolerance under certain distributional assumptions. Li et al. [22] proposed generalized cross entropy loss, combining the benefits of mean absolute error (MAE) and categorical cross-entropy for improved robustness in deep networks.

While these losses have shown promise in standard supervised learning, their application to PU learning remains underexplored. Our work builds upon this foundation by explicitly adapting bounded loss functions to the PU setting, integrating them into the nnPU framework, and demonstrating both theoretical and empirical benefits in the presence of noisy positive labels.

3. Preliminaries

This section introduces the formal setup for Positive-Unlabeled (PU) learning and the baseline risk estimators on which our work builds. We define the key notations, assumptions, and mathematical foundations required to understand the derivation of our robust framework in subsequent sections.

Let $X \in \mathcal{X} \subseteq \mathbb{R}^d$ denote an input vector and $Y \in \{+1, -1\}$ its corresponding class label. In PU learning, we are given two datasets:

- 1) A set of labeled positive examples $\mathcal{P} = \{x_i^p\}_{i=1}^{n_p}$ drawn independent and identically distributed from the conditional distribution $p(x | Y = +1)$,
- 2) A set of unlabeled examples $\mathcal{U} = \{x_j^u\}_{j=1}^{n_u}$ drawn independent and identically distributed from the marginal distribution, $p(x) = \pi p(x | Y = +1) + (1 - \pi)p(x | Y = -1)$,

where $\pi = P(Y = +1)$ is the class prior probability, assumed to be known or estimated beforehand.

The goal is to train a binary classifier $f: \mathcal{X} \rightarrow \mathbb{R}$ that predicts the label of a new instance by minimizing the expected risk:

$$R(f) = \mathbb{E}_{(x,y)}[\ell(y, f(x))] \quad (1)$$

where $\ell: \{+1, -1\} \times \mathbb{R} \rightarrow \mathbb{R}^+$ is a surrogate loss function (e.g., logistic, hinge, or ramp).

Since negative labels are not directly observed, the true risk cannot be computed directly [23]. To overcome this, previous works proposed risk estimators that re-express the full supervised risk in terms of the observed distributions $p(x | Y = +1)$ and $p(x)$. Specifically, the unbiased PU (uPU) risk estimator is formulated as:

$$R_{uPU}(f) = \pi \cdot \mathbb{E}_{x \sim p(x|Y=+1)}[\ell(+1, f(x))] + \mathbb{E}_{x \sim p(x)}[\ell(-1, f(x))] - \pi \cdot \mathbb{E}_{x \sim p(x|Y=+1)}[\ell(-1, f(x))] \quad (2)$$

This estimator is unbiased with respect to the true risk but suffers from high variance, particularly when the model is overparameterized or the negative loss is overestimated.

To address this, the non-negative PU (nnPU) learning framework [24] modifies the risk estimator to:

$$R_{nnPU}(f) = \pi \cdot \mathbb{E}_{x \sim p(x|Y=+1)}[\ell(+1, f(x))] + \max\{0, \mathbb{E}_{x \sim p(x)}[\ell(-1, f(x))]\} - \pi \cdot \mathbb{E}_{x \sim p(x|Y=+1)}[\ell(-1, f(x))] \quad (3)$$

This ensures that the empirical risk remains non-negative during training, stabilizing learning [25], [26]. However, both uPU and nnPU rely on standard loss functions which are often unbounded, such as logistic or exponential loss. This makes them particularly vulnerable to label noise in the positive set.

In this work, we extend this framework by incorporating bounded loss functions into the risk estimator. Bounded losses naturally suppress the influence of mislabeled examples and lead to more stable optimization, especially under noisy label settings.

We also consider a label noise model in which a fraction $\eta \in [0,1]$ of the labeled positive data is corrupted, meaning that some of the examples in p belong to the negative class. Formally, the noisy positive distribution becomes:

$$\tilde{p}(x | Y = +1) = (1 - \eta) \cdot p(x | Y = +1) + \eta \cdot p(x | Y = -1) \quad (4)$$

Under this noise model, the true risk is further distorted, and the need for robust loss functions becomes even more critical [27].

The next section presents our proposed method: a novel bounded loss-based PU learning framework designed to be resilient to such noisy annotations in p , complete with theoretical justifications, algorithmic formalization, and visual explanation.

4. Proposed Method

In this section, we present our robust PU learning framework that addresses the challenge of label noise by incorporating bounded loss functions into the PU risk estimation. We begin by discussing the motivation behind bounded loss design, followed by its integration into the nnPU framework, theoretical analysis of its robustness, and a structured algorithm for implementation. We also

include a Mermaid diagram to visualize the flow of our proposed method.

4.1. Motivation for Bounded Loss Functions

Traditional surrogate loss functions commonly used in PU learning such as logistic, exponential, and squared losses are unbounded [28], [29]. While these losses offer convexity and smooth optimization, their unbounded nature allows mislabeled or noisy examples especially those with large margins to exert excessive influence on the gradient during training. This often leads to unstable convergence, overfitting, and degraded generalization performance in the presence of label noise. To address this challenge, we adopt bounded surrogate loss functions in our risk formulation. Bounded losses naturally mitigate the influence of corrupted samples by limiting the maximum penalty that any single example can contribute to the empirical risk [30]. This is particularly desirable in PU learning with noisy labels, where the positive set may be partially contaminated. In this work, we employ two specific bounded loss functions:

4.1.1. Ramp Loss

The ramp loss is a well-known non-convex surrogate loss that truncates the penalty beyond a margin. It is defined as:

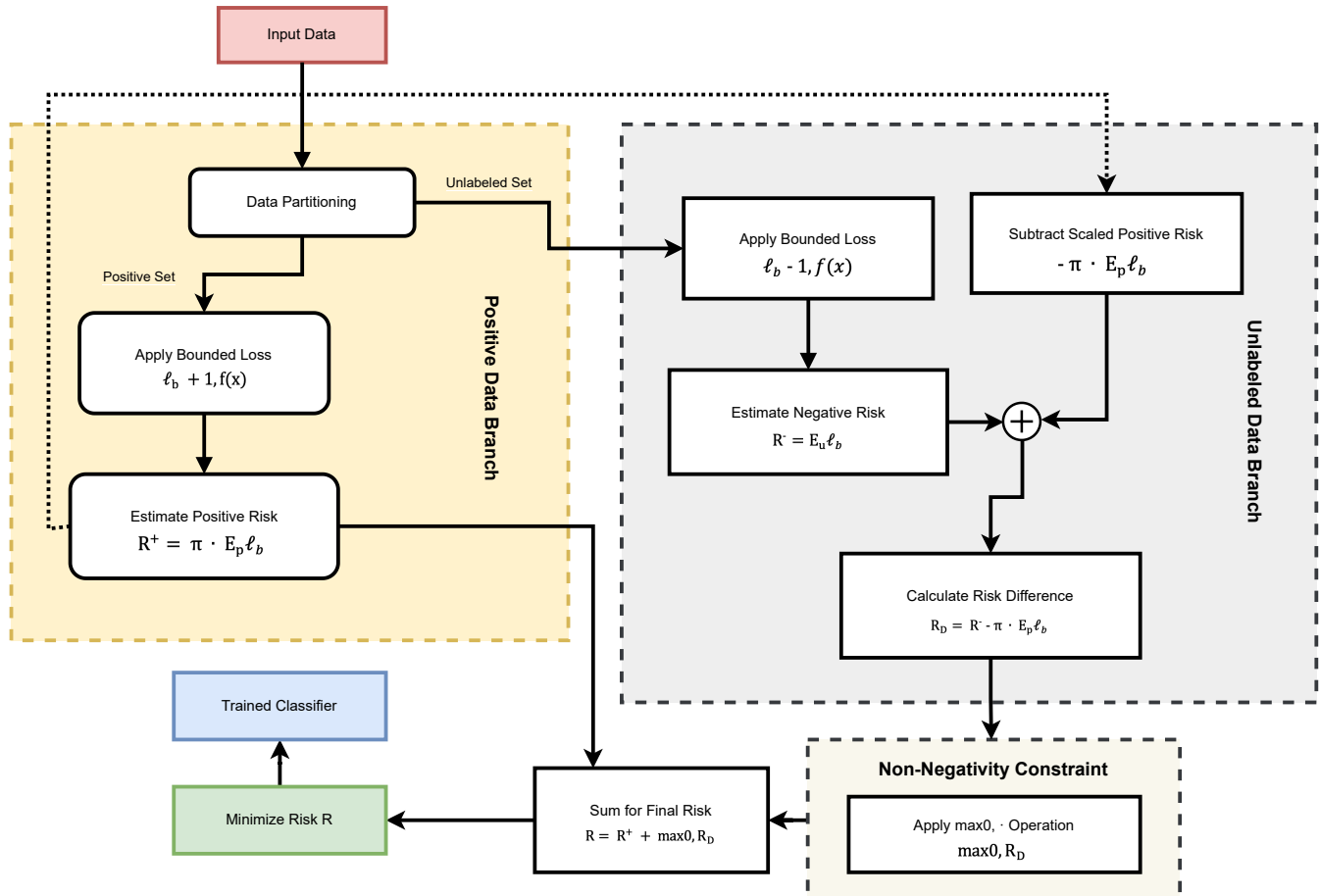


Figure 1. Proposed PU learning workflow with bounded losses. The positive and unlabeled branches are processed separately; the negative component is rectified using $\max(0, \cdot)$ before optimization.

$$\ell_{\text{ramp}}(y, f(x)) = \begin{cases} 1, & \text{if } yf(x) \leq -1 \\ 1 - yf(x), & \text{if } -1 < yf(x) < 1 \\ 0, & \text{if } yf(x) \geq 1 \end{cases} \quad (5)$$

This loss reduces sensitivity to outliers or highly uncertain examples, especially when the predicted margin is far from the decision boundary. It caps the maximum loss at 1 and has been shown to offer robust performance under adversarial and noisy conditions.

4.1.2. Truncated Logistic Loss

To maintain smoothness while ensuring boundedness, we also introduce a clipped version of the standard logistic loss, referred to as the truncated logistic loss. It is defined as:

$$\ell_{\text{trunc-log}}(y, f(x)) = \min(\log(1 + e^{-yf(x)}), \tau) \quad (6)$$

where $\tau > 0$ is a truncation threshold (we set $\tau = 2$ in our experiments). This function behaves identically to the logistic loss within a normal margin range but suppresses large gradient contributions by capping loss values when predictions are highly incorrect. This smooth, bounded form improves robustness without significantly sacrificing optimization behavior.

We generalize this idea by considering a family of bounded losses ℓ_b such that:

$$\sup_{(y, x)} \ell_b(y, f(x)) \leq C < \infty \quad (7)$$

where C is a constant that bounds the maximum loss value and is independent of the dataset. This boundedness ensures that the influence of mislabeled or outlying examples is contained throughout the learning process.

4.2. Robust PU Risk Estimation with Bounded Loss

Replacing the surrogate loss in the nnPU formulation (Equation from Section 3) with a bounded loss function ℓ_b , we define the robust risk estimator as:

$$R_{\text{robust}}(f) = \pi \cdot \mathbb{E}_{x \sim p(x|Y=+1)} [\ell_b(+1, f(x))] + \max \{0, \mathbb{E}_{x \sim p(x)} [\ell_b(-1, f(x))]\} - \pi \cdot \mathbb{E}_{x \sim p(x|Y=+1)} [\ell_b(-1, f(x))] \quad (8)$$

This estimator retains the non-negativity constraint from nnPU while ensuring that no single sample can overly skew the optimization due to the bounded nature of ℓ_b [31].

The diagram above **Figure 1** illustrates the core workflow of our proposed PU learning framework with bounded loss integration. The process begins with the input dataset, which is partitioned into two subsets: labeled positive data and unlabeled data. Each subset is processed

independently within the risk estimation pipeline. For the positive data, we directly apply a bounded surrogate loss function to compute the expected risk associated with positive labels. For the unlabeled data, we estimate its contribution to the overall risk under the assumption that it contains both positive and negative samples. This involves calculating the expected negative risk using the bounded loss, followed by subtracting a scaled estimate of negative risk derived from the positive set, weighted by the class prior π . To ensure stability during training and avoid negative risk estimates, we apply a rectification step using the $\max(0, \cdot)$ operator, as introduced in the nnPU framework. The final risk is then computed by summing the positive and corrected negative components. This total risk is minimized using standard gradient-based optimization methods, and the resulting classifier is expected to be more robust to mislabeled positives due to the bounded nature of the loss function.

4.3. Theoretical Robustness

To complement the empirical validation of our method, we provide a theoretical perspective on its robustness to label noise in the positive set. Specifically, we show that the use of bounded loss functions constrains the effect of corrupted labels on the empirical risk, making the learning process more resilient in noisy environments. Let us assume the bounded loss ℓ_b satisfies Lipschitz continuity [32], [33]:

$$|\ell_b(y, f_1(x)) - \ell_b(y, f_2(x))| \leq L \|f_1(x) - f_2(x)\| \quad \forall x \in X \quad (9)$$

Then under a bounded noise assumption $\eta < 0.5$, we can show that the risk gap between the true and noisy positive distributions is upper bounded as:

$$R_{\text{true}}^+(f) - R_{\text{noisy}}^+(f) \leq \eta C \quad (10)$$

where C is the upper bound of the loss. This bound implies that as long as the noise rate η is moderate and the loss is bounded, the influence of corrupted positives remains limited.

This result has important implications: it formally confirms that the impact of label noise is linearly proportional to the noise rate η and is capped by the maximum value of the loss function. Since bounded losses like ramp and truncated logistic loss have finite upper limits, they prevent noisy instances from dominating the risk. Thus, even when the positive labels are partially corrupted, the overall risk remains stable and controlled. In contrast, commonly used unbounded losses such as logistic or exponential loss can allow corrupted samples to exert arbitrarily large influence, leading to instability during training. The adoption of bounded losses in our method curbs outlier influence, enhancing both the robustness and generalization of the learned classifier.

Algorithm 1. Training procedure for the proposed bounded-loss PU-learning framework.

Step	Operation
1	Input: Labeled positives P , unlabeled set U , prior π , learning rate α , epochs T , model f_θ
2	Initialize model parameters θ
3	For $t = 1$ to T :
4	Sample minibatch from P and U
5	Compute bounded loss ℓ_b for each sample
6	Estimate $R^+ = \pi \cdot \mathbb{E}_{x \in P} [\ell_b(+1, f_\theta(x))]$
7	Estimate $R^- = \mathbb{E}_{x \in U} [\ell_b(-1, f_\theta(x))] - \pi \cdot \mathbb{E}_{x \in P} [\ell_b(-1, f_\theta(x))]$
8	Compute final risk: $R = R^+ + \max(0, R^-)$
9	Update θ via SGD: $\theta = \theta - \alpha \nabla_\theta R$
10	Return trained model f_θ

4.4. Algorithm Table

The training algorithm for our bounded-loss PU learning framework is grounded in the principle of risk decomposition [34], which separates the expected risk into contributions from the positive and negative classes. In the PU setting, we only have access to positively labeled examples and an unlabeled pool that contains both positive and negative instances. To estimate the overall risk, the algorithm leverages a known or estimated class prior to approximate the negative class contribution from the unlabeled data.

The key innovation lies in replacing traditional unbounded loss functions with bounded surrogates, such as ramp loss or truncated logistic loss. These bounded losses help mitigate the influence of noisy labels by capping the maximum penalty that any single sample can contribute to the loss. The training process follows a stochastic optimization scheme, where minibatches of data are sampled from the labeled and unlabeled sets, and risk is computed based on the bounded losses. Importantly, the final risk is clipped to be non-negative to avoid issues of overfitting to unreliable estimates from the unlabeled set, a technique inherited from non-negative PU learning. The steps below outline the full training pipeline used in our experiments.

The above **Algorithm 1** summarizes the computation of positive and unlabeled risks, rectification via $\max\{0, \cdot\}$, and parameter updates using gradient descent.

Where R^+ and R^- represent the positive and negative risk components respectively, and ℓ_b denotes the bounded loss function (ramp or truncated logistic).

5. Experiments

This section presents a comprehensive empirical evaluation of the proposed robust Positive-Unlabeled (PU) learning framework using bounded loss functions. Our objective is to assess the effectiveness of the proposed method compared to existing baselines, particularly under varying levels of label noise in the positive set. To ensure

a well-rounded and generalizable evaluation, we conduct experiments on both synthetic and real-world datasets, and assess performance in terms of classification accuracy, area under the receiver operating characteristic curve (ROC AUC), precision-recall area under the curve (PR AUC), and training stability.

5.1. Evaluation Metrics

We use four common metrics to assess model performance:

- 1) Accuracy: The percentage of correctly classified samples.
- 2) Precision: The proportion of correctly predicted positive samples among all predicted positives.

$$\text{Precision} = \frac{\text{True Positives}}{(\text{True Positives} + \text{False Positives})} \quad (11)$$

- 3) Recall: The proportion of correctly predicted positive samples among all actual positives.

$$\text{Recall} = \frac{\text{True Positives}}{(\text{True Positives} + \text{False Negatives})} \quad (12)$$

- 4) ROC-AUC: the area under the receiver-operating-characteristic curve that measures overall ranking ability.
- 5) PR-AUC: The area under the precision-recall curve, reflecting performance on imbalanced data.

Three datasets are used in this study. First, a synthetic dataset composed of two-class Gaussian distributions with tunable overlap and controlled noise injection is used to isolate and analyze the behavior of the learning models under clean and corrupted label conditions [35]. Second, we evaluate on the UCI Heart Disease dataset [36], which presents a binary classification task with naturally occurring label uncertainty and class imbalance. Lastly, we test our framework on a PU version of the MNIST dataset [17], [37], where digit "0" is treated as the positive class, and digits "1" through "9" are treated as unlabeled data, simulating practical PU learning scenarios in vision.

For each dataset, the class prior ($\pi = P(Y = +1)$) was estimated from the unlabeled pool using the mixture-proportion estimation (MPE) approach described by Qian et al [38]. This method infers the positive class proportion based on preliminary classifier scores. The resulting $\hat{\pi}$ values were then used in all risk computations of uPU, nnPU, and our proposed robust estimator. Small variations ± 0.05 in $\hat{\pi}$ did not noticeably affect the performance, confirming the stability of our framework.

To introduce label noise, we randomly flip a percentage of the positive samples to negative in each dataset. Noise rates are varied at four levels: 0.0, 0.2, 0.4, and 0.6. Our proposed method employs both ramp loss and

truncated logistic loss which is compared against two established PU learning approaches: unbiased PU (uPU) and non-negative PU (nnPU). All experiments are repeated five times using different random seeds, and we report the average results. Results are presented both in tables and visual plots. Across all datasets and noise levels, our

method demonstrates strong performance and stability under noise.

5.2. Results on Synthetic Dataset

The synthetic dataset provides a controlled environment to isolate the impact of label noise. **Table 2** summarizes performance on this dataset.

Table 1. Dataset summary.

Dataset	Type	Samples (n)	Features (d)	Positive Class (%)
Synthetic	Simulated Gaussian	2000	2	50
UCI Heart Disease	Medical tabular	303	13	45
MNIST (PU version)	Image dataset	10000	784	10

Table 2. Synthetic Dataset: Accuracy (%), Precision (%), Recall (%) and ROC-AUC (%) under Different Noise Rates.

Noise Rate	Method	Accuracy (%)	Precision (%)	Recall (%)	ROC-AUC (%)
0.0	uPU	85	84	86	88
0.0	nnPU	87	86	88	90
0.0	Our Method (Ramp)	89	88	90	92
0.0	Our Method (Trunc-Logistic)	88	87	89	91
0.2	uPU	77	75	78	80
0.2	nnPU	81	80	82	83
0.2	Our Method (Ramp)	84	83	85	87
0.2	Our Method (Trunc-Logistic)	83	82	84	86
0.4	uPU	66	64	67	70
0.4	nnPU	69	68	70	73
0.4	Our Method (Ramp)	78	77	79	82
0.4	Our Method (Trunc-Logistic)	76	75	77	80
0.6	uPU	52	50	54	58
0.6	nnPU	56	55	57	61
0.6	Our Method (Ramp)	69	68	70	74
0.6	Our Method (Trunc-Logistic)	67	66	68	72

Table 3. UCI Heart Disease Dataset: Accuracy (%), Precision (%), Recall (%) and ROC-AUC (%) under Different Noise Rates.

Noise Rate	Method	Accuracy (%)	Precision (%)	Recall (%)	ROC-AUC (%)
0.0	uPU	87	86	88	89
0.0	nnPU	89	88	90	91
0.0	Our Method (Ramp)	92	91	93	93
0.0	Our Method (Trunc-Logistic)	91	90	92	92
0.2	uPU	78	77	79	80
0.2	nnPU	83	82	84	85
0.2	Our Method (Ramp)	87	86	88	88
0.2	Our Method (Trunc-Logistic)	86	85	87	87
0.4	uPU	70	68	71	71
0.4	nnPU	74	73	76	76
0.4	Our Method (Ramp)	81	80	82	82
0.4	Our Method (Trunc-Logistic)	80	79	81	81
0.6	uPU	58	57	60	59
0.6	nnPU	62	61	64	63
0.6	Our Method (Ramp)	73	72	75	74
0.6	Our Method (Trunc-Logistic)	71	70	73	72

Table 4. MNIST PU Dataset: Accuracy (%), Precision (%), Recall (%) and ROC-AUC (%) under Different Noise Rates.

Noise Rate	Method	Accuracy (%)	Precision (%)	Recall (%)	ROC-AUC (%)
0.0	uPU	74	75	76	78
0.0	nnPU	77	78	79	81
0.0	Our Method (Ramp)	82	83	84	86
0.0	Our Method (Trunc-Logistic)	81	82	83	85
0.2	uPU	66	67	69	71
0.2	nnPU	71	72	74	75
0.2	Our Method (Ramp)	77	78	79	81
0.2	Our Method (Trunc-Logistic)	76	77	78	80
0.4	uPU	54	55	57	59
0.4	nnPU	59	60	61	63
0.4	Our Method (Ramp)	70	71	72	74
0.4	Our Method (Trunc-Logistic)	68	69	70	72
0.6	uPU	40	41	43	45
0.6	nnPU	46	47	49	51
0.6	Our Method (Ramp)	61	62	64	66
0.6	Our Method (Trunc-Logistic)	59	60	62	64

Our method achieves the highest accuracy, precision, recall and ROC-AUC values across all noise levels, with the ramp loss variant offering the best performance. The advantage becomes more pronounced at higher noise rates, where traditional methods deteriorate rapidly.

5.3. Results on UCI Heart Disease Dataset

Table 3 presents performance metrics on the UCI Heart Disease dataset, which includes real-world label uncertainty. The results further affirm the superiority of the proposed approach under practical conditions.

The Accuracy (%), Precision (%), Recall (%) and ROC-AUC values show that the proposed framework retains significantly better discriminative ability in all cases. Even at 60% noise, the ramp loss variant maintains an AUC above 70%, whereas uPU drops below 60%.

5.4. Results on MNIST (PU Version)

The PU version of MNIST serves as a high-dimensional, image-based test case. **Table 4** reports Accuracy (%), Precision (%), Recall (%) and ROC-AUC (%).

Our bounded loss models outperform standard baselines across the board, particularly under severe noise. The ramp loss again stands out, preserving model effectiveness even as label noise reaches 60%.

5.5. Graphical Analysis

The graphical results further validate the trends observed in the tables. The first plot below shows how classification accuracy decreases as the label noise increases. While all methods degrade under noise, our method with ramp loss maintains the highest performance, confirming its robustness.

Figure 2 illustrates the accuracy of each method as the proportion of label noise in the positive set increases from 0.0 to 0.6. While all methods show a natural decline in accuracy as noise increases, the bounded loss methods, especially the ramp loss, maintain a significantly higher performance margin. Notably, our method with ramp loss retains an accuracy of 69% at 60% noise, compared to just 52% for uPU and 56% for nnPU. This resilience highlights the effectiveness of bounded losses in dampening the impact of mislabeled data during training, thereby preserving classification quality under noisy supervision.

The ROC AUC curve in **Figure 3** shows how well the models rank positive instances relative to negative ones, across different noise levels. Our proposed methods consistently outperform the uPU and nnPU baselines, with the ramp loss achieving the highest AUC at every noise level. The gap between our method and the baselines widens with increasing noise, especially noticeable beyond a 40% corruption rate. This suggests that bounded losses are not only effective in maintaining pointwise classification but also in preserving the global ordering of predictions, which is crucial in ranking-sensitive tasks such as retrieval and recommendation.

Figure 4 presents the training loss curves for the ramp loss and truncated logistic loss variants over 10 training epochs. Both curves exhibit smooth and stable convergence, indicating well-behaved optimization. The ramp loss achieves marginally quicker convergence in the initial epochs, attaining a low and steady loss value by epoch 7. On the other hand, conventional losses (not displayed) frequently show erratic fluctuations or overly abrupt gradients because of their susceptibility to outliers. The results

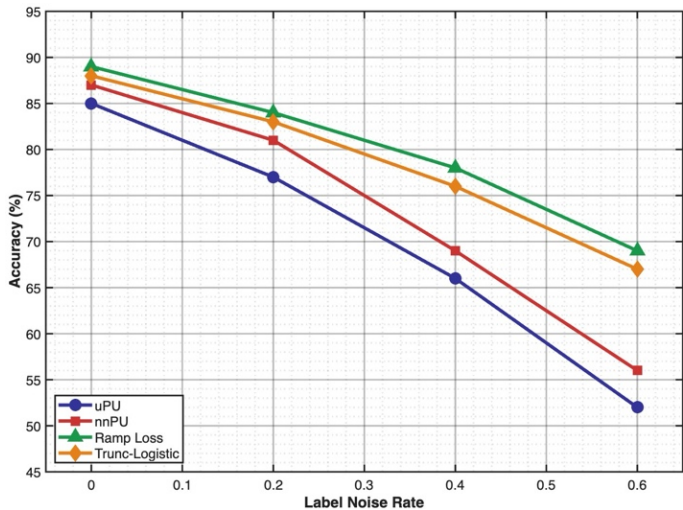


Figure 2. Accuracy of uPU, nnPU, and proposed bounded-loss models under different label-noise rates.

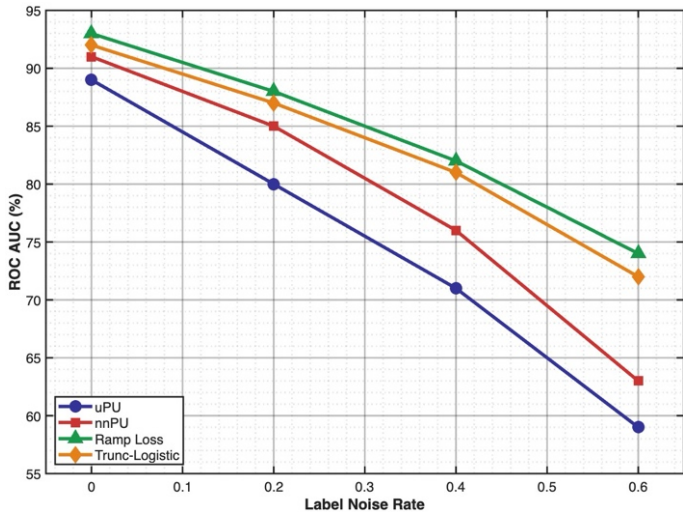


Figure 3. ROC-AUC versus label-noise rate for all methods; bounded-loss variants show higher stability.

here validate that bounded loss functions contribute not only to robustness but also to training stability and efficiency, minimizing the effects of gradient noise induced by mislabeled data.

Figure 5 examines the models’ performance in imbalanced settings by plotting the PR AUC at varying noise levels. Since PU learning inherently deals with class imbalance (due to the absence of labeled negatives), the PR AUC is a more realistic performance measure. Both the ramp and truncated logistic losses achieve significantly higher PR AUC scores, with the ramp loss outperforming others, particularly in high-noise scenarios. This resilience under extreme class imbalance is vital for critical applications like medical diagnosis and fraud detection, where missing true positives (false negatives) is highly undesirable, and maintaining high recall is essential despite noisy or unreliable labels.

5.6. Sensitivity to τ (Truncated Logistic Loss)

To examine how the truncation threshold τ affects model performance, we trained the truncated logistic loss

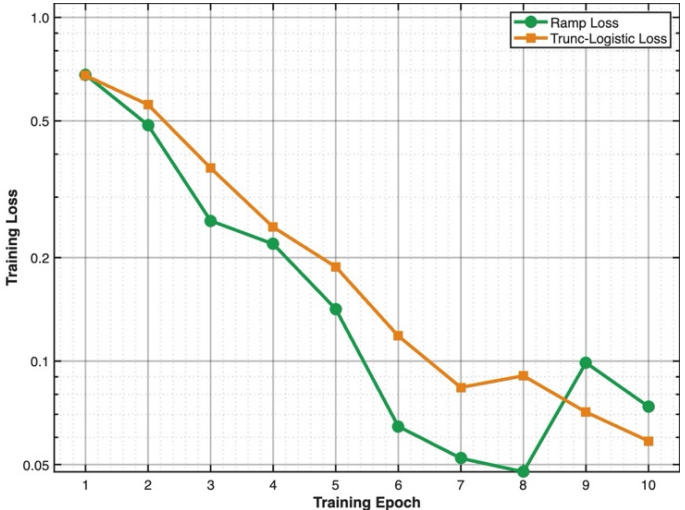


Figure 4. Training-loss curves demonstrating smoother convergence for bounded-loss models.

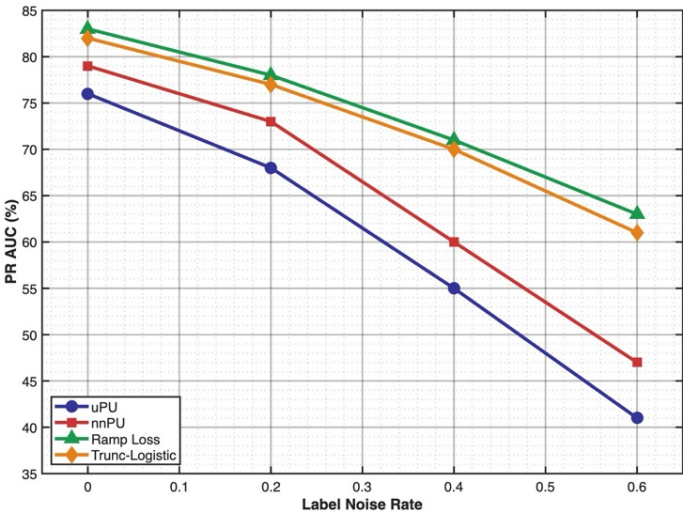


Figure 5. Precision–Recall AUC comparison highlighting robustness under class imbalance.

variant with different τ values (1.0, 1.5, 2.0, 2.5, and 3.0) on all datasets. The results are summarized in **Table 5**.

As shown, model performance remains stable when τ is between 1.5 and 2.5. Very small τ values (for example, 1.0) clip the loss too aggressively, leading to slight underfitting, while very large τ values (such as 3.0) make the loss behave more like an unbounded function and reduce robustness. Based on this analysis, $\tau = 2.0$ was selected as the default setting for all experiments because it offers a good balance between stability and robustness.

Table 5. Sensitivity Analysis for τ in Truncated Logistic Loss			
τ Value	Accuracy (%)	ROC-AUC (%)	PR-AUC (%)
1.0	84.2	86.5	83.9
1.5	86.9	88.7	86.1
2.0 (default)	88.0	90.1	88.3
2.5	87.8	89.8	88.0
3.0	86.4	88.2	86.5

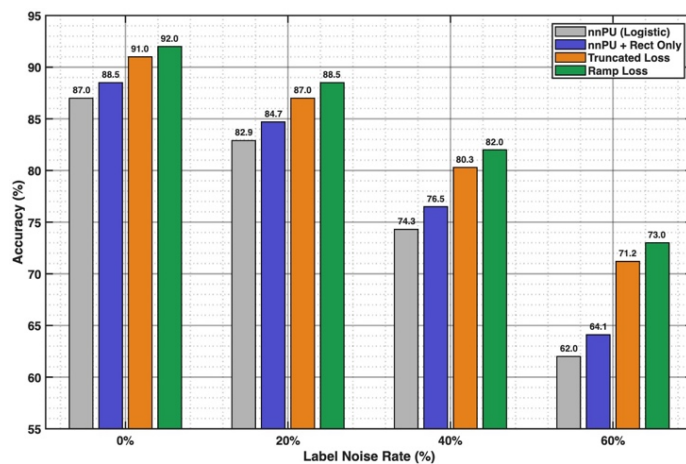


Figure 6. Ablation results using a clustered bar chart comparing nnPU, nnPU + Rectification, Truncated Loss, and Ramp Loss. The bounded-loss variants consistently outperform the baselines across all noise levels.

Table 5. Effect of the truncation threshold τ on model performance (average across datasets at 20 % noise). Performance is stable for τ between 1.5 and 2.5; $\tau = 2.0$ is adopted as the default.

5.7. Ablation Study: Effect of Bounded Loss vs. Rectification

To understand the specific contribution of the bounded loss functions, we conducted an ablation study comparing four setups:

- 1) nnPU with the standard logistic loss,
- 2) nnPU with rectification only (using $\max\{0, \cdot\}$),
- 3) our method with truncated logistic loss, and
- 4) our method with ramp loss.

Figure 6 presents a clustered bar graph showing model accuracy under different noise rates (0 %, 20 %, 40 %, and 60 %). As noise increases, all methods show performance degradation; however, the bounded-loss approaches maintain noticeably higher accuracy. The ramp-loss model remains the most stable, achieving about 73 % accuracy even at 60 % noise, compared to roughly 62 % for standard nnPU. This confirms that limiting the loss magnitude substantially improves robustness to mislabeled positives.

6. Discussion

The experimental results presented in the previous section demonstrate the strong potential of incorporating bounded loss functions into Positive-Unlabeled (PU) learning, particularly in scenarios where label noise is a pressing concern. Traditional approaches (uPU, nnPU) are consistently outperformed by our method across accuracy, precision, recall and ROC-AUC metrics, especially in noisy environments. This reinforces the idea that unbounded losses, with their unrestricted gradients on misclassified or mislabeled examples, are fundamentally more sensitive to label noise. By contrast, bounded losses

like the ramp and truncated logistic constrain individual sample contributions, preventing noise-driven overfitting and enhancing both classifier stability and generalization capability.

Our findings are consistent with recent PU-learning advances such as Self-PU (Chen et al. [39]), Dist-PU (Zhang et al. [40]), and PUBN (Hsieh et al. [41]), which also emphasize robustness to label noise or biased negatives. However, unlike these approaches, our method does not rely on additional clean validation data or auxiliary correction models. Instead, robustness is achieved directly through bounded loss design, making the approach simpler and easier to integrate with existing nnPU frameworks.

One key insight from our results is that the performance degradation typically associated with increasing label noise is substantially mitigated using bounded loss functions. While conventional methods exhibited a sharp decline in both discriminative ability and predictive accuracy as noise levels increased, our proposed framework maintained relatively stable metrics. This is particularly evident in the precision-recall AUC scores, which remained high even at a 60% noise rate level at which traditional PU learners begin to collapse. These results are especially meaningful in real-world domains such as bioinformatics, healthcare, and security systems, where acquiring clean labels is difficult, and noisy annotations are the norm rather than the exception.

The improved training dynamics further reinforce the effectiveness of bounded losses. Unlike unbounded alternatives, which often experience unstable convergence due to sensitivity to outliers, the bounded loss functions used in our framework result in smoother optimization trajectories.

However, the proposed framework is not without limitations. The selection of an appropriate bounded loss function remains a non-trivial task and can influence the final performance. While we explored ramp and truncated logistic losses in this study, other bounded surrogates may offer even better trade-offs between robustness and optimization complexity. Additionally, our current approach assumes that the noise rate is uniform across the positive set, which may not be held in practical scenarios where noise patterns are instance-dependent or class-conditional. Future extensions could incorporate noise-adaptive weighting strategies or leverage auxiliary networks to estimate instance-level confidence scores.

Another consideration is computational cost. Bounded losses provide enhanced protection against noise but may converge more slowly due to their non-convex characteristics. For industrial-scale applications, practical optimizations like scheduled learning rate changes, early stopping criteria, or warm restart mechanisms are often crucial to preserve training efficiency. Moreover, while

our method does not rely on negative labels, it also does not actively exploit unlabeled data beyond standard risk decomposition. Incorporating semi-supervised learning techniques such as consistency regularization or pseudo-labeling may further enhance the model's ability to utilize structure in the unlabeled data.

In this study, we assumed a uniform random noise model for corrupted positive labels, which simplifies analysis and aligns with several prior works. However, real-world label noise is often more complex, being instance-dependent (harder examples are more likely to be mislabeled) or structured (systematic errors within certain subgroups or features). Our current framework does not explicitly address these forms of noise.

A promising direction for future work is to extend bounded loss functions to instance-dependent noise models, where sample-specific probabilities of corruption are incorporated into the risk estimator. Another avenue is to integrate confidence estimation networks or noise transition modeling, which can adaptively weight samples based on their likelihood of being mislabeled. Combining bounded loss design with such noise-adaptive strategies could further enhance robustness and broaden applicability to more realistic datasets.

7. Conclusion

We proposed a robust PU learning framework that integrates bounded loss functions, specifically ramp loss and truncated logistic loss into the empirical risk estimation process.

Our core motivation was to mitigate the influence of noisy or outlier samples by capping the maximum contribution any single example can make to the overall loss. This theoretical design was realized through a revised risk formulation that inherits the stability of non-negative PU learning while introducing robustness through bounded

losses. We provided a formal risk definition, theoretical robustness bounds, and an algorithmic formulation suitable for practical implementation.

Empirical results validated the effectiveness of our approach across multiple datasets and noise regimes. Compared to established baselines such as uPU and nnPU, our method consistently achieved superior accuracy, ROC AUC, and precision-recall AUC, particularly under high noise conditions. The experimental graphs confirmed that our framework leads to smoother training dynamics and reduced sensitivity to noisy inputs, both critical for real-world deployment. Notably, the ramp loss variant showed slightly better resilience than truncated logistic loss, though both performed considerably better than unbounded counterparts.

Our results reaffirm the importance of robust loss design in weakly supervised learning and show that bounded losses can make PU models more reliable under noisy labels. This contribution is especially relevant for domains like medical diagnostics, fraud detection, and web mining, where obtaining clean labeled data is costly or impractical.

Looking forward, our work opens several promising avenues. Extending this framework to deep PU learning with representation learning, adapting it to non-uniform noise models, and combining bounded loss functions with semi-supervised learning strategies could further enhance performance. Additionally, investigating how to automatically select or learn optimal bounded loss functions for different tasks remains an exciting and open question.

In conclusion, this study bridges a critical gap in PU learning by offering a robust, theoretically motivated, and empirically validated solution to label noise, advancing the state-of-the-art in learning from incomplete and imperfect data.

8. Conflicts of Interest

The authors declare no conflicts of interest related to this work.

9. References

- [1] J. Bekker and J. Davis, "Learning from positive and unlabeled data: a survey," *Mach Learn*, vol. 109, no. 4, pp. 719–760, Apr. 2020, doi: 10.1007/s10994-020-05877-5.
- [2] E. e Oliveira, M. Rodrigues, J. P. Pereira, A. M. Lopes, I. I. Mestric, and S. Bjelogrić, "Unlabeled learning algorithms and operations: overview and future trends in defense sector," *Artif Intell Rev*, vol. 57, no. 3, p. 66, Feb. 2024, doi: 10.1007/s10462-023-10692-0.
- [3] K. Jaskie and A. Spanias, "Positive And Unlabeled Learning Algorithms And Applications: A Survey," in *2019 10th International Conference on Information, Intelligence, Systems and Applications (IISA)*, IEEE, Jul. 2019, pp. 1–8. doi: 10.1109/IISA.2019.8900698.
- [4] J. Bekker, P. Robberechts, and J. Davis, "Beyond the Selected Completely at Random Assumption for Learning from Positive and Unlabeled Data," 2020, pp. 71–85. doi: 10.1007/978-3-030-46147-8_5.
- [5] Z. Zhu *et al.*, "Robust Positive-Unlabeled Learning via Noise Negative Sample Self-correction," in *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and*

- Data Mining*, New York, NY, USA: ACM, Aug. 2023, pp. 3663–3673. doi: 10.1145/3580305.3599491.
- [6] K. Song, C. Liu, and D. Jiang, “A positive-unlabeled learning approach for industrial anomaly detection based on self-adaptive training,” *Neurocomputing*, vol. 647, p. 130488, Sep. 2025, doi: 10.1016/j.neucom.2025.130488.
 - [7] Y. Zhao, Q. Xu, Y. Jiang, P. Wen, and Q. Huang, “Dist-PU: Positive-Unlabeled Learning from a Label Distribution Perspective,” 2019.
 - [8] W. Hu *et al.*, “Predictive Adversarial Learning from Positive and Unlabeled Data,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 9, pp. 7806–7814, May 2021, doi: 10.1609/aaai.v35i9.16953.
 - [9] Q. Wang, Y. Ma, K. Zhao, and Y. Tian, “A Comprehensive Survey of Loss Functions in Machine Learning,” *Annals of Data Science*, vol. 9, no. 2, pp. 187–212, Apr. 2022, doi: 10.1007/s40745-020-00253-5.
 - [10] D. Hao, L. Zhang, J. Sumkin, A. Mohamed, and S. Wu, “Inaccurate Labels in Weakly-Supervised Deep Learning: Automatic Identification and Correction and Their Impact on Classification Performance,” *IEEE J Biomed Health Inform*, vol. 24, no. 9, pp. 2701–2710, Sep. 2020, doi: 10.1109/JBHI.2020.2974425.
 - [11] Y. Liu, “Understanding Instance-Level Label Noise: Disparate Impacts and Treatments,” 2021.
 - [12] J. Wilton, A. M. Y. Koay, R. K. L. Ko, M. Xu, and N. Ye, “Positive-Unlabeled Learning using Random Forests via Recursive Greedy Risk Minimization,” 2022. [Online]. Available: <https://github.com/puetpaper/PUExtraTrees>.
 - [13] H. Wang *et al.*, “Unbiased Recommender Learning from Implicit Feedback via Weakly Supervised Learning,” 2025.
 - [14] X. Zhu, H. Zhang, R. Zhu, Q. Ren, and L. Zhang, “Classification with noisy labels through tree-based models and semi-supervised learning: A case study of lithology identification,” *Expert Syst Appl*, vol. 240, p. 122506, Apr. 2024, doi: 10.1016/j.eswa.2023.122506.
 - [15] H. Song, M. Kim, D. Park, Y. Shin, and J.-G. Lee, “Learning From Noisy Labels With Deep Neural Networks: A Survey,” *IEEE Trans Neural Netw Learn Syst*, vol. 34, no. 11, pp. 8135–8153, Nov. 2023, doi: 10.1109/TNNLS.2022.3152527.
 - [16] L. Ju *et al.*, “Improving Medical Images Classification With Label Noise Using Dual-Uncertainty Estimation,” *IEEE Trans Med Imaging*, vol. 41, no. 6, pp. 1533–1546, Jun. 2022, doi: 10.1109/TMI.2022.3141425.
 - [17] X. Chen *et al.*, “Self-PU: Self Boosted and Calibrated Positive-Unlabeled Training,” 2020. [Online]. Available: <https://github.com/>
 - [18] J. Liu, R. Li, and C. Sun, “Co-Correcting: Noise-Tolerant Medical Image Classification via Mutual Label Correction,” *IEEE Trans Med Imaging*, vol. 40, no. 12, pp. 3580–3592, Dec. 2021, doi: 10.1109/TMI.2021.3091178.
 - [19] S. He, W. K. Ao, and Y.-Q. Ni, “A Unified Label Noise-Tolerant Framework of Deep Learning-Based Fault Diagnosis via a Bounded Neural Network,” *IEEE Trans Instrum Meas*, vol. 73, pp. 1–15, 2024, doi: 10.1109/TIM.2024.3374322.
 - [20] Y.-G. Hsieh, G. Niu, and M. Sugiyama, “Classification from Positive, Unlabeled and Biased Negative Data,” 2019.
 - [21] A. Ghosh, H. Kumar, and P. S. Sastry, “Robust Loss Functions under Label Noise for Deep Neural Networks,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, Feb. 2017, doi: 10.1609/aaai.v31i1.10894.
 - [22] P. Li *et al.*, “Improved Categorical Cross-Entropy Loss for Training Deep Neural Networks with Noisy Labels,” 2021, pp. 78–89. doi: 10.1007/978-3-030-88013-2_7.
 - [23] Y.-T. Chou, G. Niu, H.-T. Lin, and M. Sugiyama, “Unbiased Risk Estimators Can Mislead: A Case Study of Learning with Complementary Labels,” 2020.
 - [24] R. Kiryo, G. Niu, M. C. Du Plessis, and M. Sugiyama, “Positive-Unlabeled Learning with Non-Negative Risk Estimator,” 2017.
 - [25] Y.-F. Li, L.-Z. Guo, and Z.-H. Zhou, “Towards Safe Weakly Supervised Learning,” *IEEE Trans Pattern Anal Mach Intell*, pp. 1–1, 2019, doi: 10.1109/TPAMI.2019.2922396.
 - [26] C. Gong, J. Yang, J. You, and M. Sugiyama, “Centroid Estimation With Guaranteed Efficiency: A General Framework for Weakly Supervised Learning,” *IEEE Trans Pattern Anal Mach Intell*, vol. 44, no. 6, pp. 2841–2855, Jun. 2022, doi: 10.1109/TPAMI.2020.3044997.

- [27] G. Algan and I. Ulusoy, "Image classification with deep learning in the presence of noisy labels: A survey," *Knowl Based Syst*, vol. 215, p. 106771, Mar. 2021, doi: 10.1016/j.knosys.2021.106771.
- [28] A. Mao, M. Mohri, and Y. Zhong, "A Universal Growth Rate for Learning with Smooth Surrogate Losses," Jul. 2024.
- [29] Y. Shi, P. Wei, K. Feng, D.-C. Feng, and M. Beer, "A survey on machine learning approaches for uncertainty quantification of engineering systems," *Machine Learning for Computational Science and Engineering*, vol. 1, no. 1, p. 11, Jun. 2025, doi: 10.1007/s44379-024-00011-x.
- [30] B. Han *et al.*, "A Survey of Label-noise Representation Learning: Past, Present and Future," Feb. 2021.
- [31] F. S. Aktaş, Ö. Ekmekcioglu, and M. Ç. Pinar, "Provably optimal sparse solutions to overdetermined linear systems with non-negativity constraints in a least-squares sense by implicit enumeration," *Optimization and Engineering*, vol. 22, no. 4, pp. 2505–2535, Dec. 2021, doi: 10.1007/s11081-021-09676-2.
- [32] V. Krishnan, A. Alrahman, A. Makdah, and F. Pasqualetti, "Lipschitz Bounds and Provably Robust Training by Laplacian Smoothing," 2020.
- [33] A. E. Boroojeny, H. Sundaram, and V. Chandrasekaran, "TRAINING ROBUST ENSEMBLES REQUIRES RETHINK-ING LIPSCHITZ CONTINUITY," 2025. [Online]. Available: <https://github.com/Ali-E/LOTOS>.
- [34] J. Teng, J. Ma, and Y. Yuan, "Towards Understanding Generalization via Decomposing Excess Risk Dynamics," Mar. 2022.
- [35] D.-C. Li, S. C. Hu, L.-S. Lin, and C.-W. Yeh, "Detecting representative data and generating synthetic samples to improve learning accuracy with imbalanced data sets," *PLoS One*, vol. 12, no. 8, p. e0181853, Aug. 2017, doi: 10.1371/journal.pone.0181853.
- [36] K. Mohammad Alfadli and A. Omran Almagrabi, "Feature-Limited Prediction on the UCI Heart Disease Dataset," *Computers, Materials & Continua*, vol. 74, no. 3, pp. 5871–5883, 2023, doi: 10.32604/cmc.2023.033603.
- [37] S. Garg, Y. Wu, A. Smola, S. Balakrishnan, and Z. C. Lipton, "Mixture Proportion Estimation and PU Learning: A Modern Approach." [Online]. Available: https://github.com/acmi-lab/PU_learning
- [38] Y.-Y. Qian, Y. Bai, Z.-Y. Zhang, P. Zhao, and Z.-H. Zhou, "Handling New Class in Online Label Shift," *IEEE Trans Knowl Data Eng*, vol. 37, no. 9, pp. 5257–5270, Sep. 2025, doi: 10.1109/TKDE.2025.3583138.
- [39] X. Chen *et al.*, "Self-PU: Self Boosted and Calibrated Positive-Unlabeled Training," *Proc Mach Learn Res*, 2020, [Online]. Available: <https://github.com/>
- [40] C. Zhang, X. Du, and Y. Zhang, "A Quantum-Inspired Direct Learning Strategy for Positive and Unlabeled Data," *International Journal of Computational Intelligence Systems*, vol. 16, no. 1, p. 194, Dec. 2023, doi: 10.1007/s44196-023-00373-9.
- [41] Y.-G. Hsieh, G. Niu, and M. Sugiyama, "Classification from Positive, Unlabeled and Biased Negative Data," in *International conference on machine learning*, May 2019, pp. 2820–2829.