

# Accurate and Robust Positive-Unlabeled Learning against Adversarial Perturbations

Ryo Shibazaki<sup>1\*</sup>, Kazuhiko Kawamoto<sup>2</sup> and Hiroshi Kera<sup>2</sup>

<sup>1\*</sup>Graduate School of Science and Engineering, Chiba University, 1-33  
Yayoichō, Inage-ku, Chiba-shi, 263-8522, Chiba, Japan.

<sup>2</sup>Graduate School of Informatics, Chiba University, 1-33 Yayoichō,  
Inage-ku, Chiba-shi, 263-8522, Chiba, Japan.

\*Corresponding author(s). E-mail(s): [ryo.shibazaki0517@chiba-u.jp](mailto:ryo.shibazaki0517@chiba-u.jp);  
Contributing authors: [kawa@faculty.chiba-u.jp](mailto:kawa@faculty.chiba-u.jp); [kera@chiba-u.jp](mailto:kera@chiba-u.jp);

## Abstract

Labeling costs are high in domains such as medical image analysis, where Positive and Unlabeled (PU) learning, which trains using only positive and unlabeled data, is effective. Medical images often contain small perturbations due to sensor noise and variations in acquisition conditions, which can cause a classifier to misclassify images that should be positive as negative. Therefore, in settings where even minor misclassifications may lead to critical misdiagnoses, high robustness is required. In this study, we focus on adversarial perturbations, which are known as worst-case noise among such perturbations, and aim to improve robustness within the PU learning framework. However, directly applying standard adversarial training methods to PU learning often severely degrades standard accuracy, making the trade-off between robustness and standard accuracy more pronounced. To address this issue, we propose PU+TRADES, a new learning method that extends the TRADES framework and integrates it with PU learning. Our method introduces label-independent adversarial perturbations and optimizes the balance between robustness and standard accuracy by combining a PU loss with a Kullback–Leibler loss. Furthermore, we theoretically derive an upper bound on the estimation error for the proposed loss and clarify conditions under which PU learning can outperform supervised learning when the number of unlabeled samples is sufficiently large. Finally, experiments on multiple benchmark datasets and a medical imaging dataset demonstrate that the proposed method provides an effective framework for robust learning in PU settings.

**Keywords:** positive-unlabeled learning, adversarial robustness, risk estimation, empirical risk minimization

# 1 Introduction

In recent years, machine learning has achieved remarkable success across a wide range of tasks, driven by advances in large-scale data and high-capacity models. However, in real-world applications, it is often difficult to obtain high-quality ...; moreover, in settings where the training set contains only positive and unlabeled data, one must learn from incomplete information in which the unlabeled set is a mixture of true positives and negatives. Therefore, unlike fully supervised learning, PU learning requires careful risk estimation and additional techniques to stabilize training.

Furthermore, in practical applications, robustness is an essential requirement in addition to label scarcity. In particular, medical images and sensor data are affected by variations in acquisition conditions, noise, and device ... We focus on adversarial perturbations[1] and aim to improve robustness within the PU learning framework.

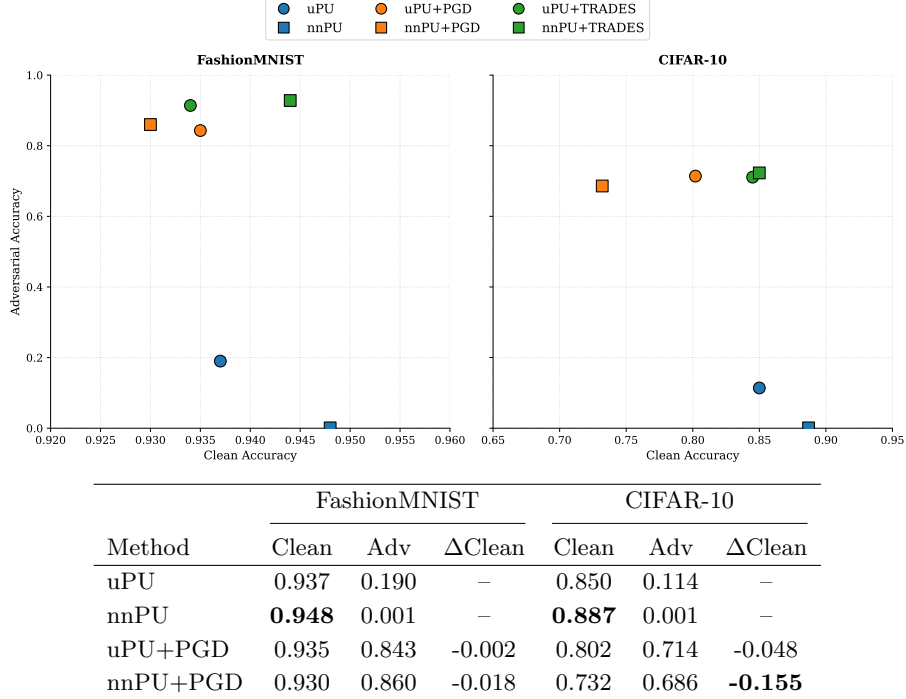
However, in preliminary experiments, directly applying adversarial training[2]—a standard approach to enhancing resilience against adversarial perturbations—to PU learning can substantially degrade classification performance on clean data. This is because, although unlabeled data contain a mixture of true positives and negatives, they are treated uniformly as negatives in the loss. As a result, the objectives of adversarial perturbation optimization and classification become misaligned, which destabilizes the updates associated with the unlabeled term. Figure 1 illustrates the relationship between clean accuracy and adversarial accuracy when PGD-based adversarial training is naively applied to PU learning. In particular, on CIFAR-10, adversarial accuracy improves while clean accuracy drops significantly. These observations indicate that, under the PU setting, introducing adversarial training may impair clean accuracy, which remains a key challenge.

In this study, we extend the TRADES framework[3], a representative adversarial training method, and propose a new learning method, PU+TRADES, by integrating it with PU learning. Our approach handles adversarial perturbations in a label-independent manner and combines a PU loss with a Kullback–Leibler loss, aiming to improve robustness while suppressing degradation in clean accuracy.

To validate the effectiveness of the proposed method, we conducted experiments on multiple benchmark datasets and medical imaging data. We evaluated both clean accuracy and ... and compared them with baseline methods. The results show that the proposed method substantially improves robustness to adversarial perturbations while maintaining accuracy on clean data.

In addition, we performed a theoretical analysis to better understand robustness in PU learning. Specifically, under binary classification with linear classifiers and adversarial perturbations, we derived upper bounds on the estimation error of the risks minimized by supervised learning and PU learning, thereby clarifying conditions under which PU learning can be advantageous over supervised learning. These conditions are consistent with practical scenarios, and they suggest that simply increasing the amount of unlabeled data can potentially achieve higher robustness than supervised learning.

Our contributions are summarized as follows.



**Fig. 1** Performance changes when PGD-based adversarial training is naively applied to PU learning (FashionMNIST / CIFAR-10). Left: a scatter plot showing the relationship between clean accuracy and adversarial accuracy for each method. Right: a numerical summary and the change in clean accuracy,  $\Delta$ Clean, before and after applying PGD (difference from the corresponding PU method without PGD). While PGD substantially improves adversarial accuracy, it degrades clean accuracy; this degradation is particularly pronounced for nnPU on CIFAR-10, where clean accuracy drops from 0.887 to 0.732 ( $\Delta$ Clean =  $-0.155$ ).

- We propose a new learning framework (PU+TRADES) that integrates TRADES-style regularization into PU learning, improving robustness to adversarial perturbations while maintaining classification accuracy on clean samples.
- Assuming linear classifiers under adversarial perturbations, we derive upper bounds on the estimation error of the risks minimized by supervised learning and PU learning, and theoretically identify conditions under which PU learning becomes more favorable than supervised learning.
- Through experiments on benchmark datasets and medical imaging data, we empirically demonstrate that the proposed method acquires robustness to adversarial samples while preserving clean accuracy.

## 2 Related Work

In this chapter, we organize representative studies on PU learning and adversarial learning to clarify the positioning of this research.

## 2.1 PU Learning

Positive-Unlabeled (PU) learning has been widely studied as a framework for building classifiers in settings where only positive and unlabeled data are available. It is particularly effective in application domains such as medical image analysis, where explicit labeling of negative examples is difficult, and many recent studies have proposed methods that achieve performance comparable to supervised learning.

In risk-estimation-based approaches, starting from uPU (unbiased PU) [4, 5], loss functions have been designed to provide unbiased estimates of the classification risk [6–12]. nnPU [7] suppresses overfitting by imposing a non-negativity constraint on this loss, enabling more stable learning. In addition, methods such as Imbalanced PU learning (ImbPU) [12] and Self-PU [6], which address class imbalance and self-training-style learning, have also been proposed. On the other hand, sample-selection-based methods extract highly reliable negative or positive examples from unlabeled data and treat them as labeled samples, thereby reducing the problem to supervised learning [13–19].

Although these studies have primarily focused on improving standard classification accuracy, they have paid little attention to robustness against adversarial perturbations or noise. Therefore, in this study, we build on risk-correction-based methods such as uPU and nnPU while introducing a robustness perspective to address a new challenge in PU learning.

## 2.2 Adversarial Training

It is well known that deep learning models are vulnerable to “adversarial examples,” in which small perturbations added to the input induce misclassification. The Fast Gradient Sign Method (FGSM) proposed by Goodfellow et al. [1] is an efficient method for generating adversarial examples by modifying the input in a single step along the gradient direction of the loss function. In contrast, Projected Gradient Descent (PGD) proposed by Madry et al. [2] iteratively applies FGSM and projects the perturbation back into the allowable region at each step, enabling stronger attacks.

As a major defense against these attacks, adversarial training, which incorporates adversarial examples into the training process, has been widely studied. Adversarial training is generally formulated as a min–max optimization problem consisting of outer parameter updates (minimization) and inner perturbation generation (maximization). In addition, TRADES [3] introduces a regularization term combining classification loss and Kullback-Leibler (KL) loss to control the trade-off between clean-data accuracy and adversarial robustness, and it has become an important foundation of modern robust learning.

However, these standard methods assume supervised learning, where all data are labeled. When they are directly applied to PU learning, which uses only positive and unlabeled data, the design of the loss function and the optimization process for adversarial perturbations become inconsistent, and preliminary experiments confirm that classification performance deteriorates significantly. In this study, we extend the TRADES framework to the PU learning setting and propose a method that achieves high robustness under the constraints specific to PU learning.

### 3 Preliminaries

In this section, we introduce PU learning, adversarial examples and representative attacks, and adversarial training. Hereafter, we refer to learning a binary classifier from fully labeled positive and negative data as Positive–Negative (PN) learning.

#### 3.1 Positive-Unlabeled (PU) Learning

We denote the input space by  $\mathcal{X} \subseteq \mathbb{R}^d$  and the label space by  $\mathcal{Y} = \{-1, +1\}$ . Let  $p(\mathbf{x}, y)$  be the joint distribution over  $(\mathcal{X}, \mathcal{Y})$ . Let the total number of samples be  $n \in \mathbb{N}$ , and let  $n_P$  and  $n_N$  denote the numbers of positive (P) and negative (N) samples, respectively. Each set is represented as follows:

$$\begin{aligned}\mathcal{X}_P &= \{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p_P(\mathbf{x}), \\ \mathcal{X}_N &= \{\mathbf{x}_i^N\}_{i=1}^{n_N} \stackrel{\text{i.i.d.}}{\sim} p_N(\mathbf{x}).\end{aligned}\tag{1}$$

Here,  $p_P(\mathbf{x})$  and  $p_N(\mathbf{x})$  denote the class-conditional densities for the positive and negative classes, respectively. The full dataset  $\mathcal{X} = \mathcal{X}_P \cup \mathcal{X}_N$  is written as

$$\begin{aligned}\mathcal{X} &= \{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}), \\ p(\mathbf{x}) &= \pi_P p_P(\mathbf{x}) + \pi_N p_N(\mathbf{x}),\end{aligned}\tag{2}$$

where  $\pi_P = p(y = +1)$  and  $\pi_N = \dots p(y = -1)$  are the class priors satisfying  $\pi_P + \pi_N = 1$ .

In PU learning, the training set consists of positive (P) samples and unlabeled (U) samples. Since the marginal distribution of unlabeled data is  $p_U(\mathbf{x}) = \pi_P p_P(\mathbf{x}) + \pi_N p_N(\mathbf{x})$ , letting  $n_U$  be the number of unlabeled samples, the unlabeled set is given by

$$\mathcal{X}_U = \{\mathbf{x}_i^U\}_{i=1}^{n_U} \stackrel{\text{i.i.d.}}{\sim} p_U(\mathbf{x}) = p(\mathbf{x}).\tag{3}$$

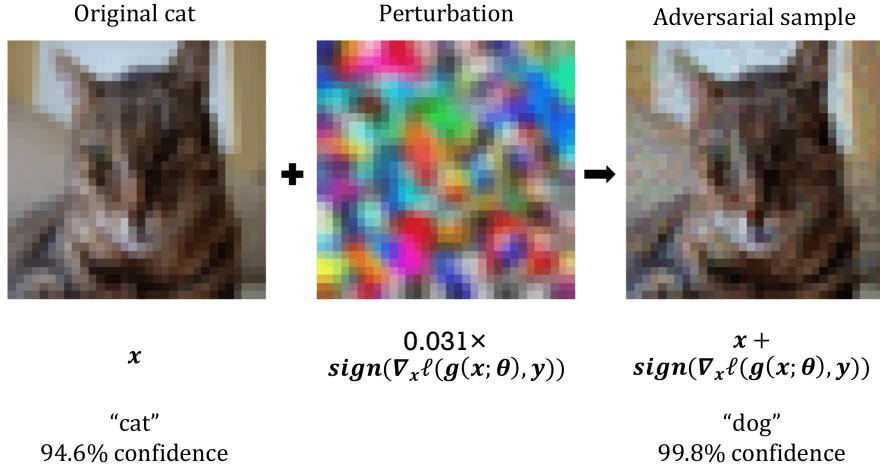
That is, the unlabeled samples are drawn i.i.d. from the marginal distribution of inputs, which is a mixture of positive and negative class-conditional distributions.

**Unbiased PU Learning (uPU).** uPU assumes that the positive class prior  $\pi_P$  is known and estimates the negative risk indirectly from the unlabeled data. Specifically, it minimizes the following empirical risk:

$$\hat{R}_{\text{uPU}}(g) = \frac{\pi_P}{n_P} \sum_{i=1}^{n_P} \tilde{\ell}(g(\mathbf{x}_i^P), +1) + \frac{1}{n_U} \sum_{i=1}^{n_U} \ell(g(\mathbf{x}_i^U), -1).\tag{4}$$

Here, the composite loss  $\tilde{\ell}(g(\mathbf{x}), y)$  is defined by  $\tilde{\ell}(g(\mathbf{x}), y) = \ell(g(\mathbf{x}), y) - \ell(g(\mathbf{x}), -y)$ .

**Non-Negative PU Learning (nnPU).** nnPU was introduced to address the overfitting issue in uPU, where the empirical risk can take negative values[7]. Specifically, when the estimated term involving the negative risk in PU learning becomes negative, nnPU clips its negative contribution to zero, yielding a non-negative risk estimator



**Fig. 2** An example of generating adversarial examples. Starting from the clean image on the left, we add a small perturbation using PGD to obtain an adversarial input... The classifier correctly predicts the clean image as a cat (confidence 94.6%), while it misclassifies the adversarial example as a dog (confidence 99.8%). This illustrates that predictions can change drastically even when the input appears almost identical to humans[1].

that keeps the overall estimate bounded below by 0. This modification prevents the empirical risk from diverging to negative values during empirical minimization and enables stable training. The empirical risk of nnPU is defined as follows:

$$\begin{aligned}
 \hat{R}_{\text{nnPU}}(g) = & \frac{\pi_P}{n_P} \sum_{i=1}^{n_P} \ell(g(x_i^P), +1) \\
 & + \max \left\{ 0, -\frac{\pi_P}{n_P} \sum_{i=1}^{n_P} \ell(g(x_i^P), -1) + \frac{1}{n_U} \sum_{i=1}^{n_U} \ell(g(x_i^U), -1) \right\}.
 \end{aligned} \tag{5}$$

### 3.2 Adversarial Examples

Adversarial examples are inputs constructed by adding small, carefully designed perturbations to an image that a classifier originally classified correctly, thereby intentionally causing misclassification...w2014ExplainingAH. For example, Fig. 2 adds a tiny perturbation to a cat image and generates an adversarial example...

### 3.3 Adversarial attack: FGSM and PGD

We refer to methods for generating adversarial examples as *adversarial attacks*, and many variants have been studied. In this work, we focus on representative first-order attacks: the Fast Gradient Sign Method (FGSM) [1] and Projected Gradient Descent

(PGD) [2]. Here,  $\text{sign} : \mathbb{R} \rightarrow [-1, 1]$  is applied element-wise to the argument vector.

$$\mathbf{x}' = \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \ell(g(\mathbf{x}; \boldsymbol{\theta}), y)) \quad (6)$$

This produces an input that increases the loss when fed into the model. A stronger iterative variant of this method is PGD [2], which updates the input in the direction that increases the loss with step size  $\alpha$ , similarly to FGSM, and then applies the projection  $\Pi_{\dots}$ . In particular,

$$\mathcal{B}_{\infty}(\mathbf{x}, \epsilon) := \{\mathbf{z} \in \mathbb{R}^d \mid \|\mathbf{z} - \mathbf{x}\|_{\infty} \leq \epsilon\}$$

Then,  $\Pi_{\mathcal{B}_{\infty}(\mathbf{x}, \epsilon)}$  denotes the projection onto the  $\ell_{\infty}$ -ball, which guarantees  $\|\mathbf{x}' - \mathbf{x}\|_{\infty} \leq \epsilon$ . Under this setting, the PGD update is given by:

$$\mathbf{x}' \leftarrow \Pi_{\mathcal{B}_{\infty}(\mathbf{x}, \epsilon)} [\mathbf{x} + \alpha \text{sign}(\nabla_{\mathbf{x}} \ell(g(\mathbf{x}; \boldsymbol{\theta}), y))]. \quad (7)$$

By repeating Eq. (7) multiple times, we can generate samples that more strongly increase the loss within the  $\epsilon$ -ball.

### 3.4 Adversarial Training

Adversarial training improves model robustness by training on adversarial examples [2]. It is typically formulated as the following min-max optimization problem: one first generates adversarial examples for each input, and then learns parameters that minimize the average loss over these adversarial inputs.

$$\min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \max_{\|\mathbf{x}'_i - \mathbf{x}_i\|_{\infty} \leq \epsilon} \ell(g(\mathbf{x}'_i; \boldsymbol{\theta}), y_i). \quad (8)$$

In addition, as a method to further enhance robustness against adversarial examples, TRADES (TRadeoff-inspired Adversarial DEfense via Surrogate-loss minimization), proposed by Zhang *et al.* [3], is a prominent approach. TRADES explicitly models the trade-off between accuracy on clean samples and robustness to adversarial samples, and aims to minimize the following loss function:

$$\begin{aligned} \mathcal{L}_{\text{TR}}(\boldsymbol{\theta}) = & \frac{1}{n} \sum_{i=1}^n \left[ \ell(g(\mathbf{x}_i; \boldsymbol{\theta}), y_i) \right. \\ & \left. + \beta \cdot \max_{\|\mathbf{x}'_i - \mathbf{x}_i\|_{\infty} \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x}_i; \boldsymbol{\theta}), g(\mathbf{x}'_i; \boldsymbol{\theta})) \right]. \end{aligned} \quad (9)$$

Here,  $\ell_{\text{KL}}(\cdot, \cdot)$  denotes the Kullback-Leibler (KL) divergence between predictive distributions, i.e.,

$$\ell_{\text{KL}}(g(\mathbf{x}_i; \boldsymbol{\theta}), g(\mathbf{x}'_i; \boldsymbol{\theta})) := \text{KL}(p_{\boldsymbol{\theta}}(\cdot \mid \mathbf{x}_i) \parallel p_{\boldsymbol{\theta}}(\cdot \mid \mathbf{x}'_i)).$$

The first term of Eq. (9),  $\ell(g(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$ , is the standard classification loss on the clean input  $\mathbf{x}_i$ .

On the other hand, the second term  $\ell_{\text{KL}}(g(\mathbf{x}_i; \boldsymbol{\theta}), g(\mathbf{x}'_i; \boldsymbol{\theta}))$  constrains the model so that the output distributions for  $\mathbf{x}_i$  and its perturbed version  $\mathbf{x}'_i$  are close, and this term plays a key role in improving robustness. Thus, TRADES is designed to enhance robustness while maintaining classification accuracy. In this study, we apply this framework to PU learning to achieve both high performance on clean samples and robustness to adversarial examples.

## 4 Accurate and Robust PU Learning

In this section, we first clarify the issues that arise when uPU learning is naively combined with PGD-based adversarial training. We then propose a new learning method, PU+TRADES, which adapts the TRADES framework to PU learning.

### 4.1 uPU+PGD

In uPU learning, the empirical risk  $\hat{R}_{\text{uPU}}(g)$  is estimated by minimizing

$$\hat{R}_{\text{uPU}}(g) = \frac{\pi_{\text{P}}}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \tilde{\ell}(g(\mathbf{x}_i^{\text{P}}), +1) + \frac{1}{n_{\text{U}}} \sum_{i=1}^{n_{\text{U}}} \ell(g(\mathbf{x}_i^{\text{U}}), -1). \quad (10)$$

Here, the loss is taken differently for P and U samples, which can be summarized as

$$\mathcal{L}(\mathbf{x}) := \begin{cases} \tilde{\ell}(g(\mathbf{x}), +1), & \mathbf{x} \in \mathcal{X}_{\text{P}}, \\ \ell(g(\mathbf{x}), -1), & \mathbf{x} \in \mathcal{X}_{\text{U}}. \end{cases} \quad (11)$$

Using this loss  $\mathcal{L}$ , we generate an adversarial example for each sample via PGD. A single PGD update step is given by

$$\mathbf{x}' \leftarrow \text{Clip}_{(\mathbf{x}-\epsilon, \mathbf{x}+\epsilon)} \left[ \mathbf{x}' + \alpha \text{sign}(\nabla_{\mathbf{x}'} \mathcal{L}(\mathbf{x}')) \right]. \quad (12)$$

### Issues with uPU+PGD

In uPU, the loss for unlabeled data is computed as if the label were always  $y = -1$ . However, in reality, the unlabeled set contains a mixture of positives and negatives. This property is incompatible with PGD-based adversarial training.

- **Negative U samples.** Since the loss  $\ell(g(\mathbf{x}^{\text{U}}), -1)$  is consistent with the true label, it pushes the input in a direction that increases the loss for the negative class. Consequently, PGD generates appropriate adversarial perturbations, contributing to improved robustness.
- **Positive U samples.** If perturbations are generated using  $\ell(g(\mathbf{x}^{\text{U}}), -1)$  even though the sample is truly positive, PGD updates the input so as to maximize the *negative-class* loss. As a result, the input may be pushed not toward the decision boundary, but rather toward a region where it is classified as positive with



higher confidence. Therefore, PGD fails to produce perturbations in the “most misclassifiable direction,” and the training can break down.

Hence, to generate adversarial perturbations appropriately in PU learning, it is essential to use a *label-independent* perturbation generation mechanism. This motivates PU+TRADES, introduced in the next section.

## 4.2 PU+TRADES

In this work, we propose **uPU+TRADES** and **nnPU+TRADES**, which integrate TRADES into uPU and nnPU, respectively. By introducing the TRADES framework into PU learning, we endow the model with robustness.

The objective function of uPU+TRADES is given by

$$\min_g \left[ \widehat{R}_{\text{uPU}}(g) + \beta \cdot \frac{1}{n} \sum_{i=1}^n \max_{\|\mathbf{x}' - \mathbf{x}_i\|_\infty \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x}_i) \parallel g(\mathbf{x}')) \right], \quad (13)$$

and the objective function of nnPU+TRADES is given by

$$\min_g \left[ \widehat{R}_{\text{nnPU}}(g) + \beta \cdot \frac{1}{n} \sum_{i=1}^n \max_{\|\mathbf{x}' - \mathbf{x}_i\|_\infty \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x}_i) \parallel g(\mathbf{x}')) \right]. \quad (14)$$

In binary classification, the network outputs a one-dimensional logit  $g(\mathbf{x}; \boldsymbol{\theta})$ . We convert it into a Bernoulli probability vector and compute the KL divergence:

$$p(\mathbf{x}) = [\sigma(g(\mathbf{x}; \boldsymbol{\theta})), 1 - \sigma(g(\mathbf{x}; \boldsymbol{\theta}))], \quad (15)$$

where  $\sigma(\cdot)$  denotes the sigmoid function. The KL loss is then defined as

$$\begin{aligned} \ell_{\text{KL}}(g(\mathbf{x}_i; \boldsymbol{\theta}), g(\mathbf{x}'_i; \boldsymbol{\theta})) &= \text{KL}(p(\mathbf{x}_i) \parallel p(\mathbf{x}'_i)) \\ &= \sum_{c \in \{0,1\}} p_c(\mathbf{x}_i) \log \frac{p_c(\mathbf{x}_i)}{p_c(\mathbf{x}'_i)}. \end{aligned} \quad (16)$$

This term encourages the model outputs to remain stable under small perturbations of  $\mathbf{x}_i$ , thereby providing robustness against adversarial perturbations.

## 5 Theoretical Analysis

In this section, we consider a binary classification problem where input data may be subject to adversarial perturbations, and derive upper bounds on the gap between the true adversarial risk and its empirical estimator for both supervised learning and PU learning. By comparing these bounds, we theoretically clarify conditions under which PU learning can be advantageous over supervised learning in the finite-sample regime.

Below we define the problem setting used throughout this section.

**Problem Setting 5.1.** (Adversarial Binary Classification Setting)

Let the input space be  $\mathcal{X} \subseteq \mathbb{R}^d$ , and the label space be  $\mathcal{Y} = \{-1, +1\}$ . Assume the input  $\mathbf{x} \in \mathcal{X}$  is bounded, i.e., there exists a constant  $C_x > 0$  such that

$$\|\mathbf{x}\|_\infty \leq C_x$$

holds.

Assume  $(\mathbf{x}, y)$  is generated from a joint distribution  $p(\mathbf{x}, y)$ , and define the class-conditional distributions as

$$p_P(\mathbf{x}) = p(\mathbf{x} \mid y = +1), \quad p_N(\mathbf{x}) = p(\mathbf{x} \mid y = -1)$$

as above. Let the class prior probabilities be  $\pi_P = p(y = +1)$ ,  $\pi_N = p(y = -1)$  and let with  $\pi_P + \pi_N = 1$ .

In supervised learning, we use labeled data from the positive (P) and negative (N) classes. The corresponding datasets are

$$\begin{aligned} \mathcal{X}_P &= \{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p_P(\mathbf{x}), \\ \mathcal{X}_N &= \{\mathbf{x}_i^N\}_{i=1}^{n_N} \stackrel{\text{i.i.d.}}{\sim} p_N(\mathbf{x}) \end{aligned} \tag{17}$$

given by

In PU learning, we use positive data and unlabeled (U) data. The marginal distribution of unlabeled data is

$$p_U(\mathbf{x}) = \pi_P p_P(\mathbf{x}) + \pi_N p_N(\mathbf{x})$$

given by and the unlabeled dataset is

$$\mathcal{X}_U = \{\mathbf{x}_i^U\}_{i=1}^{n_U} \stackrel{\text{i.i.d.}}{\sim} p_U(\mathbf{x}) \tag{18}$$

Let the classifier be  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ , and consider a linear classifier parameterized by a weight vector  $\mathbf{w} \in \mathbb{R}^d$ :

$$g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$$

For  $p \geq 1$  and  $W > 0$ , define the hypothesis class as

$$\mathcal{G} = \{g(\mathbf{x}) : \|\mathbf{w}\|_p \leq W\}$$

as above.

Let the loss function  $\ell : \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}$  take as arguments the classifier output and the true label.

Moreover, as the adversarial regularization term in TRADES we use

$$\ell_{\text{KL}}(g(\mathbf{x}), g(\mathbf{x}'))$$

where  $\ell_{\text{KL}}$  denotes the Kullback–Leibler divergence between probability distributions induced by the classifier outputs.

Throughout this section, we assume that  $\ell$  and  $\ell_{\text{KL}}$  satisfy the regularity conditions needed for the analysis (boundedness and Lipschitz continuity); see the assumptions in Theorem 5.2 for details.

## 5.1 Preliminaries (Notation and Assumptions)

To derive the estimation-error upper bounds in this section, we use the Rademacher complexity to bound the expected uniform deviation, and we use McDiarmid’s inequality to obtain high-probability guarantees that hold with probability at least  $1 - \delta$ . Below we summarize the definitions and properties used in this section.

### Rademacher Complexity

In our discussion, we need a measure of the complexity of a function class  $\mathcal{G}$ . We adopt the Rademacher complexity and recall its definition. A random variable  $\sigma$  satisfying  $\Pr(\sigma = +1) = \Pr(\sigma = -1) = 1/2$  is called a Rademacher variable. Given a set of  $n$ -dimensional vectors  $S \subseteq \mathbb{R}^n$ , the Rademacher complexity of  $S$  is defined as

$$\mathfrak{R}(S) = \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[ \sup_{(s_1, \dots, s_n) \in S} \frac{1}{n} \sum_{i=1}^n \sigma_i s_i \right]$$

as above. It can be interpreted as the expected correlation between random noise and the best-matching element of  $S$ .

Next, let  $S_n = \{\mathbf{x}_i\}_{i=1}^n$  be a dataset of size  $n$ . For a function class  $\mathcal{G}$ , define

$$\mathcal{G} \circ S_n = \{(g(\mathbf{x}_1), \dots, g(\mathbf{x}_n)) \mid g \in \mathcal{G}\}$$

Then, the empirical Rademacher complexity of  $\mathcal{G}$  on  $S_n$  is

$$\mathfrak{R}_{S_n}(\mathcal{G}) = \mathfrak{R}(\mathcal{G} \circ S_n) = \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{x}_i) \right]$$

given by Furthermore, when  $S_n = \{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \nu(\mathbf{x})$ , we define the Rademacher complexity of  $\mathcal{G}$  as follows (where  $\nu(\mathbf{x})$  denotes the distribution of  $\mathbf{x}$ ).

**Definition 5.1** (Rademacher Complexity). Let  $n$  be the sample size and let  $S_n = \{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \nu(\mathbf{x})$ . Then the Rademacher complexity of  $\mathcal{G}$  is

$$\mathfrak{R}_{n, \nu}(\mathcal{G}) = \mathbb{E}_{S_n \sim \nu^n} [\mathfrak{R}_{S_n}(\mathcal{G})] = \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_n} \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{x}_i) \right] \quad (19)$$

is defined as above.

Standard results used in this section (Talagrand’s contraction lemma, vector contraction, Rademacher complexity bounds for linear function classes, additional bounds for adversarial inputs, McDiarmid’s inequality, etc.) are summarized in Appendix A.1. (*Proofs of theorems and lemmas are provided in Appendix A.*)

## 5.2 Upper Bound on the Estimation Error of Supervised TRADES

In this section, we derive an upper bound on the estimation error of supervised TRADES. The proof proceeds by (i) deriving an upper bound on the uniform deviation (lemma), and (ii) applying the standard ERM argument to obtain the estimation-error bound (theorem), following the standard flow.

### *Supervised TRADES risk (population risk and empirical risk)*

We define the population risk of supervised TRADES as

$$R_{\text{PN-TR}}(g) := \pi_{\text{P}} \mathbb{E}_{\text{P}} \left[ \ell(g(\mathbf{x}), +1) + \beta \max_{\|\boldsymbol{\eta}\|_{\infty} \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x}), g(\mathbf{x} + \boldsymbol{\eta})) \right] \\ + \pi_{\text{N}} \mathbb{E}_{\text{N}} \left[ \ell(g(\mathbf{x}), -1) + \beta \max_{\|\boldsymbol{\eta}\|_{\infty} \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x}), g(\mathbf{x} + \boldsymbol{\eta})) \right] \quad (20)$$

as above. The corresponding empirical risk is

$$\hat{R}_{\text{PN-TR}}(g) := \frac{\pi_{\text{P}}}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \left[ \ell(g(\mathbf{x}_i^{\text{P}}), +1) + \beta \max_{\|\boldsymbol{\eta}\|_{\infty} \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x}_i^{\text{P}}), g(\mathbf{x}_i^{\text{P}} + \boldsymbol{\eta})) \right] \\ + \frac{\pi_{\text{N}}}{n_{\text{N}}} \sum_{i=1}^{n_{\text{N}}} \left[ \ell(g(\mathbf{x}_i^{\text{N}}), -1) + \beta \max_{\|\boldsymbol{\eta}\|_{\infty} \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x}_i^{\text{N}}), g(\mathbf{x}_i^{\text{N}} + \boldsymbol{\eta})) \right] \quad (21)$$

. Define the empirical risk minimizer as

$$\hat{g}_{\text{PN-TR}} := \arg \min_{g \in \mathcal{G}} \hat{R}_{\text{PN-TR}}(g)$$

Also, let

$$g^* \in \arg \min_{g \in \mathcal{G}} R_{\text{PN-TR}}(g)$$

be a population risk minimizer. With these definitions, we obtain the following upper bound on the estimation error of supervised TRADES.

**Theorem 5.2** (Upper Bound on the Estimation Error of Supervised TRADES). *Let a function class  $\mathcal{G}$  be given. Assume the following:*

- **(Boundedness of the loss)** *One of the following holds:*
  - *There exists a constant  $C_\ell > 0$  such that for any  $\hat{y} \in \mathbb{R}$  and  $y \in \mathcal{Y}$ ,  $\ell(\hat{y}, y) \leq C_\ell$  holds; or*
  - *there exists a constant  $C_g > 0$  such that  $\|g\|_\infty = \sup_{\mathbf{x} \in \mathcal{X}} |g(\mathbf{x})| \leq C_g$  holds for any  $g \in \mathcal{G}$ , and for  $|\hat{y}| \leq C_g$ ,  $\ell(\hat{y}, y) \leq C_\ell$  holds.*
- **(Lipschitz continuity of the loss)**  *$\ell(\hat{y}, y)$  is  $L_\ell$ -Lipschitz continuous with respect to  $\hat{y}$ .*
- **(Regularity of the KL term)** *The TRADES regularization term  $\ell_{\text{KL}}(g(\mathbf{x}), g(\mathbf{x}'))$  is  $L_{\text{KL}}$ -Lipschitz continuous in each argument, and is uniformly bounded:  $\ell_{\text{KL}}(u, v) \leq C_{\text{KL}}$  holds.*

*Then, for any  $\delta > 0$ , the following holds with probability at least  $1 - \delta$ :*

$$\begin{aligned} R_{\text{PN-TR}}(\hat{g}_{\text{PN-TR}}) - R_{\text{PN-TR}}(g^*) &\leq 4(L_\ell + 4\beta L_{\text{KL}}) \left( \pi_{\text{P}} \mathfrak{R}_{n_{\text{P}}, p_{\text{P}}}(\mathcal{G}) + \pi_{\text{N}} \mathfrak{R}_{n_{\text{N}}, p_{\text{N}}}(\mathcal{G}) \right) \\ &\quad + 8\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{\pi_{\text{N}}}{\sqrt{n_{\text{N}}}} \right) \\ &\quad + \sqrt{2 \ln \frac{2}{\delta}} (C_\ell + \beta C_{\text{KL}}) \left( \frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{\pi_{\text{N}}}{\sqrt{n_{\text{N}}}} \right). \end{aligned} \tag{22}$$

#### **Interpretation and implications (statistical convergence rate)**

Under the linear-in-parameters model ( $\|\mathbf{w}\|_p \leq W$ ) and bounded inputs, standard bounds in Appendix A.1 yield

$$\mathfrak{R}_{n_{\text{P}}, p_{\text{P}}}(\mathcal{G}) = \mathcal{O}\left(\frac{W}{\sqrt{n_{\text{P}}}}\right), \quad \mathfrak{R}_{n_{\text{N}}, p_{\text{N}}}(\mathcal{G}) = \mathcal{O}\left(\frac{W}{\sqrt{n_{\text{N}}}}\right)$$

Substituting this bound into Theorem 5.2 and absorbing constants such as  $\beta$  and  $\varepsilon$ , the estimation-error bound becomes

$$R_{\text{PN-TR}}(\hat{g}_{\text{PN-TR}}) - R_{\text{PN-TR}}(g^*) = \mathcal{O}_p\left(\frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{\pi_{\text{N}}}{\sqrt{n_{\text{N}}}}\right)$$

This implies an  $\mathcal{O}_p(\pi_{\text{P}}/\sqrt{n_{\text{P}}} + \pi_{\text{N}}/\sqrt{n_{\text{N}}})$  rate, and in particular the bound converges to 0 in probability as  $n_{\text{P}}, n_{\text{N}} \rightarrow \infty$ .

As preparation for Theorem 5.2, we first bound the uniform deviation for supervised TRADES.

### Auxiliary lemmas

In what follows, to bound the Rademacher complexity terms arising from the adversarial component of TRADES, we use two auxiliary lemmas: Lemma A.4 and Lemma A.2. (Hereafter,  $q$  denotes the conjugate exponent of  $p$  (i.e.,  $1/p + 1/q = 1$ ).)

**Lemma 5.3** (Upper Bound on the Uniform Deviation for Supervised TRADES). *Under the above assumptions, for any  $\delta > 0$ , with probability at least  $1 - \delta$ ,*

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left| \widehat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g) \right| &\leq 2(L_\ell + 4\beta L_{\text{KL}}) \left( \pi_{\text{P}} \mathfrak{R}_{n_{\text{P}}, p_{\text{P}}}(\mathcal{G}) + \pi_{\text{N}} \mathfrak{R}_{n_{\text{N}}, p_{\text{N}}}(\mathcal{G}) \right) \\ &\quad + 4\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{\pi_{\text{N}}}{\sqrt{n_{\text{N}}}} \right) \\ &\quad + \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} (C_\ell + \beta C_{\text{KL}}) \left( \frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{\pi_{\text{N}}}{\sqrt{n_{\text{N}}}} \right). \end{aligned} \quad (23)$$

holds.

(Proof is given in Appendix A.2.)

(Proof is given in Appendix A.3.)

### 5.3 Upper Bound on the Estimation Error of uPU+TRADES

In this section, we consider the uPU+TRADES objective based on positive (P) and unlabeled (U) data, and derive an upper bound on its estimation error.

#### uPU+TRADES risk (population risk and empirical risk)

Using the composite loss  $\tilde{\ell}(g(\mathbf{x}), y) = \ell(g(\mathbf{x}), y) - \ell(g(\mathbf{x}), -y)$  we define the population risk of uPU+TRADES as

$$\begin{aligned} R_{\text{uPU-TR}}(g) &:= \pi_{\text{P}} \mathbb{E}_{\text{P}} \left[ \tilde{\ell}(g(\mathbf{x}), +1) + \beta \max_{\|\boldsymbol{\eta}\|_\infty \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x} + \boldsymbol{\eta}), g(\mathbf{x})) \right] \\ &\quad + \mathbb{E}_{\text{U}} \left[ \ell(g(\mathbf{x}), -1) + \beta \max_{\|\boldsymbol{\eta}\|_\infty \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x} + \boldsymbol{\eta}), g(\mathbf{x})) \right] \end{aligned} \quad (24)$$

as above. The corresponding empirical risk is

$$\begin{aligned} \widehat{R}_{\text{uPU-TR}}(g) &:= \frac{\pi_{\text{P}}}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \left[ \tilde{\ell}(g(\mathbf{x}_i^{\text{P}}), +1) + \beta \max_{\|\boldsymbol{\eta}\|_\infty \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x}_i^{\text{P}} + \boldsymbol{\eta}), g(\mathbf{x}_i^{\text{P}})) \right] \\ &\quad + \frac{1}{n_{\text{U}}} \sum_{i=1}^{n_{\text{U}}} \left[ \ell(g(\mathbf{x}_i^{\text{U}}), -1) + \beta \max_{\|\boldsymbol{\eta}\|_\infty \leq \epsilon} \ell_{\text{KL}}(g(\mathbf{x}_i^{\text{U}} + \boldsymbol{\eta}), g(\mathbf{x}_i^{\text{U}})) \right] \end{aligned} \quad (25)$$

. Define the empirical risk minimizer as

$$\hat{g}_{\text{uPU-TR}} := \arg \min_{g \in \mathcal{G}} \hat{R}_{\text{uPU-TR}}(g)$$

Also, let

$$g^* \in \arg \min_{g \in \mathcal{G}} R_{\text{uPU-TR}}(g)$$

be a population risk minimizer.

**Theorem 5.4** (Upper Bound on the Estimation Error of uPU+TRADES). *Let a function class  $\mathcal{G}$  be given. Assume the following (as in the previous subsection):*

- **(Boundedness of the loss)** *There exist constants  $C_\ell, C_{\text{KL}} > 0$  such that for any  $y \in \mathcal{Y}$  and any input, the following holds:*
  - (Classification loss) *One of the following holds:*
    - \* *There exists a constant  $C_\ell > 0$  such that for any  $\hat{y} \in \mathbb{R}$ ,  $\ell(\hat{y}, y) \leq C_\ell$ , or*
    - \* *there exists a constant  $C_g > 0$  such that  $\|g\|_\infty \leq C_g$  (for  $g \in \mathcal{G}$ ) and for  $|\hat{y}| \leq C_g$ ,  $\ell(\hat{y}, y) \leq C_\ell$ .*
  - (TRADES term) *For any  $u, v$ ,  $\ell_{\text{KL}}(u, v) \leq C_{\text{KL}}$ .*
- **(Lipschitz continuity)** *There exist constants  $L_\ell, L_{\text{KL}} > 0$  such that*
  - (Classification loss)  *$\ell(\hat{y}, y)$  is  $L_\ell$ -Lipschitz continuous with respect to  $\hat{y}$ .*
  - (TRADES term)  *$\ell_{\text{KL}}(u, v)$  is  $L_{\text{KL}}$ -Lipschitz continuous in each argument*

*Then, for any  $\delta > 0$ , the following holds with probability at least  $1 - \delta$ :*

$$\begin{aligned} R_{\text{uPU-TR}}(\hat{g}_{\text{uPU-TR}}) - R_{\text{uPU-TR}}(g^*) &\leq 8\pi_{\text{P}}(L_\ell + 2\beta L_{\text{KL}})\mathfrak{R}_{n_{\text{P}}, p_{\text{P}}}(\mathcal{G}) \\ &\quad + 4(L_\ell + 4\beta L_{\text{KL}})\mathfrak{R}_{n_{\text{U}}, p_{\text{U}}}(\mathcal{G}) \\ &\quad + 8\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}} \right) \\ &\quad + \sqrt{2 \ln \frac{2}{\delta}} \left( \frac{\pi_{\text{P}}(2C_\ell + \beta C_{\text{KL}})}{\sqrt{n_{\text{P}}}} + \frac{C_\ell + \beta C_{\text{KL}}}{\sqrt{n_{\text{U}}}} \right). \end{aligned} \quad (26)$$

#### **Interpretation and implications (statistical convergence rate)**

Under the linear-in-parameters model ( $\|w\|_p \leq W$ ) and bounded inputs, standard bounds in Appendix A.1 yield

$$\mathfrak{R}_{n_{\text{P}}, p_{\text{P}}}(\mathcal{G}) = \mathcal{O}\left(\frac{W}{\sqrt{n_{\text{P}}}}\right), \quad \mathfrak{R}_{n_{\text{U}}, p_{\text{U}}}(\mathcal{G}) = \mathcal{O}\left(\frac{W}{\sqrt{n_{\text{U}}}}\right)$$

Substituting this bound into (26) and absorbing constants such as  $\beta$  and  $\varepsilon$ , the estimation-error upper bound becomes

$$R_{\text{uPU-TR}}(\hat{g}_{\text{uPU-TR}}) - R_{\text{uPU-TR}}(g^*) = \mathcal{O}_p\left(\frac{\pi_P}{\sqrt{n_P}} + \frac{1}{\sqrt{n_U}}\right)$$

This yields an upper bound that converges to 0 in probability as  $n_P, n_U \rightarrow \infty$ .

As preparation for Theorem 5.4, we present a lemma bounding the uniform deviation for uPU+TRADES.

#### Auxiliary lemmas

Below we use the auxiliary Lemma A.4 (Rademacher increase under adversarial inputs) and Lemma A.2 (vector contraction).

**Lemma 5.5** (Upper Bound on the Uniform Deviation for uPU+TRADES).  
For any  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left| \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \right| &\leq 4\pi_P (L_\ell + 2\beta L_{\text{KL}}) \mathfrak{R}_{n_P, p_P}(\mathcal{G}) \\ &\quad + 2(L_\ell + 4\beta L_{\text{KL}}) \mathfrak{R}_{n_U, p_U}(\mathcal{G}) \\ &\quad + 4\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_P}{\sqrt{n_P}} + \frac{1}{\sqrt{n_U}} \right) \\ &\quad + \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} \left( \frac{\pi_P (2C_\ell + \beta C_{\text{KL}})}{\sqrt{n_P}} + \frac{C_\ell + \beta C_{\text{KL}}}{\sqrt{n_U}} \right) \end{aligned} \quad (27)$$

(Proof is given in Appendix A.4.)

(Proof is given in Appendix A.5.)

## 5.4 Upper Bound on the Estimation Error of nnPU+TRADES

In this section, we derive an upper bound on the estimation error of nnPU+TRADES.

#### nnPU+TRADES risk (population risk and empirical risk)

First, define the adversarial regularization term in TRADES as

$$\psi(g, \mathbf{x}) := \max_{\|\boldsymbol{\eta}\|_\infty \leq \varepsilon} \ell_{\text{KL}}(g(\mathbf{x} + \boldsymbol{\eta}), g(\mathbf{x}))$$

. We define the population risk of nnPU+TRADES by adding this regularization term to the nnPU (Kiryó et al., 2017) risk estimator:

$$\begin{aligned} R_{\text{nnPU-TR}}(g) &:= \pi_P \mathbb{E}_P[\ell(g(\mathbf{x}), +1) + \beta \psi(g, \mathbf{x})] \\ &\quad + \max\left\{0, -\pi_P \mathbb{E}_P[\ell(g(\mathbf{x}), -1)] + \mathbb{E}_U[\ell(g(\mathbf{x}), -1)]\right\} \\ &\quad + \beta \mathbb{E}_U[\psi(g, \mathbf{x})] \end{aligned} \quad (28)$$



as above. The corresponding empirical risk is

$$\begin{aligned}\widehat{R}_{\text{nnPU-TR}}(g) &:= \frac{\pi_{\text{P}}}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \left[ \ell(g(\mathbf{x}_i^{\text{P}}), +1) + \beta \psi(g, \mathbf{x}_i^{\text{P}}) \right] \\ &\quad + \max \left\{ 0, -\frac{\pi_{\text{P}}}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \ell(g(\mathbf{x}_i^{\text{P}}), -1) + \frac{1}{n_{\text{U}}} \sum_{i=1}^{n_{\text{U}}} \ell(g(\mathbf{x}_i^{\text{U}}), -1) \right\} \\ &\quad + \frac{\beta}{n_{\text{U}}} \sum_{i=1}^{n_{\text{U}}} \psi(g, \mathbf{x}_i^{\text{U}})\end{aligned}\tag{29}$$

. Define the empirical risk minimizer as

$$\widehat{g}_{\text{nnPU-TR}} := \arg \min_{g \in \mathcal{G}} \widehat{R}_{\text{nnPU-TR}}(g)$$

Also, let

$$g^* \in \arg \min_{g \in \mathcal{G}} R_{\text{nnPU-TR}}(g)$$

be a population risk minimizer.

**Theorem 5.6** (Upper Bound on the Estimation Error of nnPU+TRADES). *Let a function class  $\mathcal{G}$  be given. Assume the following (as in the previous subsection):*

- **(Boundedness of the loss)** *There exist constants  $C_\ell, C_{\text{KL}} > 0$  such that for any  $y \in \mathcal{Y}$  and any input, the following holds:*
  - (Classification loss) *One of the following holds:*
    - \* *There exists a constant  $C_\ell > 0$  such that for any  $\widehat{y} \in \mathbb{R}$ ,  $\ell(\widehat{y}, y) \leq C_\ell$ , or*
    - \* *there exists a constant  $C_g > 0$  such that  $\|g\|_\infty \leq C_g$  (for  $g \in \mathcal{G}$ ) and for  $|\widehat{y}| \leq C_g$ ,  $\ell(\widehat{y}, y) \leq C_\ell$ .*
  - (TRADES term) *For any  $u, v$ ,  $\ell_{\text{KL}}(u, v) \leq C_{\text{KL}}$ .*
- **(Lipschitz continuity)** *There exist constants  $L_\ell, L_{\text{KL}} > 0$  such that*
  - (Classification loss)  *$\ell(\widehat{y}, y)$  is  $L_\ell$ -Lipschitz continuous with respect to  $\widehat{y}$ .*
  - (TRADES term)  *$\ell_{\text{KL}}(u, v)$  is  $L_{\text{KL}}$ -Lipschitz continuous in each argument*

*Then, for any  $\delta > 0$ , the following holds with probability at least  $1 - \delta$ :*

$$\begin{aligned}R_{\text{nnPU-TR}}(\widehat{g}_{\text{nnPU-TR}}) - R_{\text{nnPU-TR}}(g^*) &\leq 8\pi_{\text{P}}(L_\ell + 2\beta L_{\text{KL}})\mathfrak{R}_{n_{\text{P}}, \text{PP}}(\mathcal{G}) \\ &\quad + 4(L_\ell + 4\beta L_{\text{KL}})\mathfrak{R}_{n_{\text{U}}, \text{PU}}(\mathcal{G}) \\ &\quad + 8\beta L_{\text{KL}} \varepsilon_W d^{1/q} \left( \frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}} \right) \\ &\quad + \sqrt{2 \ln \frac{2}{\delta}} \left( \frac{\pi_{\text{P}}(2C_\ell + \beta C_{\text{KL}})}{\sqrt{n_{\text{P}}}} + \frac{C_\ell + \beta C_{\text{KL}}}{\sqrt{n_{\text{U}}}} \right).\end{aligned}\tag{30}$$

**Interpretation and implications (statistical convergence rate)**

Under the linear-in-parameters model ( $\|\mathbf{w}\|_p \leq W$ ) and bounded inputs, standard bounds in Appendix A.1 yield

$$\mathfrak{R}_{n_P, p_P}(\mathcal{G}) = \mathcal{O}\left(\frac{W}{\sqrt{n_P}}\right), \quad \mathfrak{R}_{n_U, p_U}(\mathcal{G}) = \mathcal{O}\left(\frac{W}{\sqrt{n_U}}\right)$$

Substituting this bound into (30) and absorbing constants such as  $\beta$  and  $\varepsilon$ , the estimation-error upper bound becomes

$$R_{\text{nnPU-TR}}(\hat{g}_{\text{nnPU-TR}}) - R_{\text{nnPU-TR}}(g^*) = \mathcal{O}_p\left(\frac{\pi_P}{\sqrt{n_P}} + \frac{1}{\sqrt{n_U}}\right)$$

This yields an upper bound that converges to 0 in probability as  $n_P, n_U \rightarrow \infty$ .

As preparation for Theorem 5.6, we provide a lemma bounding the uniform deviation for nnPU+TRADES.

**Auxiliary lemmas**

Below we use the auxiliary Lemma A.4 (Rademacher increase under adversarial inputs) and Lemma A.2 (vector contraction).

**Lemma 5.7** (Upper Bound on the Uniform Deviation for nnPU+TRADES).  
For any  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left| \hat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \right| &\leq 4\pi_P (L_\ell + 2\beta L_{\text{KL}}) \mathfrak{R}_{n_P, p_P}(\mathcal{G}) \\ &\quad + 2(L_\ell + 4\beta L_{\text{KL}}) \mathfrak{R}_{n_U, p_U}(\mathcal{G}) \\ &\quad + 4\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_P}{\sqrt{n_P}} + \frac{1}{\sqrt{n_U}} \right) \\ &\quad + \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} \left( \frac{\pi_P (2C_\ell + \beta C_{\text{KL}})}{\sqrt{n_P}} + \frac{C_\ell + \beta C_{\text{KL}}}{\sqrt{n_U}} \right). \end{aligned} \quad (31)$$

(Proof is given in Appendix A.6.)

(Proof is given in Appendix A.7.)

## 5.5 Sufficient Unlabeled Sample Size for PU+TRADES to Outperform Supervised TRADES

In this section, we compare the estimation-error upper bounds obtained in the previous subsections (Theorems 5.2, 5.4, and 5.6) and derive sufficient conditions on the unlabeled sample size  $n_U$  under which the bound for supervised TRADES (PN+TRADES) becomes larger than that for PU+TRADES (uPU+TRADES / nnPU+TRADES), i.e., conditions under which PU+TRADES can theoretically outperform supervised TRADES.

### *Rademacher complexity of the linear hypothesis class*

Under the assumptions in this section ( $\|\mathbf{x}\|_\infty \leq C_x$  and  $\|\mathbf{w}\|_p \leq W$ ), for any distribution  $\nu$ ,

$$\mathfrak{R}_{n,\nu}(\mathcal{G}) = \mathbb{E}_{\mathbf{x}_{1:n} \sim \nu} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\|\mathbf{w}\|_p \leq W} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbf{w}^\top \mathbf{x}_i \right] \leq \frac{W \sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_q}{\sqrt{n}} \leq \frac{WC_x d^{1/q}}{\sqrt{n}} \quad (32)$$

holds (where  $1/p + 1/q = 1$ ). In the following, for notational clarity,

$$\kappa_\delta := \sqrt{2 \ln \frac{2}{\delta}} \quad (33)$$

and we rewrite the estimation-error bounds into a  $1/\sqrt{n}$  form for comparison. Moreover, to collect constant factors,

$$\Gamma_\delta := 4(L_\ell + 4\beta L_{\text{KL}}) WC_x d^{1/q} + 8\beta L_{\text{KL}} \varepsilon W d^{1/q} + \kappa_\delta (C_\ell + \beta C_{\text{KL}}) \quad (34)$$

.

### *(1) Comparing PN+TRADES and uPU+TRADES*

First, applying (32) to Theorem 5.2 gives, with probability at least  $1 - \delta$ ,

$$R_{\text{PN-TR}}(\hat{g}_{\text{PN-TR}}) - R_{\text{PN-TR}}(g^*) \leq \Gamma_\delta \left( \frac{\pi_P}{\sqrt{n_P}} + \frac{\pi_N}{\sqrt{n_N}} \right) \quad (35)$$

is obtained.

Similarly, applying (32) to Theorem 5.4 yields, with probability at least  $1 - \delta$ ,

$$R_{\text{uPU-TR}}(\hat{g}_{\text{uPU-TR}}) - R_{\text{uPU-TR}}(g^*) \leq \frac{\pi_P}{\sqrt{n_P}} \left( \Gamma_\delta + 4L_\ell WC_x d^{1/q} + \kappa_\delta C_\ell \right) + \frac{\Gamma_\delta}{\sqrt{n_U}} \quad (36)$$

follows.

**Main result: sufficient unlabeled sample size**

The following theorem, based on the comparison of the estimation-error upper bounds, provides a sufficient condition on the unlabeled sample size  $n_U$  under which PU+TRADES can outperform supervised TRADES.

**Theorem 5.8** (Sufficient Unlabeled Sample Size for PU+TRADES to Outperform Supervised TRADES). *We compare the estimation-error upper bounds in Theorems 5.2, 5.4, and 5.6. Using (32), we reduce them to  $1/\sqrt{n}$ -type forms, and we use  $\Gamma_\delta$  and  $\kappa_\delta$  defined in (34). Then*

$$\Gamma_\delta \frac{\pi_N}{\sqrt{n_N}} > \frac{\pi_P}{\sqrt{n_P}} \left( 4L_\ell W C_x d^{1/q} + \kappa_\delta C_\ell \right) \quad (37)$$

*assume it holds.*

(i) (PN+TRADES vs uPU+TRADES) *If*

$$n_U > \left( \frac{\Gamma_\delta}{\Gamma_\delta \frac{\pi_N}{\sqrt{n_N}} - \frac{\pi_P}{\sqrt{n_P}} (4L_\ell W C_x d^{1/q} + \kappa_\delta C_\ell)} \right)^2 \quad (38)$$

*as a comparison of the estimation-error upper bounds, the bound for PN+TRADES ((35)) is larger than that for uPU+TRADES ((36)).*

(ii) (PN+TRADES vs nnPU+TRADES) *Since the right-hand side of Theorem 5.6 is identical to that of Theorem 5.4, the same sufficient condition*

$$n_U > \left( \frac{\Gamma_\delta}{\Gamma_\delta \frac{\pi_N}{\sqrt{n_N}} - \frac{\pi_P}{\sqrt{n_P}} (4L_\ell W C_x d^{1/q} + \kappa_\delta C_\ell)} \right)^2 \quad (39)$$

*implies that the PN+TRADES estimation-error upper bound is larger than the nnPU+TRADES one.*

(Proof is given in Appendix A.8.)

If condition (37) does not hold, the denominator on the right-hand side becomes non-positive, and thus there is no  $n_U$  satisfying the comparison inequality in (i); within the scope of comparisons based on the bounds derived in this section, we cannot claim that PU+TRADES outperforms supervised TRADES. On the other hand, when (37) holds, by (38) and (39), in the regime where the unlabeled sample size is sufficiently large, the superiority of PU+TRADES is theoretically guaranteed.

**Table 1** Experimental settings for each dataset (F-MNIST, CIFAR-10, CIFAR-100, and Alzheimer MRI).

Item	F-MNIST	CIFAR-10	CIFAR-100	Alzheimer
Dataset	FashionMNIST	CIFAR-10	CIFAR-100	Alzheimer MRI
Positive class	Top (0,2,4,6)	Vehicles (0,1,8,9)	organics	AD patients
# positive samples $n_p$	1,000	1,000	1,000	769
# unlabeled samples $n_u$	10,000	10,000	50,000	5,121
Class prior $\pi_p$	0.4	0.4	0.5	0.5
Model	6-layer MLP	ResNet-18	ResNet-18	ResNet-50
Perturbation budget $\epsilon$	0.3	0.031	0.031	0.031
Step size $\alpha$	0.01	0.031	0.031	0.031
# PGD iterations		10 steps		
TRADES coefficient $\beta$		6, 12, 18		
Metrics	Clean Accuracy, Adversarial Accuracy			

## 6 Experiments

In this section, we experimentally evaluate the effectiveness of PU+TRADES (uPU+TRADES and nnPU+TRADES), proposed in Section 4, on multiple benchmark datasets. Specifically, we compare our methods with existing approaches using clean accuracy (Clean Accuracy) and adversarial accuracy (Adversarial Accuracy) as evaluation metrics, and also analyze the effect of the TRADES coefficient  $\beta$ . Furthermore, we conduct additional experiments to examine how the theoretical threshold on the number of unlabeled samples, derived in Section 5, corresponds to empirical observations.

### 6.1 Experimental Setup

#### *Datasets.*

We conduct evaluation experiments on four benchmark datasets: FashionMNIST (F-MNIST [20]), CIFAR-10 [21], CIFAR-100 [21], and the Alzheimer dataset<sup>1</sup>. For F-MNIST, we consider clothing-image classification; for CIFAR-10, we focus on identifying vehicle classes; for CIFAR-100, we target **organics**; and for Alzheimer, we aim at recognizing AD patients. The definition of the positive class, the number of positive samples  $n_p$ , the number of unlabeled samples  $n_u$ , the class prior  $\pi_p$ , and the model architectures for each dataset are summarized in Table 1.

#### *Evaluation metrics.*

To compare the performance of each method, we report clean accuracy (Clean Accuracy) and adversarial accuracy (Adversarial Accuracy) on the test set. Accuracy is

<sup>1</sup>Dubey, S. *Alzheimer's Dataset* (Kaggle). <https://www.kaggle.com/tourist55/alzheimers-dataset-4-class-of-images>

**Table 2** Main results on FashionMNIST (6-layer MLP) and CIFAR-10 (ResNet-18). We compare uPU/nnPU and their adversarially trained variants (PGD training and TRADES training), and report clean accuracy (Clean) and adversarial accuracy (Adv.) under PGD attacks ( $\epsilon$  follows Table 1, 10 steps) on the test set. For each setting, we select the epoch that achieves the highest Clean Accuracy. While uPU/nnPU without adversarial training achieve almost zero Adv., TRADES improves robustness without severely degrading clean accuracy.

Method	F-MNIST		CIFAR-10	
	Clean	Adv.	Clean	Adv.
uPU	0.937	0.190	0.850	0.114
nnPU	<b>0.948</b>	0.001	<b>0.887</b>	0.001
uPU-PGD	0.935	0.843	0.802	0.714
nnPU-PGD	0.930	0.860	0.732	0.686
uPU+TRADES	0.934	0.914	0.845	0.711
nnPU+TRADES	0.944	<b>0.928</b>	0.850	<b>0.723</b>

defined as

$$\text{Accuracy} = \frac{\# \text{ correct predictions}}{\# \text{ total samples}}.$$

Adversarial accuracy is computed as the classification accuracy on adversarial examples generated by PGD; the perturbation budget  $\epsilon$ , step size  $\alpha$ , and the number of iterations (10 steps) follow Table 1.

#### *Training details.*

As listed in Table 1, we use a 6-layer MLP [22] for F-MNIST, ResNet-18 [23] for CIFAR-10/100, and ResNet-50 [23] for Alzheimer. As baselines, for F-MNIST and CIFAR-10 we use uPU and nnPU, as well as PGD-based adversarial training (uPU-PGD and nnPU-PGD) and TRADES-based adversarial training (uPU+TRADES and nnPU+TRADES). For CIFAR-100 and Alzheimer, training is unstable due to uPU-specific overfitting; thus, we restrict evaluation to nnPU and nnPU+TRADES. We set the TRADES coefficient to  $\beta \in \{6, 12, 18\}$ , and other adversarial-perturbation settings ( $\epsilon$ ,  $\alpha$ , and the number of PGD iterations) follow Table 1.

## 6.2 Main Results

Tables 2 and 3 report clean accuracy (Clean) and adversarial accuracy (Adv.) on each dataset. First, while uPU and nnPU achieve high clean accuracy, their adversarial accuracy is extremely low, indicating vulnerability to perturbations. In contrast, nnPU-PGD and (uPU/nnPU)+TRADES markedly improve adversarial accuracy, confirming gains in robustness.

Moreover, on F-MNIST and CIFAR-10, nnPU+TRADES consistently achieves higher adversarial accuracy than uPU+TRADES. On the other hand, methods with adversarial training tend to suffer a drop in clean accuracy, revealing a trade-off between accuracy and robustness. For CIFAR-100 and Alzheimer (Table 3), nnPU+TRADES maintains non-trivial adversarial accuracy while retaining reasonable clean accuracy on both datasets.

**Table 3** Main results on CIFAR-100 (ResNet-18) and Alzheimer MRI (ResNet-50) (nnPU variants only). We compare clean accuracy (Clean) and adversarial accuracy (Adv.) under PGD attacks ( $\epsilon$  follows Table 1, 10 steps) on the test set. For each setting, we select the epoch that achieves the highest Clean Accuracy. While nnPU alone attains almost zero Adv., nnPU+TRADES substantially improves robustness.

Method	CIFAR-100		Alzheimer	
	Clean	Adv.	Clean	Adv.
nnPU	<b>0.680</b>	0.001	<b>0.683</b>	0.000
nnPU+TRADES	0.645	<b>0.450</b>	0.649	<b>0.409</b>

**Table 4** Ablation on the TRADES coefficient  $\beta$  (F-MNIST / CIFAR-10). We vary  $\beta \in \{6, 12, 18\}$  and compare clean accuracy (Clean Acc.) and adversarial accuracy (Adv. Acc.) for uPU+TRADES and nnPU+TRADES. On FashionMNIST, performance changes only mildly with  $\beta$ , whereas on CIFAR-10 increasing  $\beta$  tends to improve adversarial accuracy at the cost of decreased clean accuracy.

Method	F-MNIST		CIFAR-10	
	Clean Acc.	Adv. Acc.	Clean Acc.	Adv. Acc.
uPU+TRADES ( $\beta = 6$ )	0.939	0.861	0.846	0.711
uPU+TRADES ( $\beta = 12$ )	0.935	0.848	0.831	0.735
uPU+TRADES ( $\beta = 18$ )	0.937	0.870	0.809	0.735
nnPU+TRADES ( $\beta = 6$ )	<b>0.944</b>	<b>0.877</b>	<b>0.861</b>	0.723
nnPU+TRADES ( $\beta = 12$ )	0.933	0.875	0.845	0.740
nnPU+TRADES ( $\beta = 18$ )	0.930	0.871	0.836	<b>0.743</b>

Table 4 shows performance as we vary the TRADES coefficient  $\beta$ . On CIFAR-10, adversarial accuracy tends to increase as  $\beta$  becomes larger, while clean accuracy decreases step by step. For example, for uPU+TRADES, increasing  $\beta$  from 6 to 12 improves adversarial accuracy but reduces clean accuracy, and at  $\beta = 18$  the adversarial accuracy appears to saturate. Similarly for nnPU+TRADES, increasing  $\beta$  improves adversarial accuracy but also induces a decrease in clean accuracy, indicating that  $\beta$  is a key factor controlling the accuracy–robustness trade-off.

On F-MNIST, the variation across  $\beta$  is smaller than on CIFAR-10, and in particular uPU+TRADES exhibits only limited changes in clean accuracy even as  $\beta$  increases. However, for nnPU+TRADES, setting  $\beta$  too large slightly degrades adversarial accuracy, suggesting that the optimal  $\beta$  may depend on the dataset and the learning scheme (uPU vs. nnPU).

### 6.3 Validation of Theoretical Analysis

#### *Objective.*

In Section 5, as a condition under which PU learning can become advantageous compared to supervised learning (PN), we compared the upper bounds on the estimation error of PN+TRADES and PU+TRADES, and derived a threshold on the number of unlabeled samples,  $n_U^*$ , which holds when the upper bound for PN+TRADES exceeds that of PU+TRADES (Theorem 5.8). The aim of this section is to examine to what extent this theoretical threshold aligns with experimental results as a qualitative guideline: “increasing the amount of unlabeled data can make PU methods preferable.” Note that the theoretical threshold is a sufficient condition based on upper-bound comparison; therefore, we do not generally expect it to exactly match the empirical turning point.

#### *Validation procedure.*

We validate the theory in three steps: (i) substitute constants into the theoretical expression to obtain a numerical value of  $n_U^*$ ; (ii) sweep  $n_U$  over discrete values and compare the test losses of PU+TRADES and PN+TRADES; (iii) compare the empirical turning point  $\hat{n}_U^{\text{emp}}$  with  $n_U^*$  and discuss reasons for the discrepancy.

#### *Numerical instantiation of the theoretical threshold $n_U^*$ .*

The theoretical expression contains an upper bound on the weight norm  $W$  (i.e.,  $\|w\|_2 \leq W$ ), but the threshold used here is based on comparing estimation-error upper bounds. We define the norm upper bound of the perturbed input as

$$C_x^{\text{adv}} := C_x + \varepsilon d^{1/q}, \quad C_{\text{KL}} := W C_x^{\text{adv}} = W(C_x + \varepsilon d^{1/q}) \quad (40)$$

and approximate an upper bound of the logistic loss by  $C_\ell \approx W C_x$ . Then we can rewrite  $\Gamma_\delta \approx W \bar{\Gamma}_\delta$ , and factor out  $W$  from both sides of the comparison inequality. As a result, the numerical instantiation of the threshold simplifies to a form independent of  $W$ .

Below we show intermediate steps to obtain  $\Gamma_\delta \approx W \bar{\Gamma}_\delta$ . The quantity  $\Gamma_\delta$  defined in Theorem 5.8 is

$$\Gamma_\delta = 4(L_\ell + 4\beta L_{\text{KL}}) W C_x d^{1/q} + 8\beta L_{\text{KL}} \varepsilon W d^{1/q} + \kappa_\delta (C_\ell + \beta C_{\text{KL}}) \quad (41)$$

$$\approx 4(L_\ell + 4\beta L_{\text{KL}}) W C_x d^{1/q} + 8\beta L_{\text{KL}} \varepsilon W d^{1/q} + \kappa_\delta (W C_x + \beta W C_x^{\text{adv}}) \quad (42)$$

$$= 4(L_\ell + 4\beta L_{\text{KL}}) W C_x d^{1/q} + 8\beta L_{\text{KL}} \varepsilon W d^{1/q} + \kappa_\delta W (C_x + \beta C_x^{\text{adv}}) \quad (43)$$

$$= W \left\{ 4(L_\ell + 4\beta L_{\text{KL}}) C_x d^{1/q} + 8\beta L_{\text{KL}} \varepsilon d^{1/q} + \kappa_\delta (C_x + \beta C_x^{\text{adv}}) \right\} \quad (44)$$

$$=: W \bar{\Gamma}_\delta. \quad (45)$$

Specifically, letting

$$\kappa_\delta := \sqrt{2 \ln \frac{2}{\delta}}, \quad \bar{\Gamma}_\delta := 4(L_\ell + 4\beta L_{\text{KL}}) C_x d^{1/q} + 8\beta L_{\text{KL}} \varepsilon d^{1/q} + \kappa_\delta (C_x + \beta C_x^{\text{adv}}) \quad (46)$$



**Table 5** Constants used to numerically instantiate the theoretical threshold  $n_U^*$  derived in Theorem 5.8 (FashionMNIST / CIFAR-10). We follow Table 1 for the class priors and sample sizes, and summarize the input dimension  $d$ , norm bounds  $C_x, C_x^{\text{adv}}$ , Lipschitz constants  $L_\ell, L_{\text{KL}}$ , confidence level  $\delta$ , and so on.

Symbol	Meaning	F-MNIST	CIFAR-10	Remarks
$\pi_P$	prior of positive class	0.4	0.4	assumption
$\pi_N$	prior of negative class	0.6	0.6	$\pi_N = 1 - \pi_P$
$n_P$	# positive samples	1000	1000	given
$n_N$	# negative samples	500	500	given
$\beta$	TRADES coefficient	6	6	representative
$\epsilon$	perturbation radius	0.3	0.031	attack radius
$q$	norm index	2	2	dual of $p = 2$
$d$	input dimension	784	3072	$28^2, 32^2 \times 3$
$C_x$	input norm bound	$\sqrt{d}$	$\sqrt{d}$	$\ x\ _2 \leq \sqrt{d}$
$C_x^{\text{adv}}$	perturbed bound	$C_x + \epsilon d^{1/q}$	$C_x + \epsilon d^{1/q}$	(40)
$L_\ell$	Lipschitz (loss)	1	1	logistic
$L_{\text{KL}}$	Lipschitz (KL)	1	1	Sec. 5
$\delta$	confidence level	0.05	0.05	fixed

the feasibility condition (37) in Theorem 5.8 becomes

$$\bar{\Gamma}_\delta \frac{\pi_N}{\sqrt{n_N}} > \frac{\pi_P}{\sqrt{n_P}} \left( 4L_\ell C_x d^{1/q} + \kappa_\delta C_x \right), \quad (47)$$

and under this condition we obtain

$$n_U > \left( \frac{\bar{\Gamma}_\delta}{\bar{\Gamma}_\delta \frac{\pi_N}{\sqrt{n_N}} - \frac{\pi_P}{\sqrt{n_P}} (4L_\ell C_x d^{1/q} + \kappa_\delta C_x)} \right)^2 \quad (48)$$

When this condition is satisfied, comparing the estimation-error upper bounds of PN+TRADES and PU+TRADES shows that the upper bound for PN+TRADES exceeds that of PU+TRADES. That is, in theory, beyond this threshold PU+TRADES (uPU/nnPU) can become preferable to PN+TRADES.

#### **Constants and settings.**

The main symbols and constants used for numerical instantiation are summarized in Table 5. In this section we use  $q = 2$ , and for the input norm bound we set  $C_x = \sqrt{d}$  (i.e.,  $\|x\|_2 \leq \sqrt{d}$ ). Since  $d^{1/q} = \sqrt{d}$ , we have

$$C_x d^{1/q} = \sqrt{d} \cdot \sqrt{d} = d. \quad (49)$$

Thus, we substitute  $d = 784$  for F-MNIST and  $d = 3072$  for CIFAR-10.

Substituting the above settings into (48), we obtain

$$n_U^*(\text{F-MNIST}) \approx 1443.37 (\Rightarrow n_U \geq 1444),$$

$$n_U^*(\text{CIFAR-10}) \approx 1443.25 (\Rightarrow n_U \geq 1444).$$

***Experimental protocol for validation.***

Since the theoretical comparison concerns estimation error (risk), in this section we use the loss on the test set as a proxy. Specifically, for each  $n_U$ , let  $t^*(n_U)$  denote the epoch achieving the highest Clean Accuracy, and we compare the sum of the clean loss and adversarial loss at that epoch:

$$\mathcal{L}_{\text{sum}}(n_U) := \mathcal{L}_{\text{clean}}(t^*(n_U); n_U) + \mathcal{L}_{\text{adv}}(t^*(n_U); n_U). \quad (50)$$

Here,  $\mathcal{L}_{\text{clean}}$  is the clean loss on the test set, and  $\mathcal{L}_{\text{adv}}$  is the adversarial loss computed on adversarial examples generated from the same test set. We consider the following settings: (a) FMNIST: 6-layer MLP and a linear model; (b) CIFAR-10: a linear model and ResNet-18. For each setting, we train PU+TRADES and PN+TRADES and compare how  $\mathcal{L}_{\text{sum}}$  changes with  $n_U$ .

***Empirical turning point.***

To relate to the theoretical threshold  $n_U^*$ , we define the empirical turning point as the smallest  $n_U$  at which the sign of

$$\Delta(n_U) := \mathcal{L}_{\text{sum}}^{\text{PU-TRADES}}(n_U) - \mathcal{L}_{\text{sum}}^{\text{PN-TRADES}}(n_U)$$

becomes negative. Moreover, if the sign of  $\Delta$  flips between  $n_U = U_k$  and  $U_{k+1}$ , we estimate the turning point by linear interpolation:

$$\hat{n}_U^{\text{emp}} := U_k + \frac{-\Delta(U_k)}{\Delta(U_{k+1}) - \Delta(U_k)}(U_{k+1} - U_k). \quad (51)$$

***Results.***

Figure 7 shows the trajectories of  $\mathcal{L}_{\text{sum}}$  as a function of  $n_U$ . The orange curve represents  $\mathcal{L}_{\text{sum}}(n_U)$  for PU+TRADES, the dashed gray line is the baseline value for PN+TRADES, and the blue vertical line indicates the theoretical threshold  $n_U^*$ . When a crossing is observed within the sweep range, we mark  $\hat{n}_U^{\text{emp}}$  estimated by (51) as a point in the figure.

***Discussion.***

From Fig. 7 and Table 6, for F-MNIST (both the linear model and the 6-layer MLP),  $\mathcal{L}_{\text{sum}}(n_U)$  decreases as  $n_U$  increases, and the intersection where the loss sums of PU+TRADES and PN+TRADES coincide is estimated as  $\hat{n}_U^{\text{emp}} \approx 4247.13$  (linear

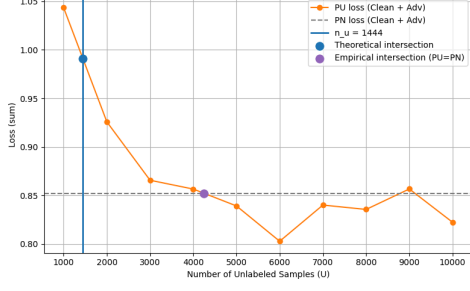


Fig. 3 \*

(a) FMNIST / Linear

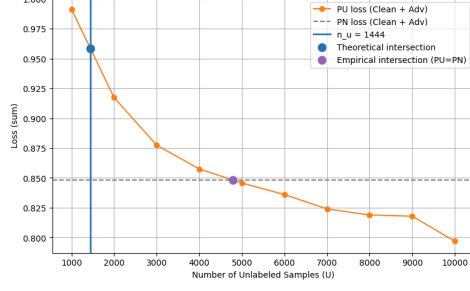


Fig. 4 \*

(b) FMNIST / 6-layer MLP

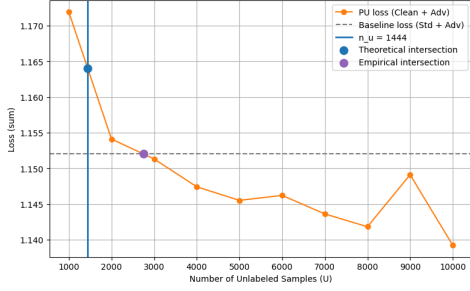


Fig. 5 \*

(c) CIFAR-10 / Linear

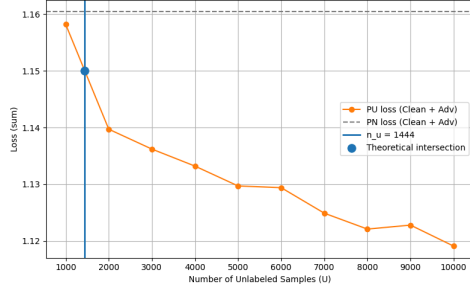


Fig. 6 \*

(d) CIFAR-10 / ResNet-18

**Fig. 7** Validation results for the theoretical threshold  $n_U^*$  of the unlabeled sample size derived in our theoretical analysis. Each panel plots the test loss sum for PU+TRADES,  $\mathcal{L}_{\text{sum}}(n_U) = \mathcal{L}_{\text{clean}}(t^*(n_U); n_U) + \mathcal{L}_{\text{adv}}(t^*(n_U); n_U)$  (Eq. (50)), against  $n_U$ , and compares it with the PN+TRADES baseline value (which does not depend on  $n_U$ ). Here,  $t^*(n_U)$  denotes the epoch achieving the highest Clean Accuracy. The intersection  $\hat{n}_U^{\text{emp}}$  is estimated by linear interpolation.

model) and  $\hat{n}_U^{\text{emp}} \approx 4788.14$  (6-layer MLP). These values are larger than the theoretical threshold  $n_U^* = 1444$ , which is consistent with the fact that  $n_U^*$  is a sufficient condition derived from upper-bound comparison (and hence can be conservative).

For CIFAR-10, the results depend on the model: (i) **CIFAR-10 / Linear**: a crossing is observed at  $\hat{n}_U^{\text{emp}} \approx 2750.00$  (between  $U = 2000$  and  $3000$ ). (ii) **CIFAR-10 / ResNet-18**: using the PN baseline value 1.1604, the PU loss sum remains smaller throughout the range  $U = 1000$ – $10000$ , and thus there is no intersection within this range (i.e., PU is already better at the smallest observed value  $U = 1000$ ). Overall, while the theoretical threshold  $n_U^*$  is not an exact predictor of the intersection location, the experiments also indicate that with sufficiently many unlabeled samples, PU methods can outperform PN methods.

**Table 6** Comparison between the theoretical threshold  $n_U^*$  (Theorem 5.8) and the empirical turning point  $\hat{n}_U^{\text{emp}}$ .  $\hat{n}_U^{\text{emp}}$  is estimated by linear interpolation from Fig. 7 as the value of  $n_U$  where the loss sums  $\mathcal{L}_{\text{sum}}$  of PU+TRADES and PN+TRADES coincide. “—” indicates that no intersection existed within the range of  $n_U$  evaluated in this study.

Setting	$n_U^*$ (theory)	$\hat{n}_U^{\text{emp}}$ (emp.)
FMNIST / 6-layer MLP	1444	4788.14
FMNIST / Linear	1444	4247.13
CIFAR-10 / Linear	1444	2750.00
CIFAR-10 / ResNet-18	1444	—

## 7 Conclusion

In this study, we focus on Positive-Unlabeled (PU) learning, which trains a classifier using only positive and unlabeled data. Motivated by applications such as medical image analysis—where fully annotating data is difficult and even small perturbations can lead to critical misclassifications—we aim to improve robustness against adversarial perturbations. First, we confirm that naively applying adversarial training to PU learning can make optimization unstable and significantly degrade clean performance, because unlabeled data, despite containing a mixture of positive and negative samples, tends to be treated uniformly as negative in the loss. To address this issue, we propose **PU+TRADES**, an extension of TRADES, a representative adversarial training framework, by integrating a PU loss with a label-independent regularization term. Experiments on multiple benchmark datasets and medical imaging data demonstrate that while standard PU learning can achieve high clean accuracy, its adversarial accuracy may drop drastically; in contrast, the proposed method improves adversarial accuracy without substantially sacrificing clean accuracy. Moreover, we show that adjusting the TRADES coefficient allows one to control the trade-off between robustness and clean performance. In addition, we derive an upper bound on the estimation error for binary classification under adversarial perturbations, and by comparing supervised learning with PU learning, we theoretically clarify conditions under which PU learning is advantageous, expressed as an inequality involving the numbers of positive (P), negative (N), and unlabeled (U) samples. Future work includes designing robust learning schemes that incorporate estimation error when the class prior is unknown, and extending the theory to deep models beyond linear assumptions.

## References

- [1] Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. CoRR (2014)
- [2] Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083

(2017)

- [3] Zhang, H., Yu, Y., Jiao, J., Xing, E.P., El Ghaoui, L., Jordan, M.I.: Theoretically principled trade-off between robustness and accuracy. In: Proceedings of the 36th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 97, pp. 7472–7482 (2019)
- [4] Plessis, M.C., Niu, G., Sugiyama, M.: Analysis of learning from positive and unlabeled data. In: Advances in Neural Information Processing Systems, vol. 27, pp. 703–711 (2014). <https://api.semanticscholar.org/CorpusID:14878885>
- [5] Plessis, M.C., Niu, G., Sugiyama, M.: Convex formulation for learning from positive and unlabeled data. In: Proceedings of the 32nd International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 37, pp. 1386–1394 (2015). <https://api.semanticscholar.org/CorpusID:15313053>
- [6] Chen, X., Chen, W., Chen, T., Yuan, Y., Gong, C., Chen, K., Wang, Z.: Self-pu: Self boosted and calibrated positive-unlabeled training. In: Proceedings of the 37th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 119, pp. 1510–1519. PMLR, ??? (2020). <https://proceedings.mlr.press/v119/chen20b.html>
- [7] Kiryo, R., Niu, G., Plessis, M.C., Sugiyama, M.: Positive-unlabeled learning with non-negative risk estimator. arXiv preprint arXiv:1703.00593 (2017)
- [8] Kato, M., Teshima, T., Honda, J.: Learning from positive and unlabeled data with a selection bias. In: International Conference on Learning Representations (2019). <https://openreview.net/forum?id=rJzLciCqKm>
- [9] Sakai, T., Shimizu, N.: Covariate shift adaptation on learning from positive and unlabeled data. In: AAAI Conference on Artificial Intelligence, pp. 4838–4845 (2019). <https://api.semanticscholar.org/CorpusID:70299136>
- [10] Bekker, J., Robberechts, P., Davis, J.: Beyond the selected completely at random assumption for learning from positive and unlabeled data. In: Machine Learning and Knowledge Discovery in Databases (ECML PKDD 2019), Proceedings, Part II. Lecture Notes in Computer Science (LNAI), vol. 11907, pp. 71–85. Springer, ??? (2020). [https://doi.org/10.1007/978-3-030-46147-8\\_5](https://doi.org/10.1007/978-3-030-46147-8_5)
- [11] Hammoudeh, Z., Lowd, D.: Learning from positive and unlabeled data with arbitrary positive shift. In: Advances in Neural Information Processing Systems, vol. 33, pp. 13088–13099 (2020). <https://arxiv.org/abs/2002.10261>
- [12] Dorigatti, E., Schweisthal, J., Bischl, B., Rezaei, M.: Robust and Efficient Imbalanced Positive-Unlabeled Learning with Self-supervision (2022). <https://arxiv.org/abs/2209.02459>

- [13] Liu, B., Lee, W.S., Yu, P.S., Li, X.: Partially supervised classification of text documents. In: Proceedings of the Nineteenth International Conference on Machine Learning (ICML 2002), pp. 387–394 (2002)
- [14] Hou, M., Chaib-draa, B., Li, C., Zhao, Q.: Generative adversarial positive-unlabelled learning. In: Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence (IJCAI-18), pp. 2255–2261 (2018). <https://doi.org/10.24963/ijcai.2018/312> . <https://www.ijcai.org/proceedings/2018/312>
- [15] Zhao, Y., Xu, Q., Jiang, Y., Wen, P., Huang, Q.: Dist-pu: Positive-unlabeled learning from a label distribution perspective. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 14441–14450 (2022)
- [16] Yu, H., Han, J., Chang, K.: Pebl: Web page classification without negative examples. IEEE Transactions on Knowledge and Data Engineering **16**(1), 70–81 (2004) <https://doi.org/10.1109/TKDE.2004.1264823>
- [17] Hsieh, Y.-G., Niu, G., Sugiyama, M.: Classification from positive, unlabeled and biased negative data. In: Proceedings of the 36th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 97, pp. 2820–2829. PMLR, ??? (2019). <https://proceedings.mlr.press/v97/hsieh19c.html>
- [18] Luo, C., Zhao, P., Chen, C., Qiao, B., Du, C., Zhang, H., Wu, W., Cai, S., He, B., Rajmohan, S., Lin, Q.: Pulns: Positive-unlabeled learning with effective negative sample selector. Proceedings of the AAAI Conference on Artificial Intelligence **35**, 8784–8792 (2021) <https://doi.org/10.1609/aaai.v35i10.17064>
- [19] Hu, W., Le, R., Liu, B., Ji, F., Ma, J., Zhao, D., Yan, R.: Predictive adversarial learning from positive and unlabeled data. Proceedings of the AAAI Conference on Artificial Intelligence **35**, 7806–7814 (2021) <https://doi.org/10.1609/aaai.v35i9.16953>
- [20] Xiao, H., Rasul, K., Vollgraf, R.: Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. arXiv preprint arXiv:1708.07747 (2017)
- [21] Krizhevsky, A.: Learning multiple layers of features from tiny images. Technical report, University of Toronto (April 2009). Available at: cs.toronto.edu
- [22] Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. MIT Press, Cambridge, MA (2016). <http://www.deeplearningbook.org>
- [23] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–778 (2016). <https://doi.org/10.1109/CVPR.2016.90>

## Appendix A Proofs of Theoretical Analysis

### A.1 Background / Tool Box (Standard Results)

In this appendix, we summarize standard lemmas and inequalities used in the proofs of this paper.

**Lemma A.1** (Talagrand's Contraction Lemma). *Let  $S_n = \{\mathbf{x}_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} \nu(\mathbf{x})$ . Suppose  $f : \mathbb{R} \rightarrow \mathbb{R}$  is  $L_f$ -Lipschitz. Then, for the Rademacher complexities of  $\mathcal{G}$  and  $f \circ \mathcal{G}$ ,*

$$\mathfrak{R}_{n,\nu}(f \circ \mathcal{G}) \leq L_f \mathfrak{R}_{n,\nu}(\mathcal{G}) \quad (\text{A1})$$

*holds.*

**Lemma A.2** (Vector Contraction). *If  $\ell(u, v)$  is  $L_\ell$ -Lipschitz in each argument, then for any pair of functions  $(f_h, g_h)$ ,*

$$\begin{aligned} \mathbb{E}_\sigma \left[ \sup_h \frac{1}{n} \sum_{i=1}^n \sigma_i \ell(f_h(\mathbf{x}_i), g_h(\mathbf{x}_i)) \right] &\leq 2L_\ell \left\{ \mathbb{E}_\sigma \left[ \sup_h \frac{1}{n} \sum_{i=1}^n \sigma_i f_h(\mathbf{x}_i) \right] \right. \\ &\quad \left. + \mathbb{E}_\sigma \left[ \sup_h \frac{1}{n} \sum_{i=1}^n \sigma_i g_h(\mathbf{x}_i) \right] \right\}. \end{aligned} \quad (\text{A2})$$

**Theorem A.3** (Upper Bound on Rademacher Complexity). *Assume the input satisfies  $\|\mathbf{x}\|_\infty \leq C_x$ , and let the class of linear classifiers be  $\mathcal{G} = \{\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} \mid \|\mathbf{w}\|_\infty \leq W\}$ . Then,*

$$\mathfrak{R}_{n,\nu}(\mathcal{G}) \leq \frac{C_x W}{\sqrt{n}} \quad (\text{A3})$$

*holds.*

**Lemma A.4** (Adversarial Rademacher Additive Term). *For the linear class  $\mathcal{G} = \{\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} : \|\mathbf{w}\|_p \leq W\}$ , for any distribution  $\nu$  and any  $n \in \mathbb{N}$ ,*

$$\mathfrak{R}_{n,\nu} \left( \left\{ \mathbf{x} \mapsto g(\mathbf{x} + \boldsymbol{\eta}) : \|\boldsymbol{\eta}\|_\infty \leq \varepsilon, g \in \mathcal{G} \right\} \right) \leq \mathfrak{R}_{n,\nu}(\mathcal{G}) + \frac{\varepsilon W d^{1/q}}{\sqrt{n}}.$$

**Theorem A.5** (McDiarmid’s Inequality). *Let  $X_1, \dots, X_n$  be independent random variables taking values in a set  $\mathcal{X}$ , and consider a function  $f : \mathcal{X}^n \rightarrow \mathbb{R}$ . Assume that there exist constants  $c_1, \dots, c_n$  such that, for each  $i = 1, \dots, n$  and any  $\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}'_i \in \mathcal{X}$ ,*

$$|f(\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n) - f(\mathbf{x}_1, \dots, \mathbf{x}'_i, \dots, \mathbf{x}_n)| \leq c_i \quad (\text{A4})$$

*holds. Then, for any  $t > 0$ ,*

$$\Pr(f(X_1, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X_n)] \geq t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right) \quad (\text{A5})$$

*holds.*

## A.2 Proof of the Lemma “Uniform Deviation Bound for Supervised TRADES” (Lemma 5.3)

*Proof* We divide the proof into the following steps (i)–(iv). We evaluate  $\sup_{g \in \mathcal{G}} \{\widehat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g)\}$  and  $\sup_{g \in \mathcal{G}} \{R_{\text{PN-TR}}(g) - \widehat{R}_{\text{PN-TR}}(g)\}$  with probability at least  $1 - \delta/2$ , respectively, and then combine them via a union bound to obtain an upper bound on the absolute value. Below, we only show  $\sup_{g \in \mathcal{G}} \{\widehat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g)\}$  (the reverse direction can be proved analogously).

### (i) McDiarmid’s inequality

Since  $\ell(\cdot) \leq C_\ell$  and  $\ell_{\text{KL}}(\cdot) \leq C_{\text{KL}}$ , replacing one sample on the P side changes  $\widehat{R}_{\text{PN-TR}}(g)$  by at most  $\frac{\pi_{\text{P}}}{n_{\text{P}}}(C_\ell + \beta C_{\text{KL}})$ , and replacing one sample on the N side changes it by at most  $\frac{\pi_{\text{N}}}{n_{\text{N}}}(C_\ell + \beta C_{\text{KL}})$ . Therefore, by McDiarmid’s inequality (A.5), for any  $t > 0$ ,

$$\begin{aligned} \Pr\left(\sup_{g \in \mathcal{G}} \{\widehat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g)\} - \mathbb{E}\left[\sup_{g \in \mathcal{G}} \{\widehat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g)\}\right] \geq t\right) \\ \leq \exp\left(-\frac{2t^2}{(C_\ell + \beta C_{\text{KL}})^2(\pi_{\text{P}}^2/n_{\text{P}} + \pi_{\text{N}}^2/n_{\text{N}})}\right). \end{aligned} \quad (\text{A6})$$

Setting  $\delta/2 = \exp\left(-2t^2/((C_\ell + \beta C_{\text{KL}})^2(\pi_{\text{P}}^2/n_{\text{P}} + \pi_{\text{N}}^2/n_{\text{N}}))\right)$  and solving for  $t$ , and then using the subadditivity of the square root, we obtain, with probability at least  $1 - \delta/2$ ,

$$\begin{aligned} \sup_{g \in \mathcal{G}} \{\widehat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g)\} &\leq \mathbb{E}\left[\sup_{g \in \mathcal{G}} \{\widehat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g)\}\right] \\ &\quad + \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} (C_\ell + \beta C_{\text{KL}}) \left(\frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{\pi_{\text{N}}}{\sqrt{n_{\text{N}}}}\right) \end{aligned} \quad (\text{A7})$$

as desired.



**(ii) Ghost sampling and symmetrization**

We use the symmetrization technique in statistical learning theory [24]. First, in

$$\mathbb{E} \left[ \sup_{g \in \mathcal{G}} \hat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g) \right],$$

the expectation  $\mathbb{E}$  is taken over repeated sampling of  $(\mathcal{X}_P, \mathcal{X}_N)$  used to evaluate  $\hat{R}_{\text{PN-TR}}(g)$ . To make this explicit, we write

$$\hat{R}_{\text{PN-TR}}(g) = \hat{R}_{\text{PN-TR}}(g; \mathcal{X}_P, \mathcal{X}_N).$$

Let  $(\mathcal{X}'_P, \mathcal{X}'_N)$  be a ghost sample (an independent copy with the same distribution) independent of  $(\mathcal{X}_P, \mathcal{X}_N)$ . Since the true risk can be written as the expectation over the ghost sample,

$$R_{\text{PN-TR}}(g) = \mathbb{E}_{(\mathcal{X}'_P, \mathcal{X}'_N)} [\hat{R}_{\text{PN-TR}}(g; \mathcal{X}'_P, \mathcal{X}'_N)],$$

we obtain

$$\begin{aligned} & \mathbb{E}_{(\mathcal{X}_P, \mathcal{X}_N)} \left[ \sup_{g \in \mathcal{G}} \hat{R}_{\text{PN-TR}}(g; \mathcal{X}_P, \mathcal{X}_N) - R_{\text{PN-TR}}(g) \right] \\ &= \mathbb{E}_{(\mathcal{X}_P, \mathcal{X}_N)} \left[ \sup_{g \in \mathcal{G}} \hat{R}_{\text{PN-TR}}(g; \mathcal{X}_P, \mathcal{X}_N) - \mathbb{E}_{(\mathcal{X}'_P, \mathcal{X}'_N)} \hat{R}_{\text{PN-TR}}(g; \mathcal{X}'_P, \mathcal{X}'_N) \right] \\ &\leq \mathbb{E}_{(\mathcal{X}_P, \mathcal{X}_N), (\mathcal{X}'_P, \mathcal{X}'_N)} \left[ \sup_{g \in \mathcal{G}} \hat{R}_{\text{PN-TR}}(g; \mathcal{X}_P, \mathcal{X}_N) - \hat{R}_{\text{PN-TR}}(g; \mathcal{X}'_P, \mathcal{X}'_N) \right], \end{aligned} \quad (\text{A8})$$

where we used Jensen's inequality, since  $\sup$  is a convex function.

Next, we decompose the difference into the P and N parts. For notational simplicity, define

$$\phi_y(g, \mathbf{x}) := \ell(g(\mathbf{x}), y) + \beta \max_{\|\boldsymbol{\eta}\|_\infty \leq \varepsilon} \ell_{\text{KL}}(g(\mathbf{x}), g(\mathbf{x} + \boldsymbol{\eta})).$$

Then,

$$\begin{aligned} & \hat{R}_{\text{PN-TR}}(g; \mathcal{X}_P, \mathcal{X}_N) - \hat{R}_{\text{PN-TR}}(g; \mathcal{X}'_P, \mathcal{X}'_N) \\ &= \frac{\pi_P}{n_P} \sum_{i=1}^{n_P} (\phi_{+1}(g, \mathbf{x}_i^P) - \phi_{+1}(g, \mathbf{x}_i^{P'})) + \frac{\pi_N}{n_N} \sum_{i=1}^{n_N} (\phi_{-1}(g, \mathbf{x}_i^N) - \phi_{-1}(g, \mathbf{x}_i^{N'})). \end{aligned} \quad (\text{A9})$$

By the subadditivity of  $\sup$ ,

$$\begin{aligned} & \mathbb{E}_{(\mathcal{X}_P, \mathcal{X}_N), (\mathcal{X}'_P, \mathcal{X}'_N)} \left[ \sup_{g \in \mathcal{G}} \hat{R}_{\text{PN-TR}}(g; \mathcal{X}_P, \mathcal{X}_N) - \hat{R}_{\text{PN-TR}}(g; \mathcal{X}'_P, \mathcal{X}'_N) \right] \\ &\leq \frac{\pi_P}{n_P} \mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (\phi_{+1}(g, \mathbf{x}_i^P) - \phi_{+1}(g, \mathbf{x}_i^{P'})) \right] \\ &\quad + \frac{\pi_N}{n_N} \mathbb{E}_{\mathcal{X}_N, \mathcal{X}'_N} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_N} (\phi_{-1}(g, \mathbf{x}_i^N) - \phi_{-1}(g, \mathbf{x}_i^{N'})) \right]. \end{aligned} \quad (\text{A10})$$

Now consider, for example, the P side. Since  $\mathbf{x}_i^P$  and  $\mathbf{x}_i^{P'}$  are independent and identically distributed (both from  $p_P$ ), the two differences  $\phi_{+1}(g, \mathbf{x}_i^P) - \phi_{+1}(g, \mathbf{x}_i^{P'})$  and  $\phi_{+1}(g, \mathbf{x}_i^{P'}) - \phi_{+1}(g, \mathbf{x}_i^P)$  have the same distribution. Therefore, letting  $L_{2, n_P} := \sum_{i=2}^{n_P} (\phi_{+1}(g, \mathbf{x}_i^P) - \phi_{+1}(g, \mathbf{x}_i^{P'}))$ , we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (\phi_{+1}(g, \mathbf{x}_i^P) - \phi_{+1}(g, \mathbf{x}_i^{P'})) \right] \\ &= \mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} (\phi_{+1}(g, \mathbf{x}_1^P) - \phi_{+1}(g, \mathbf{x}_1^{P'})) + L_{2, n_P} \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} (\phi_{+1}(g, \mathbf{x}_1^P) - \phi_{+1}(g, \mathbf{x}_1^{P'})) + L_{2, n_P} \right] \\
&\quad + \frac{1}{2} \mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} (\phi_{+1}(g, \mathbf{x}_1^{P'}) - \phi_{+1}(g, \mathbf{x}_1^P)) + L_{2, n_P} \right] \\
&= \mathbb{E}_{\sigma_1, \mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sigma_1 (\phi_{+1}(g, \mathbf{x}_1^P) - \phi_{+1}(g, \mathbf{x}_1^{P'})) + L_{2, n_P} \right]. \tag{A11}
\end{aligned}$$

Repeating the same argument  $n_P$  times yields

$$\begin{aligned}
&\mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (\phi_{+1}(g, \mathbf{x}_i^P) - \phi_{+1}(g, \mathbf{x}_i^{P'})) \right] \\
&= \mathbb{E}_{\sigma, \mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} \sigma_i (\phi_{+1}(g, \mathbf{x}_i^P) - \phi_{+1}(g, \mathbf{x}_i^{P'})) \right] \\
&\leq \mathbb{E}_{\sigma, \mathcal{X}_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} \sigma_i \phi_{+1}(g, \mathbf{x}_i^P) \right] + \mathbb{E}_{\sigma, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (-\sigma_i) \phi_{+1}(g, \mathbf{x}_i^{P'}) \right] \\
&= 2 \mathbb{E}_{\sigma, \mathcal{X}_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} \sigma_i \phi_{+1}(g, \mathbf{x}_i^P) \right] = 2n_P \mathfrak{R}_{n_P, p_P}(\phi_{+1} \circ \mathcal{G}). \tag{A12}
\end{aligned}$$

Similarly, for the N side,

$$\mathbb{E}_{\mathcal{X}_N, \mathcal{X}'_N} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_N} (\phi_{-1}(g, \mathbf{x}_i^N) - \phi_{-1}(g, \mathbf{x}_i^{N'})) \right] \leq 2n_N \mathfrak{R}_{n_N, p_N}(\phi_{-1} \circ \mathcal{G})$$

is obtained.

Finally, we upper bound  $\mathfrak{R}(\phi_y \circ \mathcal{G})$  by  $\mathfrak{R}(\mathcal{G})$ . Write  $\phi_y = \ell(\cdot, y) + \beta \psi$ , where

$$\psi(g, \mathbf{x}) := \max_{\|\boldsymbol{\eta}\|_\infty \leq \varepsilon} \ell_{\text{KL}}(g(\mathbf{x}), g(\mathbf{x} + \boldsymbol{\eta})).$$

Then, by the subadditivity of sup,

$$\mathfrak{R}_{n, q}(\phi_y \circ \mathcal{G}) \leq \mathfrak{R}_{n, q}(\ell(\cdot, y) \circ \mathcal{G}) + \beta \mathfrak{R}_{n, q}(\psi \circ \mathcal{G}).$$

For the classification loss term, by the contraction lemma,

$$\mathfrak{R}_{n, q}(\ell(\cdot, y) \circ \mathcal{G}) \leq L_\ell \mathfrak{R}_{n, q}(\mathcal{G}).$$

For the KL term, applying Lemma A.2 to  $f_g(\mathbf{x}) = g(\mathbf{x})$  and  $g_{g, \boldsymbol{\eta}}(\mathbf{x}) = g(\mathbf{x} + \boldsymbol{\eta})$ , and then using Lemma A.4, we obtain

$$\begin{aligned}
\mathfrak{R}_{n, q}(\psi \circ \mathcal{G}) &\leq 2L_{\text{KL}} \left( \mathfrak{R}_{n, q}(\mathcal{G}) + \mathfrak{R}_{n, q}(\{ \mathbf{x} \mapsto g(\mathbf{x} + \boldsymbol{\eta}) : \|\boldsymbol{\eta}\|_\infty \leq \varepsilon, g \in \mathcal{G} \}) \right) \\
&\leq 2L_{\text{KL}} \left( 2\mathfrak{R}_{n, q}(\mathcal{G}) + \frac{\varepsilon W d^{1/q}}{\sqrt{n}} \right). \tag{A13}
\end{aligned}$$

Substituting these bounds into (A8)–(A10), we get

$$\begin{aligned}
&\mathbb{E} \left[ \sup_{g \in \mathcal{G}} \hat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g) \right] \\
&\leq \frac{\pi_P}{n_P} \cdot 2n_P \mathfrak{R}_{n_P, p_P}(\phi_{+1} \circ \mathcal{G}) + \frac{\pi_N}{n_N} \cdot 2n_N \mathfrak{R}_{n_N, p_N}(\phi_{-1} \circ \mathcal{G}) \\
&= 2\pi_P \mathfrak{R}_{n_P, p_P}(\phi_{+1} \circ \mathcal{G}) + 2\pi_N \mathfrak{R}_{n_N, p_N}(\phi_{-1} \circ \mathcal{G}) \\
&\leq 2\pi_P \left\{ L_\ell \mathfrak{R}_{n_P, p_P}(\mathcal{G}) + \beta \mathfrak{R}_{n_P, p_P}(\psi \circ \mathcal{G}) \right\} + 2\pi_N \left\{ L_\ell \mathfrak{R}_{n_N, p_N}(\mathcal{G}) + \beta \mathfrak{R}_{n_N, p_N}(\psi \circ \mathcal{G}) \right\}
\end{aligned}$$

$$\begin{aligned}
&\leq 2\pi_P \left\{ L_\ell \mathfrak{R}_{n_P, p_P}(\mathcal{G}) + \beta \cdot 2L_{\text{KL}} \left( 2\mathfrak{R}_{n_P, p_P}(\mathcal{G}) + \frac{\varepsilon W d^{1/q}}{\sqrt{n_P}} \right) \right\} \\
&\quad + 2\pi_N \left\{ L_\ell \mathfrak{R}_{n_N, p_N}(\mathcal{G}) + \beta \cdot 2L_{\text{KL}} \left( 2\mathfrak{R}_{n_N, p_N}(\mathcal{G}) + \frac{\varepsilon W d^{1/q}}{\sqrt{n_N}} \right) \right\} \\
&= 2(L_\ell + 4\beta L_{\text{KL}}) \left( \pi_P \mathfrak{R}_{n_P, p_P}(\mathcal{G}) + \pi_N \mathfrak{R}_{n_N, p_N}(\mathcal{G}) \right) + 4\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_P}{\sqrt{n_P}} + \frac{\pi_N}{\sqrt{n_N}} \right),
\end{aligned} \tag{A14}$$

where  $\phi_y(g, \mathbf{x}) := \ell(g(\mathbf{x}), y) + \beta \max_{\|\boldsymbol{\eta}\|_\infty \leq \varepsilon} \ell_{\text{KL}}(g(\mathbf{x}), g(\mathbf{x} + \boldsymbol{\eta}))$  and  $\psi(g, \mathbf{x}) := \max_{\|\boldsymbol{\eta}\|_\infty \leq \varepsilon} \ell_{\text{KL}}(g(\mathbf{x}), g(\mathbf{x} + \boldsymbol{\eta}))$ .

**(iii) Combining the results of (i) and (ii)**

Substituting the expectation bound in (ii), i.e., (A14), into McDiarmid's inequality result in (i), i.e., (A7), we obtain, with probability at least  $1 - \delta/2$ ,

$$\begin{aligned}
\sup_{g \in \mathcal{G}} \left\{ \widehat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g) \right\} &\leq 2(L_\ell + 4\beta L_{\text{KL}}) \left( \pi_P \mathfrak{R}_{n_P, p_P}(\mathcal{G}) + \pi_N \mathfrak{R}_{n_N, p_N}(\mathcal{G}) \right) \\
&\quad + 4\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_P}{\sqrt{n_P}} + \frac{\pi_N}{\sqrt{n_N}} \right) \\
&\quad + \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} (C_\ell + \beta C_{\text{KL}}) \left( \frac{\pi_P}{\sqrt{n_P}} + \frac{\pi_N}{\sqrt{n_N}} \right)
\end{aligned} \tag{A15}$$

as claimed.

**(iv) Reverse direction and union bound**

By the same argument,  $\sup_{g \in \mathcal{G}} \{R_{\text{PN-TR}}(g) - \widehat{R}_{\text{PN-TR}}(g)\}$  is also bounded by the same right-hand side with probability at least  $1 - \delta/2$ . Therefore, by the union bound, (23) holds with probability at least  $1 - \delta$ .  $\square$

### A.3 Proof of Theorem 5.2

*Proof* Since  $\widehat{g}_{\text{PN-TR}}$  is an empirical risk minimizer, we have  $\widehat{R}_{\text{PN-TR}}(\widehat{g}_{\text{PN-TR}}) \leq \widehat{R}_{\text{PN-TR}}(g^*)$ . Moreover, by Lemma 5.3, the following holds with probability at least  $1 - \delta$ :

$$\begin{aligned}
&R_{\text{PN-TR}}(\widehat{g}_{\text{PN-TR}}) - R_{\text{PN-TR}}(g^*) \\
&= \left( R_{\text{PN-TR}}(\widehat{g}_{\text{PN-TR}}) - \widehat{R}_{\text{PN-TR}}(\widehat{g}_{\text{PN-TR}}) \right) + \left( \widehat{R}_{\text{PN-TR}}(\widehat{g}_{\text{PN-TR}}) - \widehat{R}_{\text{PN-TR}}(g^*) \right) \\
&\quad + \left( \widehat{R}_{\text{PN-TR}}(g^*) - R_{\text{PN-TR}}(g^*) \right) \\
&\leq \sup_{g \in \mathcal{G}} (R_{\text{PN-TR}}(g) - \widehat{R}_{\text{PN-TR}}(g)) + 0 + \sup_{g \in \mathcal{G}} (\widehat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g)) \\
&\leq 2 \sup_{g \in \mathcal{G}} \left| \widehat{R}_{\text{PN-TR}}(g) - R_{\text{PN-TR}}(g) \right| \\
&\leq 4(L_\ell + 4\beta L_{\text{KL}}) \left( \pi_P \mathfrak{R}_{n_P, p_P}(\mathcal{G}) + \pi_N \mathfrak{R}_{n_N, p_N}(\mathcal{G}) \right) \\
&\quad + 8\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_P}{\sqrt{n_P}} + \frac{\pi_N}{\sqrt{n_N}} \right)
\end{aligned}$$

$$+ \sqrt{2 \ln \frac{2}{\delta}} (C_\ell + \beta C_{\text{KL}}) \left( \frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{\pi_{\text{N}}}{\sqrt{n_{\text{N}}}} \right).$$

Thus, the claim of Theorem 5.2 follows.  $\square$

#### A.4 Proof of the Lemma “Uniform Deviation Bound for uPU+TRADES” (Lemma 5.5)

*Proof* As in the previous section, (i) we first show concentration of the uniform deviation using McDiarmid’s inequality, (ii) evaluate the expectation via ghost sampling and Rademacher complexity, (iii) combine the results of (i) and (ii) to obtain a one-sided upper bound, and (iv) combine the reverse direction via a union bound to obtain the uniform deviation bound (the lemma). Finally, in (v), we derive the estimation error bound (the theorem) by the standard decomposition for an empirical risk minimizer.

##### (i) McDiarmid’s inequality

Since  $\ell(\cdot) \leq C_\ell$  and  $\ell_{\text{KL}}(\cdot) \leq C_{\text{KL}}$ , replacing one sample on the P side changes  $\hat{R}_{\text{uPU-TR}}(g)$  by at most  $\frac{\pi_{\text{P}}}{n_{\text{P}}} (2C_\ell + \beta C_{\text{KL}})$ . Similarly, replacing one sample on the U side changes  $\hat{R}_{\text{uPU-TR}}(g)$  by at most  $\frac{1}{n_{\text{U}}} (C_\ell + \beta C_{\text{KL}})$ . Therefore, by McDiarmid’s inequality, for any  $t > 0$ ,

$$\begin{aligned} & \Pr \left( \sup_{g \in \mathcal{G}} \left\{ \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \right\} - \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \left\{ \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \right\} \right] \geq t \right) \\ & \leq \exp \left( - \frac{2t^2}{\frac{\pi_{\text{P}}^2 (2C_\ell + \beta C_{\text{KL}})^2}{n_{\text{P}}} + \frac{(C_\ell + \beta C_{\text{KL}})^2}{n_{\text{U}}}} \right). \end{aligned} \quad (\text{A16})$$

Setting the right-hand side equal to  $\delta/2$ , we obtain, with probability at least  $1 - \delta/2$ ,

$$\begin{aligned} \sup_{g \in \mathcal{G}} \{ \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \} & \leq \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \{ \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \} \right] \\ & + \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} \sqrt{\frac{\pi_{\text{P}}^2 (2C_\ell + \beta C_{\text{KL}})^2}{n_{\text{P}}} + \frac{(C_\ell + \beta C_{\text{KL}})^2}{n_{\text{U}}}}. \end{aligned} \quad (\text{A17})$$

Furthermore, by using the subadditivity of the square root, we obtain

$$\begin{aligned} \sup_{g \in \mathcal{G}} \{ \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \} & \leq \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \{ \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \} \right] \\ & + \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} \left( \frac{\pi_{\text{P}} (2C_\ell + \beta C_{\text{KL}})}{\sqrt{n_{\text{P}}}} + \frac{C_\ell + \beta C_{\text{KL}}}{\sqrt{n_{\text{U}}}} \right) \end{aligned} \quad (\text{A18})$$

as desired.

##### (ii) Ghost sampling

Here, we upper bound  $\mathbb{E} \left[ \sup_{g \in \mathcal{G}} \{ \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \} \right]$  by the Rademacher complexity. The expectation in  $\mathbb{E} [\sup_{g \in \mathcal{G}} \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g)]$  is taken over repeated sampling of  $(\mathcal{X}_{\text{P}}, \mathcal{X}_{\text{U}})$  used to evaluate  $\hat{R}_{\text{uPU-TR}}(g)$ . That is,  $\mathbb{E} = \mathbb{E}_{(\mathcal{X}_{\text{P}}, \mathcal{X}_{\text{U}})}$ , and  $\hat{R}_{\text{uPU-TR}}(g)$  can also

be written as  $\hat{R}_{\text{uPU-TR}}(g; \mathcal{X}_P, \mathcal{X}_U)$ . Introducing an independent ghost sample  $(\mathcal{X}'_P, \mathcal{X}'_U)$ , we obtain

$$\begin{aligned} & \mathbb{E}_{(\mathcal{X}_P, \mathcal{X}_U)} \left[ \sup_{g \in \mathcal{G}} \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \right] \\ &= \mathbb{E}_{(\mathcal{X}_P, \mathcal{X}_U)} \left[ \sup_{g \in \mathcal{G}} \hat{R}_{\text{uPU-TR}}(g) - \mathbb{E}_{(\mathcal{X}'_P, \mathcal{X}'_U)} \hat{R}_{\text{uPU-TR}}(g; \mathcal{X}'_P, \mathcal{X}'_U) \right] \\ &\leq \mathbb{E}_{(\mathcal{X}_P, \mathcal{X}_U), (\mathcal{X}'_P, \mathcal{X}'_U)} \left[ \sup_{g \in \mathcal{G}} \hat{R}_{\text{uPU-TR}}(g; \mathcal{X}_P, \mathcal{X}_U) - \hat{R}_{\text{uPU-TR}}(g; \mathcal{X}'_P, \mathcal{X}'_U) \right] \end{aligned} \quad (\text{A19})$$

The last inequality follows from Jensen's inequality, since sup is convex.

Next, we decompose the difference:

$$\phi_P(g, \mathbf{x}) := \ell(g(\mathbf{x}), +1) - \ell(g(\mathbf{x}), -1) + \beta\psi(g, \mathbf{x}), \quad \phi_U(g, \mathbf{x}) := \ell(g(\mathbf{x}), -1) + \beta\psi(g, \mathbf{x}),$$

where

$$\psi(g, \mathbf{x}) := \max_{\|\boldsymbol{\eta}\|_\infty \leq \varepsilon} \ell_{\text{KL}}(g(\mathbf{x} + \boldsymbol{\eta}), g(\mathbf{x})).$$

Then,

$$\begin{aligned} & \hat{R}_{\text{uPU-TR}}(g; \mathcal{X}_P, \mathcal{X}_U) - \hat{R}_{\text{uPU-TR}}(g; \mathcal{X}'_P, \mathcal{X}'_U) \\ &= \frac{\pi_P}{n_P} \sum_{i=1}^{n_P} (\phi_P(g, \mathbf{x}_i^P) - \phi_P(g, \mathbf{x}_i^{P'})) + \frac{1}{n_U} \sum_{i=1}^{n_U} (\phi_U(g, \mathbf{x}_i^U) - \phi_U(g, \mathbf{x}_i^{U'})). \end{aligned} \quad (\text{A20})$$

Therefore, by the subadditivity of sup,

$$\begin{aligned} & \mathbb{E}_{(\mathcal{X}_P, \mathcal{X}_U), (\mathcal{X}'_P, \mathcal{X}'_U)} \left[ \sup_{g \in \mathcal{G}} \hat{R}_{\text{uPU-TR}}(g; \mathcal{X}_P, \mathcal{X}_U) - \hat{R}_{\text{uPU-TR}}(g; \mathcal{X}'_P, \mathcal{X}'_U) \right] \\ &\leq \frac{\pi_P}{n_P} \mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (\phi_P(g, \mathbf{x}_i^P) - \phi_P(g, \mathbf{x}_i^{P'})) \right] \\ &\quad + \frac{1}{n_U} \mathbb{E}_{\mathcal{X}_U, \mathcal{X}'_U} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_U} (\phi_U(g, \mathbf{x}_i^U) - \phi_U(g, \mathbf{x}_i^{U'})) \right]. \end{aligned} \quad (\text{A21})$$

Next, we symmetrize the difference terms appearing on the right-hand side of (A21) and bound them by the Rademacher complexity. For example, consider the P side. Since  $\mathbf{x}_i^P$  and  $\mathbf{x}_i^{P'}$  are independent and identically distributed (both from  $p_P$ ), the two differences  $\phi_P(g, \mathbf{x}_i^P) - \phi_P(g, \mathbf{x}_i^{P'})$  and  $\phi_P(g, \mathbf{x}_i^{P'}) - \phi_P(g, \mathbf{x}_i^P)$  have the same distribution. Therefore, letting  $L_{2, n_P} := \sum_{i=2}^{n_P} (\phi_P(g, \mathbf{x}_i^P) - \phi_P(g, \mathbf{x}_i^{P'}))$ , we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (\phi_P(g, \mathbf{x}_i^P) - \phi_P(g, \mathbf{x}_i^{P'})) \right] \\ &= \frac{1}{2} \mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} (\phi_P(g, \mathbf{x}_1^P) - \phi_P(g, \mathbf{x}_1^{P'})) + L_{2, n_P} \right] \\ &\quad + \frac{1}{2} \mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} (\phi_P(g, \mathbf{x}_1^{P'}) - \phi_P(g, \mathbf{x}_1^P)) + L_{2, n_P} \right] \\ &= \mathbb{E}_{\sigma_1, \mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sigma_1 (\phi_P(g, \mathbf{x}_1^P) - \phi_P(g, \mathbf{x}_1^{P'})) + L_{2, n_P} \right]. \end{aligned} \quad (\text{A22})$$

Repeating the same argument  $n_P$  times, and using independent Rademacher variables  $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_{n_P})$ , we obtain

$$\mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (\phi_P(g, \mathbf{x}_i^P) - \phi_P(g, \mathbf{x}_i^{P'})) \right]$$

$$\begin{aligned}
&= \mathbb{E}_{\boldsymbol{\sigma}, \mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} \sigma_i (\phi_P(g, \mathbf{x}_i^P) - \phi_P(g, \mathbf{x}_i^{P'})) \right] \\
&\leq \mathbb{E}_{\boldsymbol{\sigma}, \mathcal{X}_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} \sigma_i \phi_P(g, \mathbf{x}_i^P) \right] + \mathbb{E}_{\boldsymbol{\sigma}, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (-\sigma_i) \phi_P(g, \mathbf{x}_i^{P'}) \right] \\
&= 2 \mathbb{E}_{\boldsymbol{\sigma}, \mathcal{X}_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} \sigma_i \phi_P(g, \mathbf{x}_i^P) \right] = 2n_P \mathfrak{R}_{n_P, P_P}(\phi_P \circ \mathcal{G}). \tag{A23}
\end{aligned}$$

Similarly, for the U side,

$$\mathbb{E}_{\mathcal{X}_U, \mathcal{X}'_U} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_U} (\phi_U(g, \mathbf{x}_i^U) - \phi_U(g, \mathbf{x}_i^{U'})) \right] \leq 2n_U \mathfrak{R}_{n_U, P_U}(\phi_U \circ \mathcal{G}) \tag{A24}$$

is obtained.

Finally, we upper bound  $\mathfrak{R}(\phi_P \circ \mathcal{G})$  and  $\mathfrak{R}(\phi_U \circ \mathcal{G})$  by  $\mathfrak{R}(\mathcal{G})$ . By the subadditivity of sup,

$$\begin{aligned}
\mathfrak{R}_{n,q}(\phi_P \circ \mathcal{G}) &\leq \mathfrak{R}_{n,q}(\ell(\cdot, +1) \circ \mathcal{G}) + \mathfrak{R}_{n,q}(\ell(\cdot, -1) \circ \mathcal{G}), \\
\mathfrak{R}_{n,q}(\phi_U \circ \mathcal{G}) &\leq \mathfrak{R}_{n,q}(\ell(\cdot, -1) \circ \mathcal{G}) + \beta \mathfrak{R}_{n,q}(\psi \circ \mathcal{G}).
\end{aligned}$$

For the classification loss term, the contraction lemma gives

$$\mathfrak{R}_{n,q}(\ell(\cdot, y) \circ \mathcal{G}) \leq L_\ell \mathfrak{R}_{n,q}(\mathcal{G}) \quad (y \in \mathcal{Y}).$$

For the KL term, applying Lemma A.2 to  $f_g(\mathbf{x}) = g(\mathbf{x})$  and  $g_{g,\boldsymbol{\eta}}(\mathbf{x}) = g(\mathbf{x} + \boldsymbol{\eta})$ , and then using Lemma A.4, we obtain

$$\begin{aligned}
\mathfrak{R}_{n,q}(\psi \circ \mathcal{G}) &\leq 2L_{\text{KL}} \left( \mathfrak{R}_{n,q}(\mathcal{G}) + \mathfrak{R}_{n,q}(\{ \mathbf{x} \mapsto g(\mathbf{x} + \boldsymbol{\eta}) : \|\boldsymbol{\eta}\|_\infty \leq \varepsilon, g \in \mathcal{G} \}) \right) \\
&\leq 2L_{\text{KL}} \left( 2\mathfrak{R}_{n,q}(\mathcal{G}) + \frac{\varepsilon W d^{1/q}}{\sqrt{n}} \right). \tag{A25}
\end{aligned}$$

Substituting these bounds into (A19)–(A21), we obtain

$$\begin{aligned}
\mathbb{E} \left[ \sup_{g \in \mathcal{G}} \{ \widehat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \} \right] &\leq 2\pi_P \mathfrak{R}_{n_P, P_P}(\phi_P \circ \mathcal{G}) + 2\mathfrak{R}_{n_U, P_U}(\phi_U \circ \mathcal{G}) \\
&\leq 4\pi_P (L_\ell + 2\beta L_{\text{KL}}) \mathfrak{R}_{n_P, P_P}(\mathcal{G}) \\
&\quad + 2(L_\ell + 4\beta L_{\text{KL}}) \mathfrak{R}_{n_U, P_U}(\mathcal{G}) \\
&\quad + 4\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_P}{\sqrt{n_P}} + \frac{1}{\sqrt{n_U}} \right). \tag{A26}
\end{aligned}$$

This completes the bound.

### (iii) Combining the results of (i) and (ii)

Substituting the expectation bound in (ii), i.e., (A26), into the McDiarmid inequality result in (i), i.e., (A18), we obtain, with probability at least  $1 - \delta/2$ ,

$$\begin{aligned}
\sup_{g \in \mathcal{G}} \{ \widehat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \} &\leq 4\pi_P (L_\ell + 2\beta L_{\text{KL}}) \mathfrak{R}_{n_P, P_P}(\mathcal{G}) + 2(L_\ell + 4\beta L_{\text{KL}}) \mathfrak{R}_{n_U, P_U}(\mathcal{G}) \\
&\quad + 4\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_P}{\sqrt{n_P}} + \frac{1}{\sqrt{n_U}} \right) \\
&\quad + \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} \left( \frac{\pi_P(2C_\ell + \beta C_{\text{KL}})}{\sqrt{n_P}} + \frac{C_\ell + \beta C_{\text{KL}}}{\sqrt{n_U}} \right). \tag{A27}
\end{aligned}$$

This completes the one-sided bound.  $\square$

## A.5 Proof of Theorem 5.4

*Proof* Since  $\hat{g}_{\text{uPU-TR}}$  is an empirical risk minimizer, we have

$$\hat{R}_{\text{uPU-TR}}(\hat{g}_{\text{uPU-TR}}) \leq \hat{R}_{\text{uPU-TR}}(g^*)$$

Moreover, on the event where Lemma 5.5 holds (which occurs with probability at least  $1 - \delta$ ),

$$\begin{aligned} & R_{\text{uPU-TR}}(\hat{g}_{\text{uPU-TR}}) - R_{\text{uPU-TR}}(g^*) \\ &= \left( R_{\text{uPU-TR}}(\hat{g}_{\text{uPU-TR}}) - \hat{R}_{\text{uPU-TR}}(\hat{g}_{\text{uPU-TR}}) \right) \\ &+ \left( \hat{R}_{\text{uPU-TR}}(\hat{g}_{\text{uPU-TR}}) - \hat{R}_{\text{uPU-TR}}(g^*) \right) \\ &+ \left( \hat{R}_{\text{uPU-TR}}(g^*) - R_{\text{uPU-TR}}(g^*) \right) \\ &\leq \sup_{g \in \mathcal{G}} (R_{\text{uPU-TR}}(g) - \hat{R}_{\text{uPU-TR}}(g)) + 0 + \sup_{g \in \mathcal{G}} (\hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g)) \\ &\leq 2 \sup_{g \in \mathcal{G}} \left| \hat{R}_{\text{uPU-TR}}(g) - R_{\text{uPU-TR}}(g) \right|. \end{aligned}$$

Therefore, multiplying (27) by 2 yields (26).  $\square$

## A.6 Proof of the Lemma “Uniform Deviation Bound for nnPU+TRADES” (Lemma 5.7)

*Proof*

### (i) McDiarmid’s inequality

Since  $\ell(\cdot) \leq C_\ell$  and  $\ell_{\text{KL}}(\cdot) \leq C_{\text{KL}}$ , replacing one sample on the P side changes the sum in the first line by at most  $\frac{\pi_P}{n_P}(C_\ell + \beta C_{\text{KL}})$ . In addition, since the truncation term  $\max\{0, \cdot\}$  in nnPU is 1-Lipschitz, it can be upper bounded by the change in its argument. Inside the truncation term,  $\ell(g(\mathbf{x}), -1)$  on the P side is replaced at only one point, so the change is at most  $\frac{\pi_P}{n_P}C_\ell$ . Therefore, replacing one P-side sample changes  $\hat{R}_{\text{nnPU-TR}}(g)$  by at most  $\frac{\pi_P}{n_P}(2C_\ell + \beta C_{\text{KL}})$  in total. Similarly, replacing one U-side sample affects  $\ell(g(\mathbf{x}), -1)$  inside the truncation term and  $\psi$  in the third line at only one point, and hence changes  $\hat{R}_{\text{nnPU-TR}}(g)$  by at most  $\frac{1}{n_U}(C_\ell + \beta C_{\text{KL}})$ . Therefore, by McDiarmid’s inequality, for any  $t > 0$ ,

$$\begin{aligned} & \Pr \left( \sup_{g \in \mathcal{G}} \left\{ \hat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \right\} - \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \left\{ \hat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \right\} \right] \geq t \right) \\ & \leq \exp \left( - \frac{2t^2}{\frac{\pi_P^2(2C_\ell + \beta C_{\text{KL}})^2}{n_P} + \frac{(C_\ell + \beta C_{\text{KL}})^2}{n_U}} \right). \end{aligned} \quad (\text{A28})$$

Setting the right-hand side equal to  $\delta/2$ , we obtain, with probability at least  $1 - \delta/2$ ,

$$\begin{aligned} \sup_{g \in \mathcal{G}} \{ \hat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \} &\leq \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \{ \hat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \} \right] \\ &+ \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} \sqrt{\frac{\pi_P^2(2C_\ell + \beta C_{\text{KL}})^2}{n_P} + \frac{(C_\ell + \beta C_{\text{KL}})^2}{n_U}}. \end{aligned} \quad (\text{A29})$$

Furthermore, by using the subadditivity of the square root, we obtain

$$\sup_{g \in \mathcal{G}} \{ \hat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \} \leq \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \{ \hat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \} \right]$$

$$+ \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} \left( \frac{\pi_P(2C_\ell + \beta C_{\text{KL}})}{\sqrt{n_P}} + \frac{C_\ell + \beta C_{\text{KL}}}{\sqrt{n_U}} \right) \quad (\text{A30})$$

as desired.

**(ii) Ghost sampling**

Here, we upper bound  $\mathbb{E}[\sup_{g \in \mathcal{G}} \{\widehat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g)\}]$  by the Rademacher complexity. Unlike the uPU case,  $\widehat{R}_{\text{nnPU-TR}}(g)$  contains the nnPU truncation term  $\max\{0, \cdot\}$ , and therefore, in general,  $R_{\text{nnPU-TR}}(g) = \mathbb{E}[\widehat{R}_{\text{nnPU-TR}}(g)]$  does not hold. Thus, using the 1-Lipschitz property of  $\max\{0, \cdot\}$ ,

$$|\max\{0, a\} - \max\{0, b\}| \leq |a - b|,$$

we first reduce the problem to a sum of differences between empirical and population averages, and then evaluate each term by ghost sampling and symmetrization, as in uPU+TRADES.

First, define

$$\phi_+(g, \mathbf{x}) := \ell(g(\mathbf{x}), +1) + \beta \psi(g, \mathbf{x}), \quad \phi_-(g, \mathbf{x}) := \ell(g(\mathbf{x}), -1), \quad \phi_\psi(g, \mathbf{x}) := \psi(g, \mathbf{x}),$$

where

$$\psi(g, \mathbf{x}) := \max_{\|\boldsymbol{\eta}\|_\infty \leq \varepsilon} \ell_{\text{KL}}(g(\mathbf{x} + \boldsymbol{\eta}), g(\mathbf{x})).$$

Then, writing the inside of the truncation term as

$$\begin{aligned} \widehat{s}(g) &:= -\frac{\pi_P}{n_P} \sum_{i=1}^{n_P} \phi_-(g, \mathbf{x}_i^P) + \frac{1}{n_U} \sum_{i=1}^{n_U} \phi_-(g, \mathbf{x}_i^U), \\ s(g) &:= -\pi_P \mathbb{E}_P[\phi_-(g, \mathbf{x})] + \mathbb{E}_U[\phi_-(g, \mathbf{x})], \end{aligned}$$

we have

$$\begin{aligned} \widehat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) &= \pi_P \left\{ \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_+(g, \mathbf{x}_i^P) - \mathbb{E}_P[\phi_+(g, \mathbf{x})] \right\} \\ &\quad + \beta \left\{ \frac{1}{n_U} \sum_{i=1}^{n_U} \phi_\psi(g, \mathbf{x}_i^U) - \mathbb{E}_U[\phi_\psi(g, \mathbf{x})] \right\} \\ &\quad + \max\{0, \widehat{s}(g)\} - \max\{0, s(g)\} \\ &\leq \pi_P \left\{ \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_+(g, \mathbf{x}_i^P) - \mathbb{E}_P[\phi_+(g, \mathbf{x})] \right\} \\ &\quad + \beta \left\{ \frac{1}{n_U} \sum_{i=1}^{n_U} \phi_\psi(g, \mathbf{x}_i^U) - \mathbb{E}_U[\phi_\psi(g, \mathbf{x})] \right\} \\ &\quad + |\widehat{s}(g) - s(g)| \\ &\leq \pi_P \left\{ \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_+(g, \mathbf{x}_i^P) - \mathbb{E}_P[\phi_+(g, \mathbf{x})] \right\} \\ &\quad + \beta \left\{ \frac{1}{n_U} \sum_{i=1}^{n_U} \phi_\psi(g, \mathbf{x}_i^U) - \mathbb{E}_U[\phi_\psi(g, \mathbf{x})] \right\} \\ &\quad + \pi_P \left| \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_-(g, \mathbf{x}_i^P) - \mathbb{E}_P[\phi_-(g, \mathbf{x})] \right| \end{aligned}$$



$$+ \left| \frac{1}{n_U} \sum_{i=1}^{n_U} \phi_{-}(g, \mathbf{x}_i^U) - \mathbb{E}_U[\phi_{-}(g, \mathbf{x})] \right|. \quad (\text{A31})$$

Therefore, by the subadditivity of sup and the triangle inequality,

$$\begin{aligned} \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \{ \widehat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \} \right] &\leq \pi_P \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_{+}(g, \mathbf{x}_i^P) - \mathbb{E}_P[\phi_{+}(g, \mathbf{x})] \right\} \right] \\ &\quad + \beta \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n_U} \sum_{i=1}^{n_U} \phi_{\psi}(g, \mathbf{x}_i^U) - \mathbb{E}_U[\phi_{\psi}(g, \mathbf{x})] \right\} \right] \\ &\quad + \pi_P \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \left| \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_{-}(g, \mathbf{x}_i^P) - \mathbb{E}_P[\phi_{-}(g, \mathbf{x})] \right| \right] \\ &\quad + \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \left| \frac{1}{n_U} \sum_{i=1}^{n_U} \phi_{-}(g, \mathbf{x}_i^U) - \mathbb{E}_U[\phi_{-}(g, \mathbf{x})] \right| \right]. \end{aligned} \quad (\text{A32})$$

Next, we evaluate each term by ghost sampling and symmetrization (as in uPU+TRADES). For illustration, consider the first term on the P side. Let  $\mathcal{X}'_P = \{\mathbf{x}_i^{P'}\}_{i=1}^{n_P}$  be an independent ghost sample on the P side. Then,

$$\begin{aligned} &\mathbb{E}_{\mathcal{X}_P} \left[ \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_{+}(g, \mathbf{x}_i^P) - \mathbb{E}_P[\phi_{+}(g, \mathbf{x})] \right\} \right] \\ &= \mathbb{E}_{\mathcal{X}_P} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_{+}(g, \mathbf{x}_i^P) - \mathbb{E}_{\mathcal{X}'_P} \left\{ \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_{+}(g, \mathbf{x}_i^{P'}) \right\} \right] \\ &\leq \mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n_P} \sum_{i=1}^{n_P} \left( \phi_{+}(g, \mathbf{x}_i^P) - \phi_{+}(g, \mathbf{x}_i^{P'}) \right) \right] \end{aligned} \quad (\text{A33})$$

where the last inequality follows from Jensen's inequality, since sup is convex.

Here, since  $\mathbf{x}_i^P$  and  $\mathbf{x}_i^{P'}$  are independent and identically distributed (both from  $p_P$ ), the two differences  $\phi_{+}(g, \mathbf{x}_i^P) - \phi_{+}(g, \mathbf{x}_i^{P'})$  and  $\phi_{+}(g, \mathbf{x}_i^{P'}) - \phi_{+}(g, \mathbf{x}_i^P)$  have the same distribution. Therefore, letting  $L_{2, n_P} := \sum_{i=2}^{n_P} (\phi_{+}(g, \mathbf{x}_i^P) - \phi_{+}(g, \mathbf{x}_i^{P'}))$ , we have

$$\begin{aligned} &\mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (\phi_{+}(g, \mathbf{x}_i^P) - \phi_{+}(g, \mathbf{x}_i^{P'})) \right] \\ &= \mathbb{E}_{\sigma_1, \mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sigma_1 (\phi_{+}(g, \mathbf{x}_1^P) - \phi_{+}(g, \mathbf{x}_1^{P'})) + L_{2, n_P} \right]. \end{aligned} \quad (\text{A34})$$

Repeating the same argument  $n_P$  times, and using independent Rademacher variables  $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_{n_P})$ , we obtain

$$\begin{aligned} &\mathbb{E}_{\mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (\phi_{+}(g, \mathbf{x}_i^P) - \phi_{+}(g, \mathbf{x}_i^{P'})) \right] \\ &= \mathbb{E}_{\boldsymbol{\sigma}, \mathcal{X}_P, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} \sigma_i (\phi_{+}(g, \mathbf{x}_i^P) - \phi_{+}(g, \mathbf{x}_i^{P'})) \right] \\ &\leq \mathbb{E}_{\boldsymbol{\sigma}, \mathcal{X}_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} \sigma_i \phi_{+}(g, \mathbf{x}_i^P) \right] + \mathbb{E}_{\boldsymbol{\sigma}, \mathcal{X}'_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} (-\sigma_i) \phi_{+}(g, \mathbf{x}_i^{P'}) \right] \end{aligned}$$

$$= 2 \mathbb{E}_{\boldsymbol{\sigma}, \mathcal{X}_P} \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n_P} \sigma_i \phi_+(g, \mathbf{x}_i^P) \right] = 2n_P \mathfrak{R}_{n_P, p_P}(\phi_+ \circ \mathcal{G}). \quad (\text{A35})$$

Hence, combining this with (A33), we obtain

$$\mathbb{E}_{\mathcal{X}_P} \left[ \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_+(g, \mathbf{x}_i^P) - \mathbb{E}_P[\phi_+(g, \mathbf{x})] \right\} \right] \leq 2 \mathfrak{R}_{n_P, p_P}(\phi_+ \circ \mathcal{G}).$$

Similarly, the second term on the U side and the third and fourth terms with absolute values can also be symmetrized, yielding

$$\begin{aligned} \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n_U} \sum_{i=1}^{n_U} \phi_\psi(g, \mathbf{x}_i^U) - \mathbb{E}_U[\phi_\psi(g, \mathbf{x})] \right\} \right] &\leq 2 \mathfrak{R}_{n_U, p_U}(\phi_\psi \circ \mathcal{G}), \\ \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \left| \frac{1}{n_P} \sum_{i=1}^{n_P} \phi_-(g, \mathbf{x}_i^P) - \mathbb{E}_P[\phi_-(g, \mathbf{x})] \right| \right] &\leq 2 \mathfrak{R}_{n_P, p_P}(\phi_- \circ \mathcal{G}), \\ \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \left| \frac{1}{n_U} \sum_{i=1}^{n_U} \phi_-(g, \mathbf{x}_i^U) - \mathbb{E}_U[\phi_-(g, \mathbf{x})] \right| \right] &\leq 2 \mathfrak{R}_{n_U, p_U}(\phi_- \circ \mathcal{G}) \end{aligned} \quad (\text{A36})$$

Substituting these bounds into (A32), we obtain

$$\begin{aligned} \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \{ \widehat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \} \right] &\leq 2\pi_P \mathfrak{R}_{n_P, p_P}(\phi_+ \circ \mathcal{G}) \\ &\quad + 2\beta \mathfrak{R}_{n_U, p_U}(\phi_\psi \circ \mathcal{G}) \\ &\quad + 2\pi_P \mathfrak{R}_{n_P, p_P}(\phi_- \circ \mathcal{G}) \\ &\quad + 2 \mathfrak{R}_{n_U, p_U}(\phi_- \circ \mathcal{G}). \end{aligned} \quad (\text{A37})$$

Finally, we upper bound each Rademacher complexity term by  $\mathfrak{R}(\mathcal{G})$ . By the subadditivity of sup,

$$\mathfrak{R}_{n,p}(\phi_+ \circ \mathcal{G}) \leq \mathfrak{R}_{n,p}(\ell(\cdot, +1) \circ \mathcal{G}) + \beta \mathfrak{R}_{n,p}(\psi \circ \mathcal{G}), \quad \mathfrak{R}_{n,p}(\phi_- \circ \mathcal{G}) = \mathfrak{R}_{n,p}(\ell(\cdot, -1) \circ \mathcal{G}).$$

For the classification loss, the contraction lemma implies

$$\mathfrak{R}_{n,p}(\ell(\cdot, y) \circ \mathcal{G}) \leq L_\ell \mathfrak{R}_{n,p}(\mathcal{G}) \quad (y \in \mathcal{Y}).$$

For  $\psi$ , applying Lemma A.2 and Lemma A.4, we obtain

$$\begin{aligned} \mathfrak{R}_{n,p}(\psi \circ \mathcal{G}) &\leq 2L_{\text{KL}} \left( \mathfrak{R}_{n,p}(\mathcal{G}) + \mathfrak{R}_{n,p}(\{ \mathbf{x} \mapsto g(\mathbf{x} + \boldsymbol{\eta}) : \|\boldsymbol{\eta}\|_\infty \leq \varepsilon, g \in \mathcal{G} \}) \right) \\ &\leq 2L_{\text{KL}} \left( 2\mathfrak{R}_{n,p}(\mathcal{G}) + \frac{\varepsilon W d^{1/q}}{\sqrt{n}} \right). \end{aligned} \quad (\text{A38})$$

Substituting these bounds into (A37) and simplifying yields

$$\begin{aligned} \mathbb{E} \left[ \sup_{g \in \mathcal{G}} \{ \widehat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \} \right] &\leq 4\pi_P (L_\ell + 2\beta L_{\text{KL}}) \mathfrak{R}_{n_P, p_P}(\mathcal{G}) \\ &\quad + 2(L_\ell + 4\beta L_{\text{KL}}) \mathfrak{R}_{n_U, p_U}(\mathcal{G}) \\ &\quad + 4\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_P}{\sqrt{n_P}} + \frac{1}{\sqrt{n_U}} \right). \end{aligned} \quad (\text{A39})$$

as desired.

**(iii) Combining the results of (i) and (ii)**

Substituting the expectation bound in (ii), i.e., (A39), into the McDiarmid inequality result in (i), i.e., (A30), we obtain, with probability at least  $1 - \delta/2$ ,

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left\{ \hat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \right\} &\leq 4\pi_{\text{P}}(L_{\ell} + 2\beta L_{\text{KL}})\mathfrak{R}_{n_{\text{P}}, p_{\text{P}}}(\mathcal{G}) + 2(L_{\ell} + 4\beta L_{\text{KL}})\mathfrak{R}_{n_{\text{U}}, p_{\text{U}}}(\mathcal{G}) \\ &\quad + 4\beta L_{\text{KL}} \varepsilon W d^{1/q} \left( \frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}} \right) \\ &\quad + \sqrt{\frac{1}{2} \ln \frac{2}{\delta}} \left( \frac{\pi_{\text{P}}(2C_{\ell} + \beta C_{\text{KL}})}{\sqrt{n_{\text{P}}}} + \frac{C_{\ell} + \beta C_{\text{KL}}}{\sqrt{n_{\text{U}}}} \right). \end{aligned} \tag{A40}$$

This establishes the one-sided bound.

**(iv) Opposite direction and union bound** Similarly,  $\sup_{g \in \mathcal{G}} \{R_{\text{nnPU-TR}}(g) - \hat{R}_{\text{nnPU-TR}}(g)\}$  is also bounded by the right-hand side of (A40) with probability at least  $1 - \delta/2$ . Therefore, by the union bound, (31) holds with probability at least  $1 - \delta$ .  $\square$

## A.7 Proof of Theorem 5.6

*Proof* Since  $\hat{g}_{\text{nnPU-TR}}$  is an empirical risk minimizer, we have

$$\hat{R}_{\text{nnPU-TR}}(\hat{g}_{\text{nnPU-TR}}) \leq \hat{R}_{\text{nnPU-TR}}(g^*)$$

Moreover, on the event where Lemma 5.7 holds (which occurs with probability at least  $1 - \delta$ ),

$$\begin{aligned} &R_{\text{nnPU-TR}}(\hat{g}_{\text{nnPU-TR}}) - R_{\text{nnPU-TR}}(g^*) \\ &= \left( R_{\text{nnPU-TR}}(\hat{g}_{\text{nnPU-TR}}) - \hat{R}_{\text{nnPU-TR}}(\hat{g}_{\text{nnPU-TR}}) \right) \\ &\quad + \left( \hat{R}_{\text{nnPU-TR}}(\hat{g}_{\text{nnPU-TR}}) - \hat{R}_{\text{nnPU-TR}}(g^*) \right) \\ &\quad + \left( \hat{R}_{\text{nnPU-TR}}(g^*) - R_{\text{nnPU-TR}}(g^*) \right) \\ &\leq \sup_{g \in \mathcal{G}} \left( R_{\text{nnPU-TR}}(g) - \hat{R}_{\text{nnPU-TR}}(g) \right) + 0 + \sup_{g \in \mathcal{G}} \left( \hat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \right) \\ &\leq 2 \sup_{g \in \mathcal{G}} \left| \hat{R}_{\text{nnPU-TR}}(g) - R_{\text{nnPU-TR}}(g) \right|. \end{aligned}$$

Therefore, multiplying (31) by 2 yields (30).  $\square$

## A.8 Proof of the Theorem “Condition on the Number of Unlabeled Samples for PU+TRADES to Outperform Supervised TRADES” (Theorem 5.8)

*Proof* (i) Comparing (35) and (36), and canceling  $\Gamma_{\delta}(\pi_{\text{P}}/\sqrt{n_{\text{P}}})$  from both sides, we obtain

$$\Gamma_{\delta} \frac{\pi_{\text{N}}}{\sqrt{n_{\text{N}}}} > \frac{\pi_{\text{P}}}{\sqrt{n_{\text{P}}}} \left( 4L_{\ell} W C_x d^{1/q} + \kappa_{\delta} C_{\ell} \right) + \frac{\Gamma_{\delta}}{\sqrt{n_{\text{U}}}}.$$

Rearranging by moving the first term on the right-hand side, and solving for  $\Gamma_{\delta}/\sqrt{n_{\text{U}}}$ , we obtain (38) under condition (37). Part (ii) follows immediately from the fact that the bound for nnPU+TRADES has the same form as that for uPU+TRADES.  $\square$