
Positive-Unlabeled Learning with Non-Negative Risk Estimator

Ryuichi Kiryo^{1,2} Gang Niu^{1,2} Marthinus C. du Plessis Masashi Sugiyama^{2,1}

¹The University of Tokyo, 7-3-1 Hongo, Tokyo 113-0033, Japan

²RIKEN, 1-4-1 Nihonbashi, Tokyo 103-0027, Japan

{ kiryo@ms., gang@ms., sugi@ }k.u-tokyo.ac.jp

Abstract

From only *positive* (P) and *unlabeled* (U) data, a binary classifier could be trained with PU learning, in which the state of the art is *unbiased PU learning*. However, if its model is very flexible, empirical risks on training data will go negative, and we will suffer from serious overfitting. In this paper, we propose a *non-negative risk estimator* for PU learning: when getting minimized, it is more robust against overfitting, and thus we are able to use very flexible models (such as deep neural networks) given limited P data. Moreover, we analyze the *bias*, *consistency*, and *mean-squared-error reduction* of the proposed risk estimator, and bound the *estimation error* of the resulting *empirical risk minimizer*. Experiments demonstrate that our risk estimator fixes the overfitting problem of its unbiased counterparts.

1 Introduction

Positive-unlabeled (PU) *learning* can be dated back to [1, 2, 3] and has been well studied since then. It mainly focuses on binary classification applied to retrieval and novelty or outlier detection tasks [4, 5, 6, 7], while it also has applications in matrix completion [8] and sequential data [9, 10].

Existing PU methods can be divided into two categories based on how U data is handled. The first category (e.g., [11, 12]) identifies possible *negative* (N) data in U data, and then performs ordinary supervised (PN) learning; the second (e.g., [13, 14]) regards U data as N data with smaller weights. The former heavily relies on the heuristics for identifying N data; the latter heavily relies on good choices of the weights of U data, which is computationally expensive to tune.

In order to avoid tuning the weights, *unbiased PU learning* comes into play as a subcategory of the second category. The milestone is [4], which regards a U data as weighted P and N data simultaneously. It might lead to *unbiased risk estimators*, if we unrealistically assume that the class-posterior probability is one for all P data.¹ A breakthrough in this direction is [15] for proposing the first unbiased risk estimator, and a more general estimator was suggested in [16] as a common foundation. The former is unbiased but non-convex for loss functions satisfying some symmetric condition; the latter is always unbiased, and it is further convex for loss functions meeting some linear-odd condition [17, 18]. PU learning based on these unbiased risk estimators is the current state of the art.

However, the unbiased risk estimators will give negative empirical risks, if the model being trained is very flexible. For the general estimator in [16], there exist three partial risks in the total risk (see Eq. (2) defined later), especially it has a negative risk regarding P data as N data to cancel the bias caused by regarding U data as N data. The worst case is that the model can realize any measurable function and the loss function is not upper bounded, so that the empirical risk is not lower bounded. This needs to be fixed since the original risk, which is the target to be estimated, is non-negative.

¹It implies the P and N class-conditional densities have disjoint support sets, and then any P and N data (as the test data) can be perfectly separated by a fixed classifier that is sufficiently flexible.

To this end, we propose a novel *non-negative risk estimator* that follows and improves on the state-of-the-art unbiased risk estimators mentioned above. This estimator can be used for two purposes: first, given some validation data (which are also PU data), we can use our estimator to evaluate the risk—for this case it is *biased yet optimal*, and for some symmetric losses, the *mean-squared-error reduction* is guaranteed; second, given some training data, we can use our estimator to train binary classifiers—for this case its risk minimizer possesses an *estimation error bound* of the same order as the risk minimizers corresponding to its unbiased counterparts [15, 16, 19].

In addition, we propose a *large-scale* PU learning algorithm for minimizing the unbiased and non-negative risk estimators. This algorithm accepts any *surrogate loss* and is based on *stochastic optimization*, e.g., [20]. Note that [21] is the only existing large-scale PU algorithm, but it only accepts a single surrogate loss from [16] and is based on *sequential minimal optimization* [22].

The rest of this paper is organized as follows. In Section 2 we review unbiased PU learning, and in Section 3 we propose non-negative PU learning. Theoretical analyses are carried out in Section 4, and experimental results are discussed in Section 5. Conclusions are given in Section 6.

2 Unbiased PU learning

In this section, we review unbiased PU learning [15, 16].

Problem settings Let $X \in \mathbb{R}^d$ and $Y \in \{\pm 1\}$ ($d \in \mathbb{N}$) be the input and output random variables. Let $p(x, y)$ be the *underlying joint density* of (X, Y) , $p_p(x) = p(x | Y = +1)$ and $p_n(x) = p(x | Y = -1)$ be the *P and N marginals* (a.k.a. the P and N class-conditional densities), $p(x)$ be the *U marginal*, $\pi_p = p(Y = +1)$ be the *class-prior probability*, and $\pi_n = p(Y = -1) = 1 - \pi_p$. π_p is assumed known throughout the paper; it can be estimated from P and U data [23, 24, 25, 26].

Consider the *two-sample problem setting* of PU learning [5]: two sets of data are sampled independently from $p_p(x)$ and $p(x)$ as $\mathcal{X}_p = \{x_i^p\}_{i=1}^{n_p} \sim p_p(x)$ and $\mathcal{X}_u = \{x_i^u\}_{i=1}^{n_u} \sim p(x)$, and a classifier needs to be trained from \mathcal{X}_p and \mathcal{X}_u .² If it is PN learning as usual, $\mathcal{X}_n = \{x_i^n\}_{i=1}^{n_n} \sim p_n(x)$ rather than \mathcal{X}_u would be available and a classifier could be trained from \mathcal{X}_p and \mathcal{X}_n .

Risk estimators Unbiased PU learning relies on unbiased risk estimators. Let $g : \mathbb{R}^d \rightarrow \mathbb{R}$ be an arbitrary *decision function*, and $\ell : \mathbb{R} \times \{\pm 1\} \rightarrow \mathbb{R}$ be the *loss function*, such that the value $\ell(t, y)$ means the loss incurred by predicting an output t when the ground truth is y . Denote by $R_p^+(g) = \mathbb{E}_p[\ell(g(X), +1)]$ and $R_n^-(g) = \mathbb{E}_n[\ell(g(X), -1)]$, where $\mathbb{E}_p[\cdot] = \mathbb{E}_{X \sim p_p}[\cdot]$ and $\mathbb{E}_n[\cdot] = \mathbb{E}_{X \sim p_n}[\cdot]$. Then, the *risk* of g is $R(g) = \mathbb{E}_{(X, Y) \sim p(x, y)}[\ell(g(X), Y)] = \pi_p R_p^+(g) + \pi_n R_n^-(g)$. In PN learning, thanks to the availability of \mathcal{X}_p and \mathcal{X}_n , $R(g)$ can be approximated directly by

$$\hat{R}_{pn}(g) = \pi_p \hat{R}_p^+(g) + \pi_n \hat{R}_n^-(g), \quad (1)$$

where $\hat{R}_p^+(g) = (1/n_p) \sum_{i=1}^{n_p} \ell(g(x_i^p), +1)$ and $\hat{R}_n^-(g) = (1/n_n) \sum_{i=1}^{n_n} \ell(g(x_i^n), -1)$. In PU learning, \mathcal{X}_n is unavailable, but $R_n^-(g)$ can be approximated indirectly, as shown in [15, 16]. Denote by $R_p^-(g) = \mathbb{E}_p[\ell(g(X), -1)]$ and $R_u^-(g) = \mathbb{E}_{X \sim p(x)}[\ell(g(X), -1)]$. As $\pi_n p_n(x) = p(x) - \pi_p p_p(x)$, we can obtain that $\pi_n R_n^-(g) = R_u^-(g) - \pi_p R_p^-(g)$, and $R(g)$ can be approximated indirectly by

$$\hat{R}_{pu}(g) = \pi_p \hat{R}_p^+(g) - \pi_p \hat{R}_p^-(g) + \hat{R}_u^-(g), \quad (2)$$

where $\hat{R}_p^-(g) = (1/n_p) \sum_{i=1}^{n_p} \ell(g(x_i^p), -1)$ and $\hat{R}_u^-(g) = (1/n_u) \sum_{i=1}^{n_u} \ell(g(x_i^u), -1)$.

The *empirical risk estimators* in Eqs. (1) and (2) are *unbiased* and *consistent* w.r.t. all popular loss functions.³ When they are used for evaluating the risk (e.g., in cross-validation), ℓ is by default the *zero-one loss*, namely $\ell_{01}(t, y) = (1 - \text{sign}(ty))/2$; when used for training, ℓ_{01} is replaced with a *surrogate loss* [27]. In particular, [15] showed that if ℓ satisfies a *symmetric condition*:

$$\ell(t, +1) + \ell(t, -1) = 1, \quad (3)$$

² \mathcal{X}_p is a set of independent data and so is \mathcal{X}_u , but $\mathcal{X}_p \cup \mathcal{X}_u$ does not need to be such a set.

³The consistency here means for fixed g , $\hat{R}_{pn}(g) \rightarrow R(g)$ and $\hat{R}_{pu}(g) \rightarrow R(g)$ as $n_p, n_n, n_u \rightarrow \infty$.

we will have

$$\hat{R}_{\text{pu}}(g) = 2\pi_{\text{p}}\hat{R}_{\text{p}}^+(g) + \hat{R}_{\text{u}}^-(g) - \pi_{\text{p}}, \quad (4)$$

which can be minimized by separating \mathcal{X}_{p} and \mathcal{X}_{u} with ordinary cost-sensitive learning. An issue is $\hat{R}_{\text{pu}}(g)$ in (4) must be non-convex in g , since no $\ell(t, y)$ in (3) can be convex in t . [16] showed that $\hat{R}_{\text{pu}}(g)$ in (2) is convex in g , if $\ell(t, y)$ is convex in t and meets a *linear-odd condition* [17, 18]:

$$\ell(t, +1) - \ell(t, -1) = -t. \quad (5)$$

Let g be parameterized by θ , then (5) leads to a convex optimization problem so long as g is linear in θ , for which the globally optimal solution can be obtained. Eq. (5) is not only sufficient but also necessary for the convexity, if ℓ is unary, i.e., $\ell(t, -1) = \ell(-t, +1)$.

Justification Thanks to the unbiasedness, we can study *estimation error bounds* (EEB). Let \mathcal{G} be the *function class*, and \hat{g}_{pn} and \hat{g}_{pu} be the *empirical risk minimizers* of $\hat{R}_{\text{pn}}(g)$ and $\hat{R}_{\text{pu}}(g)$. [19] proved EEB of \hat{g}_{pu} is tighter than EEB of \hat{g}_{pn} when $\pi_{\text{p}}/\sqrt{n_{\text{p}}} + 1/\sqrt{n_{\text{u}}} < \pi_{\text{n}}/\sqrt{n_{\text{n}}}$, if (a) ℓ satisfies (3) and is *Lipschitz continuous*; (b) the *Rademacher complexity* of \mathcal{G} decays in $\mathcal{O}(1/\sqrt{n})$ for data of size n drawn from $p(x)$, $p_{\text{p}}(x)$ or $p_{\text{n}}(x)$.⁴ In other words, under mild conditions, PU learning is likely to outperform PN learning when $\pi_{\text{p}}/\sqrt{n_{\text{p}}} + 1/\sqrt{n_{\text{u}}} < \pi_{\text{n}}/\sqrt{n_{\text{n}}}$. This phenomenon has been observed in experiments [19] and is illustrated in Figure 1(a).

3 Non-negative PU learning

In this section, we propose the non-negative risk estimator and the large-scale PU algorithm.

3.1 Motivation

Let us look inside the aforementioned justification of unbiased PU (uPU) learning. Intuitively, the advantage comes from the transformation $\pi_{\text{n}}R_{\text{n}}^-(g) = R_{\text{u}}^-(g) - \pi_{\text{p}}R_{\text{p}}^-(g)$. When we approximate $\pi_{\text{n}}R_{\text{n}}^-(g)$ from N data $\{x_i^{\text{n}}\}_{i=1}^{n_{\text{n}}}$, the convergence rate is $\mathcal{O}_p(\pi_{\text{n}}/\sqrt{n_{\text{n}}})$, where \mathcal{O}_p denotes the order in probability; when we approximate $R_{\text{u}}^-(g) - \pi_{\text{p}}R_{\text{p}}^-(g)$ from P data $\{x_i^{\text{p}}\}_{i=1}^{n_{\text{p}}}$ and U data $\{x_i^{\text{u}}\}_{i=1}^{n_{\text{u}}}$, the convergence rate becomes $\mathcal{O}_p(\pi_{\text{p}}/\sqrt{n_{\text{p}}} + 1/\sqrt{n_{\text{u}}})$. As a result, we might benefit from a tighter *uniform deviation bound* when $\pi_{\text{p}}/\sqrt{n_{\text{p}}} + 1/\sqrt{n_{\text{u}}} < \pi_{\text{n}}/\sqrt{n_{\text{n}}}$.

However, the critical assumption on the Rademacher complexity is indispensable, otherwise it will be difficult for EEB of \hat{g}_{pu} to be tighter than EEB of \hat{g}_{pn} . If $\mathcal{G} = \{g \mid \|g\|_{\infty} \leq C_g\}$ where $C_g > 0$ is a constant, i.e., it has all measurable functions with some bounded norm, then $\mathfrak{R}_{n,q}(\mathcal{G}) = \mathcal{O}(1)$ for any n and $q(x)$ and all bounds become trivial; moreover if ℓ is not bounded from above, $\hat{R}_{\text{pu}}(g)$ becomes not bounded from below, i.e., it may diverge to $-\infty$. Thus, in order to obtain high-quality \hat{g}_{pu} , \mathcal{G} cannot be too complex, or equivalently the model of g cannot be too flexible.

This argument is supported by an experiment as illustrated in Figure 1(b). A *multilayer perceptron* was trained for separating the even and odd digits of MNIST hand-written digits [29]. This model is so flexible that the number of parameters is 500 times more than the total number of P and N data. From Figure 1(b) we can see:

- (A) on training data, the risks of uPU and PN both decrease, and uPU is faster than PN;
- (B) on test data, the risk of PN decreases, whereas the risk of uPU does not; the risk of uPU is lower at the beginning but higher at the end than that of PN.

To sum up, the overfitting problem of uPU is serious, which evidences that in order to obtain high-quality \hat{g}_{pu} , the model of g cannot be too flexible.

3.2 Non-negative risk estimator

Nevertheless, we have no choice sometimes: we are interested in using flexible models, while labeling more data is out of our control. Can we alleviate the overfitting problem with neither changing the model nor labeling more data?

⁴Let $\sigma_1, \dots, \sigma_n$ be n Rademacher variables, the Rademacher complexity of \mathcal{G} for \mathcal{X} of size n drawn from $q(x)$ is defined by $\mathfrak{R}_{n,q}(\mathcal{G}) = \mathbb{E}_{\mathcal{X}} \mathbb{E}_{\sigma_1, \dots, \sigma_n} [\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{x_i \in \mathcal{X}} \sigma_i g(x_i)]$ [28]. For any fixed \mathcal{G} and q , $\mathfrak{R}_{n,q}(\mathcal{G})$ still depends on n and should decrease with n .

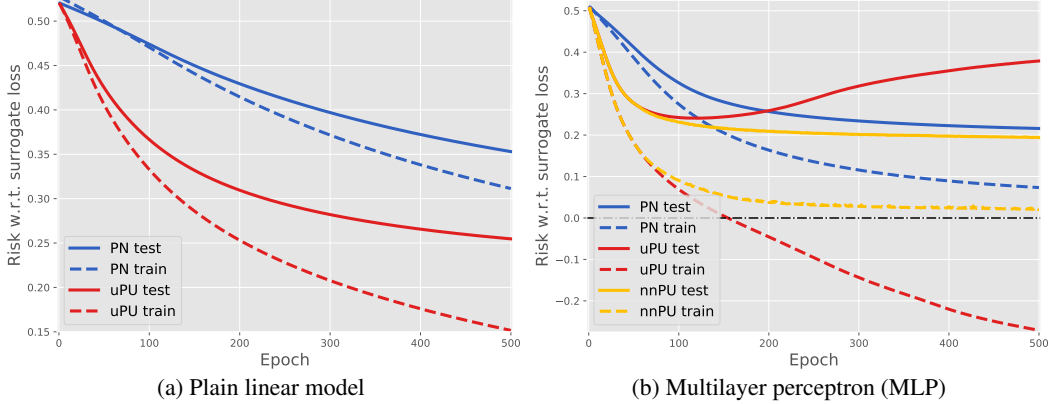


Figure 1: Illustrative experimental results.

The answer is affirmative. Note that $\hat{R}_{\text{pu}}(\hat{g}_{\text{pu}})$ keeps decreasing and goes negative. This should be fixed since $R(g) \geq 0$ for any g . Specifically, it holds that $R_{\text{u}}^-(g) - \pi_{\text{p}} R_{\text{p}}^-(g) = \pi_{\text{n}} R_{\text{n}}^-(g) \geq 0$, but $\hat{R}_{\text{u}}^-(g) - \pi_{\text{p}} \hat{R}_{\text{p}}^-(g) \geq 0$ is not always true, which is a potential reason for uPU to overfit. Based on this key observation, we propose a *non-negative risk estimator* for PU learning:

$$\tilde{R}_{\text{pu}}(g) = \pi_{\text{p}} \hat{R}_{\text{p}}^+(g) + \max \left\{ 0, \hat{R}_{\text{u}}^-(g) - \pi_{\text{p}} \hat{R}_{\text{p}}^-(g) \right\}. \quad (6)$$

Let $\tilde{g}_{\text{pu}} = \arg \min_{g \in \mathcal{G}} \tilde{R}_{\text{pu}}(g)$ be the empirical risk minimizer of $\tilde{R}_{\text{pu}}(g)$. We refer to the process of obtaining \tilde{g}_{pu} as *non-negative PU (nnPU) learning*. The implementation of nnPU will be given in Section 3.3, and theoretical analyses of $\tilde{R}_{\text{pu}}(g)$ and \tilde{g}_{pu} will be given in Section 4.

Again, from Figure 1(b) we can see:

- (A) on training data, the risk of nnPU first decreases and then becomes more and more flat, so that the risk of nnPU is closer to the risk of PN and farther from that of uPU; in short, the risk of nnPU does not go down with uPU after a certain epoch;
- (B) on test data, the tendency is similar, but the risk of nnPU does not go up with uPU;
- (C) at the end, nnPU achieves the lowest risk on test data.

In summary, nnPU works by explicitly constraining the training risk of uPU to be non-negative.

3.3 Implementation

A list of popular loss functions and their properties is shown in Table 1. Let g be parameterized by θ . If g is linear in θ , the losses satisfying (5) result in convex optimizations. However, if g needs to be flexible, it will be highly nonlinear in θ ; then the losses satisfying (5) are not advantageous over others, since the optimizations are anyway non-convex. In [15], the *ramp loss* was used and $\hat{R}_{\text{pu}}(g)$ was minimized by the *concave-convex procedure* [30]. This solver is fairly sophisticated, and if we replace $\hat{R}_{\text{pu}}(g)$ with $\tilde{R}_{\text{pu}}(g)$, it will be more difficult to implement. To this end, we propose to use the *sigmoid loss* $\ell_{\text{sig}}(t, y) = 1/(1 + \exp(ty))$: its gradient is everywhere non-zero and $\tilde{R}_{\text{pu}}(g)$ can be minimized by off-the-shelf gradient methods.

In front of big data, we should scale PU learning up by stochastic optimization. Minimizing $\hat{R}_{\text{pu}}(g)$ is *embarrassingly parallel* while minimizing $\tilde{R}_{\text{pu}}(g)$ is not, since $\hat{R}_{\text{pu}}(g)$ is *point-wise* but $\tilde{R}_{\text{pu}}(g)$ is not due to the max operator. That being said, $\max\{0, \hat{R}_{\text{u}}^-(g; \mathcal{X}_{\text{u}}) - \pi_{\text{p}} \hat{R}_{\text{p}}^-(g; \mathcal{X}_{\text{p}})\}$ is no greater than $(1/N) \sum_{i=1}^N \max\{0, \hat{R}_{\text{u}}^-(g; \mathcal{X}_{\text{u}}^i) - \pi_{\text{p}} \hat{R}_{\text{p}}^-(g; \mathcal{X}_{\text{p}}^i)\}$, where $(\mathcal{X}_{\text{p}}^i, \mathcal{X}_{\text{u}}^i)$ is the i -th mini-batch, and hence the corresponding upper bound of $\tilde{R}_{\text{pu}}(g)$ can easily be minimized in parallel.

Table 1: Loss functions for PU learning and their properties.

| Name | Definition | (3) | (5) | Bounded | Lipschitz | $\ell'(z) \neq 0$ |
|-------------------|-----------------------------------|-----|-----|---------|-----------|-----------------------|
| Zero-one loss | $(1 - \text{sign}(z))/2$ | ✓ | × | ✓ | × | $z = 0$ |
| Ramp loss | $\max\{0, \min\{1, (1 - z)/2\}\}$ | ✓ | × | ✓ | ✓ | $z \in [-1, +1]$ |
| Squared loss | $(z - 1)^2/4$ | × | ✓ | × | × | $z \in \mathbb{R}$ |
| Logistic loss | $\ln(1 + \exp(-z))$ | × | ✓ | × | ✓ | $z \in \mathbb{R}$ |
| Hinge loss | $\max\{0, 1 - z\}$ | × | × | × | ✓ | $z \in (-\infty, +1]$ |
| Double hinge loss | $\max\{0, (1 - z)/2, -z\}$ | × | ✓ | × | ✓ | $z \in (-\infty, +1]$ |
| Sigmoid loss | $1/(1 + \exp(z))$ | ✓ | × | ✓ | ✓ | $z \in \mathbb{R}$ |

All loss functions are unary, such that $\ell(t, y) = \ell(z)$ with $z = ty$. The ramp loss comes from [15]; the double hinge loss is from [16], in which the squared, logistic and hinge losses were discussed as well. The ramp and squared losses are scaled to satisfy (3) or (5). The sigmoid loss is a horizontally mirrored *logistic function*; the logistic loss is the negative logarithm of the logistic function.

Algorithm 1 Large-scale PU learning based on stochastic optimization

Input: training data $(\mathcal{X}_p, \mathcal{X}_u)$;
hyperparameters $0 \leq \beta \leq \pi_p \sup_t \max_y \ell(t, y)$ and $0 \leq \gamma \leq 1$
Output: model parameter θ for $\hat{g}_{pu}(x; \theta)$ or $\tilde{g}_{pu}(x; \theta)$

- 1: Let \mathcal{A} be an external SGD-like stochastic optimization algorithm such as [20] or [31]
- 2: **while** no stopping criterion has been met:
- 3: Shuffle $(\mathcal{X}_p, \mathcal{X}_u)$ into N mini-batches, and denote by $(\mathcal{X}_p^i, \mathcal{X}_u^i)$ the i -th mini-batch
- 4: **for** $i = 1$ **to** N :
- 5: **if** $\hat{R}_u^-(g; \mathcal{X}_u^i) - \pi_p \hat{R}_p^-(g; \mathcal{X}_p^i) \geq -\beta$:
- 6: Set gradient $\nabla_{\theta} \hat{R}_{pu}(g; \mathcal{X}_p^i, \mathcal{X}_u^i)$
- 7: Update θ by \mathcal{A} with its current step size η
- 8: **else**:
- 9: Set gradient $\nabla_{\theta} (\pi_p \hat{R}_p^-(g; \mathcal{X}_p^i) - \hat{R}_u^-(g; \mathcal{X}_u^i))$
- 10: Update θ by \mathcal{A} with a discounted step size $\gamma\eta$

The large-scale PU algorithm is described in Algorithm 1. Let $r_i = \hat{R}_u^-(g; \mathcal{X}_u^i) - \pi_p \hat{R}_p^-(g; \mathcal{X}_p^i)$. In practice, we may tolerate $r_i \geq -\beta$ where $0 \leq \beta \leq \pi_p \sup_t \max_y \ell(t, y)$, as r_i comes from a single mini-batch. The degree of tolerance is controlled by β : there is zero tolerance if $\beta = 0$, and we are minimizing $\hat{R}_{pu}(g)$ if $\beta = \pi_p \sup_t \max_y \ell(t, y)$. Otherwise if $r_i < -\beta$, we go along $-\nabla_{\theta} r_i$ with a step size discounted by γ where $0 \leq \gamma \leq 1$, to make this mini-batch less overfitted. Algorithm 1 is insensitive to the choice of γ , if the optimization algorithm \mathcal{A} is adaptive such as [20] or [31].

4 Theoretical analyses

In this section, we analyze the risk estimator (6) and its minimizer (all proofs are in Appendix B).

4.1 Bias and consistency

Fix g , $\tilde{R}_{pu}(g) \geq \hat{R}_{pu}(g)$ for any $(\mathcal{X}_p, \mathcal{X}_u)$ but $\hat{R}_{pu}(g)$ is unbiased, which implies $\tilde{R}_{pu}(g)$ is biased in general. A fundamental question is then whether $\tilde{R}_{pu}(g)$ is consistent. From now on, we prove this consistency. To begin with, partition all possible $(\mathcal{X}_p, \mathcal{X}_u)$ into $\mathfrak{D}^+(g) = \{(\mathcal{X}_p, \mathcal{X}_u) \mid \hat{R}_u^-(g) - \pi_p \hat{R}_p^-(g) \geq 0\}$ and $\mathfrak{D}^-(g) = \{(\mathcal{X}_p, \mathcal{X}_u) \mid \hat{R}_u^-(g) - \pi_p \hat{R}_p^-(g) < 0\}$. Assume there are $C_g > 0$ and $C_{\ell} > 0$ such that $\sup_{g \in \mathcal{G}} \|g\|_{\infty} \leq C_g$ and $\sup_{|t| \leq C_g} \max_y \ell(t, y) \leq C_{\ell}$.

Lemma 1. *The following three conditions are equivalent: (A) the probability measure of $\mathfrak{D}^-(g)$ is non-zero; (B) $\tilde{R}_{pu}(g)$ differs from $\hat{R}_{pu}(g)$ with a non-zero probability over repeated sampling of $(\mathcal{X}_p, \mathcal{X}_u)$; (C) the bias of $\tilde{R}_{pu}(g)$ is positive. In addition, by assuming that there is $\alpha > 0$ such that $R_u^-(g) \geq \alpha$, the probability measure of $\mathfrak{D}^-(g)$ can be bounded by*

$$\Pr(\mathfrak{D}^-(g)) \leq \exp(-2(\alpha/C_{\ell})^2/(\pi_p^2/n_p + 1/n_u)). \quad (7)$$

Based on Lemma 1, we can show the exponential decay of the bias and also the consistency. For convenience, denote by $\chi_{n_p, n_u} = 2\pi_p/\sqrt{n_p} + 1/\sqrt{n_u}$.

Theorem 2 (Bias and consistency). *Assume that $R_n^-(g) \geq \alpha > 0$ and denote by Δ_g the right-hand side of Eq. (7). As $n_p, n_u \rightarrow \infty$, the bias of $\tilde{R}_{pu}(g)$ decays exponentially:*

$$0 \leq \mathbb{E}_{\mathcal{X}_p, \mathcal{X}_u}[\tilde{R}_{pu}(g)] - R(g) \leq C_\ell \pi_p \Delta_g. \quad (8)$$

Moreover, for any $\delta > 0$, let $C_\delta = C_\ell \sqrt{\ln(2/\delta)/2}$, then we have with probability at least $1 - \delta$,

$$|\tilde{R}_{pu}(g) - R(g)| \leq C_\delta \cdot \chi_{n_p, n_u} + C_\ell \pi_p \Delta_g, \quad (9)$$

and with probability at least $1 - \delta - \Delta_g$,

$$|\tilde{R}_{pu}(g) - R(g)| \leq C_\delta \cdot \chi_{n_p, n_u}. \quad (10)$$

Either (9) or (10) in Theorem 2 indicates for fixed g , $\tilde{R}_{pu}(g) \rightarrow R(g)$ in $\mathcal{O}_p(\pi_p/\sqrt{n_p} + 1/\sqrt{n_u})$. This convergence rate is optimal according to the *central limit theorem* [32], which means the proposed estimator is a biased yet optimal estimator to the risk.

4.2 Mean squared error

After introducing the bias, $\tilde{R}_{pu}(g)$ tends to overestimate $R(g)$. It is not a *shrinkage estimator* [33, 34] so that its *mean squared error* (MSE) is not necessarily smaller than that of $\hat{R}_{pu}(g)$. However, we can still characterize this reduction in MSE.

Theorem 3 (MSE reduction). *It holds that $\text{MSE}(\tilde{R}_{pu}(g)) < \text{MSE}(\hat{R}_{pu}(g))$,⁵ if and only if*

$$\int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} (\hat{R}_{pu}(g) + \tilde{R}_{pu}(g) - 2R(g))(\hat{R}_u^-(g) - \pi_p \hat{R}_p^-(g)) dF(\mathcal{X}_p, \mathcal{X}_u) > 0, \quad (11)$$

where $dF(\mathcal{X}_p, \mathcal{X}_u) = \prod_{i=1}^{n_p} p_p(x_i^p) dx_i^p \cdot \prod_{i=1}^{n_u} p_u(x_i^u) dx_i^u$. Eq. (11) is valid, if the following conditions are met: (a) $\Pr(\mathfrak{D}^-(g)) > 0$; (b) ℓ satisfies Eq. (3); (c) $R_n^-(g) \geq \alpha > 0$; (d) $n_u \gg n_p$, such that we have $R_u^-(g) - \hat{R}_u^-(g) \leq 2\alpha$ almost surely on $\mathfrak{D}^-(g)$. In fact, given these four conditions, we have for any $0 \leq \beta \leq C_\ell \pi_p$,

$$\text{MSE}(\hat{R}_{pu}(g)) - \text{MSE}(\tilde{R}_{pu}(g)) \geq 3\beta^2 \Pr\{\tilde{R}_{pu}(g) - \hat{R}_{pu}(g) > \beta\}. \quad (12)$$

The assumption (d) in Theorem 3 is explained as follows. Since U data can be much cheaper than P data in practice, it would be natural to assume n_u is much larger and grows much faster than n_p , hence $\Pr\{\hat{R}_u^-(g) - \hat{R}_u^-(g) \geq \alpha\} / \Pr\{\hat{R}_p^-(g) - \hat{R}_p^-(g) \geq \alpha/\pi_p\} \propto \exp(n_p - n_u)$ asymptotically.⁶ This means the contribution of \mathcal{X}_u is negligible for making $(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)$ so that $\Pr(\mathfrak{D}^-(g))$ exhibits exponential decay mainly in n_p . As $\Pr\{\hat{R}_u^-(g) - \hat{R}_u^-(g) \geq 2\alpha\}$ has stronger exponential decay in n_u than $\Pr\{\hat{R}_u^-(g) - \hat{R}_u^-(g) \geq \alpha\}$ as well as $n_u \gg n_p$, we made the assumption (d).

4.3 Estimation error

While Theorems 2 and 3 addressed the use of (6) for evaluating the risk, we are likewise interested in its use for training classifiers. In what follows, we analyze the estimation error $R(\tilde{g}_{pu}) - R(g^*)$, where g^* is the true risk minimizer in \mathcal{G} , i.e., $g^* = \arg \min_{g \in \mathcal{G}} R(g)$. As a common practice [28], assume that $\ell(t, y)$ is Lipschitz continuous in t for all $|t| \leq C_g$ with a Lipschitz constant L_ℓ .

Theorem 4 (Estimation error bound). *Assume that (a) $\inf_{g \in \mathcal{G}} R_n^-(g) \geq \alpha > 0$ and denote by Δ the right-hand side of Eq. (7); (b) \mathcal{G} is closed under negation, i.e., $g \in \mathcal{G}$ if and only if $-g \in \mathcal{G}$. Then, for any $\delta > 0$, with probability at least $1 - \delta$,*

$$R(\tilde{g}_{pu}) - R(g^*) \leq 16L_\ell \pi_p \mathfrak{R}_{n_p, p_p}(\mathcal{G}) + 8L_\ell \mathfrak{R}_{n_u, p}(\mathcal{G}) + 2C'_\delta \cdot \chi_{n_p, n_u} + 2C_\ell \pi_p \Delta, \quad (13)$$

where $C'_\delta = C_\ell \sqrt{\ln(1/\delta)/2}$, and $\mathfrak{R}_{n_p, p_p}(\mathcal{G})$ and $\mathfrak{R}_{n_u, p}(\mathcal{G})$ are the Rademacher complexities of \mathcal{G} for the sampling of size n_p from $p_p(x)$ and of size n_u from $p(x)$, respectively.

⁵Here, $\text{MSE}(\cdot)$ is over repeated sampling of $(\mathcal{X}_p, \mathcal{X}_u)$.

⁶This can be derived as $n_p, n_u \rightarrow \infty$ by applying the *central limit theorem* to the two differences and then *L'Hôpital's rule* to the ratio of complementary error functions [32].

Theorem 4 ensures that learning with (6) is also consistent: as $n_p, n_u \rightarrow \infty$, $R(\tilde{g}_{pu}) \rightarrow R(g^*)$ and if ℓ satisfies (5), all optimizations are convex and $\tilde{g}_{pu} \rightarrow g^*$. For linear-in-parameter models with a bounded norm, $\mathfrak{R}_{n_p, p_p}(\mathcal{G}) = \mathcal{O}(1/\sqrt{n_p})$ and $\mathfrak{R}_{n_u, p}(\mathcal{G}) = \mathcal{O}(1/\sqrt{n_u})$, and thus $R(\tilde{g}_{pu}) \rightarrow R(g^*)$ in $\mathcal{O}_p(\pi_p/\sqrt{n_p} + 1/\sqrt{n_u})$.

For comparison, $R(\hat{g}_{pu}) - R(g^*)$ can be bounded using a *different proof technique* [19]:

$$R(\hat{g}_{pu}) - R(g^*) \leq 8L_\ell \pi_p \mathfrak{R}_{n_p, p_p}(\mathcal{G}) + 4L_\ell \mathfrak{R}_{n_u, p}(\mathcal{G}) + 2C_\delta \cdot \chi_{n_p, n_u}, \quad (14)$$

where $C_\delta = C_\ell \sqrt{\ln(2/\delta)/2}$. The differences of (13) and (14) are completely from the differences of the corresponding uniform deviation bounds, i.e., the following lemma and Lemma 8 of [19].

Lemma 5. *Under the assumptions of Theorem 4, for any $\delta > 0$, with probability at least $1 - \delta$,*

$$\sup_{g \in \mathcal{G}} |\tilde{R}_{pu}(g) - R(g)| \leq 8L_\ell \pi_p \mathfrak{R}_{n_p, p_p}(\mathcal{G}) + 4L_\ell \mathfrak{R}_{n_u, p}(\mathcal{G}) + C'_\delta \cdot \chi_{n_p, n_u} + C_\ell \pi_p \Delta. \quad (15)$$

Notice that $\hat{R}_{pu}(g)$ is point-wise while $\tilde{R}_{pu}(g)$ is not due to the maximum, which makes Lemma 5 much more difficult to prove than Lemma 8 of [19]. The key trick is that after *symmetrization*, we employ $|\max\{0, z\} - \max\{0, z'\}| \leq |z - z'|$, making three differences of partial risks point-wise (see (18) in the proof). As a consequence, we have to use a different Rademacher complexity *with the absolute value inside the supremum* [35, 36], whose *contraction* makes the coefficients of (15) doubled compared with Lemma 8 of [19]; moreover, we have to assume \mathcal{G} is closed under negation to change back to the standard Rademacher complexity *without the absolute value* [28]. Therefore, the differences of (13) and (14) are mainly due to different proof techniques and cannot reflect the intrinsic differences of empirical risk minimizers.

5 Experiments

In this section, we compare PN, unbiased PU (uPU) and non-negative PU (nnPU) learning experimentally. We focus on training deep neural networks, as uPU learning usually does not overfit if a linear-in-parameter model is used [19] and nothing needs to be fixed.

Table 2 describes the specification of benchmark datasets. MNIST, 20News and CIFAR-10 have 10, 7 and 10 classes originally, and we constructed the P and N classes from them as follows: MNIST was preprocessed in such a way that 0, 2, 4, 6, 8 constitute the P class, while 1, 3, 5, 7, 9 constitute the N class; for 20News, ‘alt.’, ‘comp.’, ‘misc.’ and ‘rec.’ make up the P class, and ‘sci.’, ‘soc.’ and ‘talk.’ make up the N class; for CIFAR-10, the P class is formed by ‘airplane’, ‘automobile’, ‘ship’ and ‘truck’, and the N class is formed by ‘bird’, ‘cat’, ‘deer’, ‘dog’, ‘frog’ and ‘horse’. The dataset epsilon has 2 classes and such a construction is unnecessary.

Three learning methods were set up as follows: (A) for PN, $n_p = 1,000$ and $n_n = (\pi_n/2\pi_p)^2 n_p$; (B) for uPU, $n_p = 1,000$ and n_u is the total number of training data; (C) for nnPU, n_p and n_u are exactly same as uPU. For uPU and nnPU, P and U data were dependent, because neither $\hat{R}_{pu}(g)$ in Eq. (2) nor $\tilde{R}_{pu}(g)$ in Eq. (6) requires them to be independent. The choice of n_n was motivated by [19] and may make nnPU potentially better than PN as $n_u \rightarrow \infty$ (whether $n_p < \infty$ or $n_p \leq n_u$).

The model for MNIST was a 6-layer *multilayer perceptron* (MLP) with ReLU [40] (more specifically, d -300-300-300-300-1). For epsilon, the model was similar while the activation was replaced with Softsign [41] for better performance. For 20News, we borrowed the pre-trained word embeddings from GloVe [42], and the model can be written as d -avg_pool(word_emb(d ,300))-300-300-1,

Table 2: Specification of benchmark datasets, models, and optimization algorithms.

| Name | # Train | # Test | # Feature | π_p | Model $g(x; \theta)$ | Opt. alg. \mathcal{A} |
|---------------|---------|---------|-----------|---------|---------------------------|-------------------------|
| MNIST [29] | 60,000 | 10,000 | 784 | 0.49 | 6-layer MLP with ReLU | Adam [20] |
| epsilon [37] | 400,000 | 100,000 | 2,000 | 0.50 | 6-layer MLP with Softsign | Adam [20] |
| 20News [38] | 11,314 | 7,532 | 61,188 | 0.44 | 5-layer MLP with Softsign | AdaGrad [31] |
| CIFAR-10 [39] | 50,000 | 10,000 | 3,072 | 0.40 | 13-layer CNN with ReLU | Adam [20] |

See <http://yann.lecun.com/exdb/mnist/> for MNIST, <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html> for epsilon, <http://qwone.com/~jason/20Newsgroups/> for 20News, and <https://www.cs.toronto.edu/~kriz/cifar.html> for CIFAR-10.

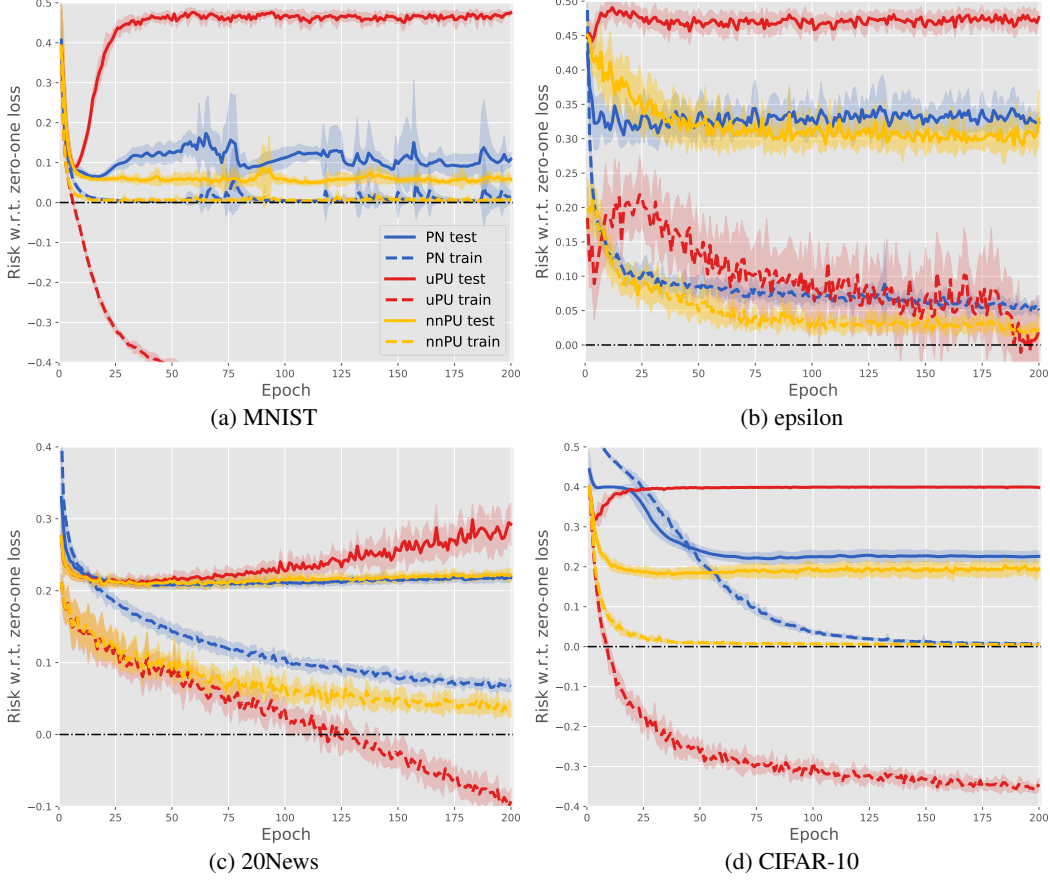


Figure 2: Experimental results of training deep neural networks.

where $\text{word_emb}(d,300)$ retrieves 300-dimensional word embeddings for all words in a document, avg_pool executes average pooling, and the resulting vector is fed to a 4-layer MLP with Softsign. The model for CIFAR-10 was an *all convolutional net* [43]: $(32*32*3)-[C(3*3,96)]*2-C(3*3,96,2)-[C(3*3,192)]*2-C(3*3,192,2)-C(3*3,192)-C(1*1,192)-C(1*1,10)-1000-1000-1$, where the input is a $32*32$ RGB image, $C(3*3,96)$ means 96 channels of $3*3$ convolutions followed by ReLU, $[\cdot]*2$ means there are two such layers, $C(3*3,96,2)$ means a similar layer but with stride 2, etc.; it is one of the best architectures for CIFAR-10. Batch normalization [44] was applied before hidden layers. Furthermore, the sigmoid loss ℓ_{sig} was used as the surrogate loss and an ℓ_2 -regularization was also added. The resulting objectives were minimized by Adam [20] on MNIST, epsilon and CIFAR-10, and by AdaGrad [31] on 20News; we fixed $\beta = 0$ and $\gamma = 1$ for simplicity.

The experimental results are reported in Figure 2, where means and standard deviations of training and test risks based on the same 10 random samplings are shown. We can see that uPU overfitted training data and nnPU fixed this problem. Additionally, given limited N data, nnPU outperformed PN on MNIST, epsilon and CIFAR-10 and was comparable to it on 20News. In summary, with the proposed non-negative risk estimator, we are able to use very flexible models given limited P data.

We further tried some cases where π_p is misspecified, in order to simulate PU learning in the wild, where we must suffer from errors in estimating π_p . More specifically, we tested nnPU learning by replacing π_p with $\pi'_p \in \{0.8\pi_p, 0.9\pi_p, \dots, 1.2\pi_p\}$ and giving π'_p to the learning method, so that it would regard π'_p as π_p during the entire training process. The experimental setup was exactly same as before except the replacement of π_p .

The experimental results are reported in Figure 3, where means of test risks of nnPU based on the same 10 random samplings are shown, and the best test risks are identified (horizontal lines are the best mean test risks and vertical lines are the epochs when they were achieved). We can see that on MNIST, the more misspecification was, the worse nnPU performed, while under-misspecification hurt more than over-misspecification; on epsilon, the cases where π'_p equals to π_p , $1.1\pi_p$ and $1.2\pi_p$

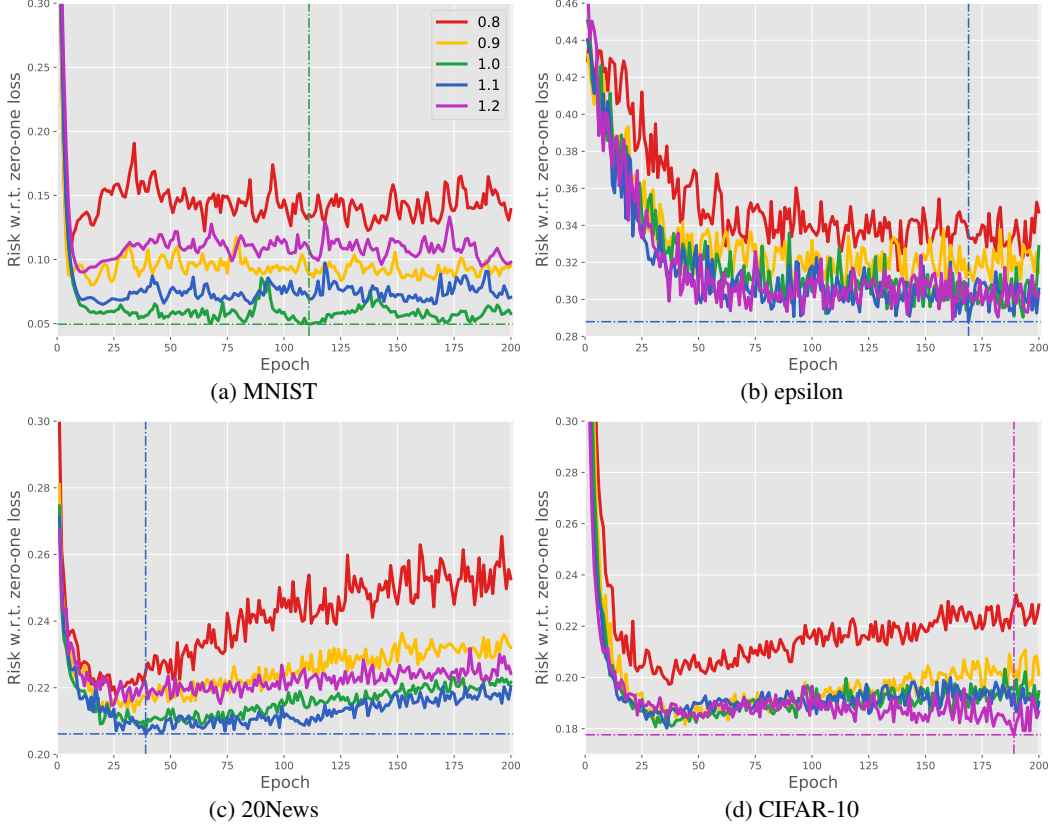


Figure 3: Experimental results given $\pi'_p \in \{0.8\pi_p, 0.9\pi_p, \dots, 1.2\pi_p\}$.

were comparable, but the best was $\pi'_p = 1.1\pi_p$ rather than $\pi'_p = \pi_p$; on 20News, these three cases became different, such that $\pi'_p = \pi_p$ was superior to $\pi'_p = 1.2\pi_p$ but inferior to $\pi'_p = 1.1\pi_p$; at last on CIFAR-10, $\pi'_p = \pi_p$ and $\pi'_p = 1.1\pi_p$ were comparable again, and $\pi'_p = 1.2\pi_p$ was the winner.

In all the experiments, we have fixed $\beta = 0$, which may explain this phenomenon. Recall that uPU overfitted seriously on all the benchmark datasets, and note that the larger π'_p is, the more different nnPU is from uPU. Therefore, the replacement of π_p with some $\pi'_p > \pi_p$ introduces additional bias of $\tilde{R}_{pu}(g)$ in estimating $R(g)$, but it also pushes $\tilde{R}_{pu}(g)$ away from $\hat{R}_{pu}(g)$ and then pushes nnPU away from uPU. This may result in lower test risks given some π'_p slightly larger than π_p as shown in Figure 3. This is also why under-misspecified π'_p hurt more than over-misspecified π'_p .

All the experiments were done with *Chainer* [45], and our implementation based on it is available at <https://github.com/kiryor/nnPUlearning>.

6 Conclusions

We proposed a non-negative risk estimator for PU learning that follows and improves on the state-of-the-art unbiased risk estimators. No matter how flexible the model is, it will not go negative as its unbiased counterparts. It is more robust against overfitting when being minimized, and training very flexible models such as deep neural networks given limited P data becomes possible. We also developed a large-scale PU learning algorithm. Extensive theoretical analyses were presented, and the usefulness of our non-negative PU learning was verified by intensive experiments. A promising future direction is extending the current work to semi-supervised learning along [46].

Acknowledgments

GN and MS were supported by JST CREST JPMJCR1403 and GN was also partially supported by Microsoft Research Asia.

References

- [1] F. Denis. PAC learning from positive statistical queries. In *ALT*, 1998.
- [2] F. De Comit  , F. Denis, R. Gilleron, and F. Letouzey. Positive and unlabeled examples help learning. In *ALT*, 1999.
- [3] F. Letouzey, F. Denis, and R. Gilleron. Learning from positive and unlabeled examples. In *ALT*, 2000.
- [4] C. Elkan and K. Noto. Learning classifiers from only positive and unlabeled data. In *KDD*, 2008.
- [5] G. Ward, T. Hastie, S. Barry, J. Elith, and J. Leathwick. Presence-only data and the EM algorithm. *Biometrics*, 65(2):554–563, 2009.
- [6] C. Scott and G. Blanchard. Novelty detection: Unlabeled data definitely help. In *AISTATS*, 2009.
- [7] G. Blanchard, G. Lee, and C. Scott. Semi-supervised novelty detection. *Journal of Machine Learning Research*, 11:2973–3009, 2010.
- [8] C.-J. Hsieh, N. Natarajan, and I. S. Dhillon. PU learning for matrix completion. In *ICML*, 2015.
- [9] X. Li, P. S. Yu, B. Liu, and S.-K. Ng. Positive unlabeled learning for data stream classification. In *SDM*, 2009.
- [10] M. N. Nguyen, X. Li, and S.-K. Ng. Positive unlabeled learning for time series classification. In *IJCAI*, 2011.
- [11] B. Liu, W. S. Lee, P. S. Yu, and X. Li. Partially supervised classification of text documents. In *ICML*, 2002.
- [12] X. Li and B. Liu. Learning to classify texts using positive and unlabeled data. In *IJCAI*, 2003.
- [13] W. S. Lee and B. Liu. Learning with positive and unlabeled examples using weighted logistic regression. In *ICML*, 2003.
- [14] B. Liu, Y. Dai, X. Li, W. S. Lee, and P. S. Yu. Building text classifiers using positive and unlabeled examples. In *ICDM*, 2003.
- [15] M. C. du Plessis, G. Niu, and M. Sugiyama. Analysis of learning from positive and unlabeled data. In *NIPS*, 2014.
- [16] M. C. du Plessis, G. Niu, and M. Sugiyama. Convex formulation for learning from positive and unlabeled data. In *ICML*, 2015.
- [17] N. Natarajan, I. S. Dhillon, P. Ravikumar, and A. Tewari. Learning with noisy labels. In *NIPS*, 2013.
- [18] G. Patrini, F. Nielsen, R. Nock, and M. Carioni. Loss factorization, weakly supervised learning and label noise robustness. In *ICML*, 2016.
- [19] G. Niu, M. C. du Plessis, T. Sakai, Y. Ma, and M. Sugiyama. Theoretical comparisons of positive-unlabeled learning against positive-negative learning. In *NIPS*, 2016.
- [20] D. P. Kingma and J. L. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [21] E. Sansone, F. G. B. De Natale, and Z.-H. Zhou. Efficient training for positive unlabeled learning. *arXiv preprint arXiv:1608.06807*, 2016.
- [22] J. C. Platt. Fast training of support vector machines using sequential minimal optimization. In B. Sch  lkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods*, pages 185–208. MIT Press, 1999.
- [23] A. Menon, B. Van Rooyen, C. S. Ong, and B. Williamson. Learning from corrupted binary labels via class-probability estimation. In *ICML*, 2015.
- [24] H. G. Ramaswamy, C. Scott, and A. Tewari. Mixture proportion estimation via kernel embedding of distributions. In *ICML*, 2016.
- [25] S. Jain, M. White, and P. Radivojac. Estimating the class prior and posterior from noisy positives and unlabeled data. In *NIPS*, 2016.
- [26] M. C. du Plessis, G. Niu, and M. Sugiyama. Class-prior estimation for learning from positive and unlabeled data. *Machine Learning*, 106(4):463–492, 2017.
- [27] P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- [28] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. MIT Press, 2012.
- [29] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [30] A. L. Yuille and A. Rangarajan. The concave-convex procedure (CCCP). In *NIPS*, 2001.

- [31] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- [32] K.-L. Chung. *A Course in Probability Theory*. Academic Press, 1968.
- [33] C. Stein. Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. In *Proc. 3rd Berkeley Symposium on Mathematical Statistics and Probability*, 1956.
- [34] W. James and C. Stein. Estimation with quadratic loss. In *Proc. 4th Berkeley Symposium on Mathematical Statistics and Probability*, 1961.
- [35] V. Koltchinskii. Rademacher penalties and structural risk minimization. *IEEE Transactions on Information Theory*, 47(5):1902–1914, 2001.
- [36] P. L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- [37] G.-X. Yuan, C.-H. Ho, and C.-J. Lin. An improved GLMNET for l1-regularized logistic regression. *Journal of Machine Learning Research*, 13:1999–2030, 2012.
- [38] K. Lang. Newsweeder: Learning to filter netnews. In *ICML*, 1995.
- [39] A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [40] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.
- [41] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010.
- [42] J. Pennington, R. Socher, and C. D. Manning. GloVe: Global vectors for word representation. In *EMNLP*, 2014.
- [43] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller. Striving for simplicity: The all convolutional net. In *ICLR*, 2015.
- [44] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [45] S. Tokui, K. Oono, S. Hido, and J. Clayton. Chainer: a next-generation open source framework for deep learning. In *Machine Learning Systems Workshop at NIPS*, 2015.
- [46] T. Sakai, M. C. du Plessis, G. Niu, and M. Sugiyama. Semi-supervised classification based on classification from positive and unlabeled data. In *ICML*, 2017.
- [47] C. McDiarmid. On the method of bounded differences. In J. Siemons, editor, *Surveys in Combinatorics*, pages 148–188. Cambridge University Press, 1989.
- [48] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*. Springer, 1991.
- [49] S. Shalev-Shwartz and S. Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.
- [50] V. N. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, 1998.

A Supplementary experimental results

Due to limited space, we considered the surrogate loss without the zero-one loss in Figure 1. Here, we include the zero-one loss and show the extended version of Figure 1 in Figure 4. In general, the curves of risks w.r.t. ℓ_{01} look quite similar to (but less smooth than) those w.r.t. ℓ_{sig} . Therefore, the curves of risks w.r.t. ℓ_{sig} are more visually appealing as the illustrative experimental results.

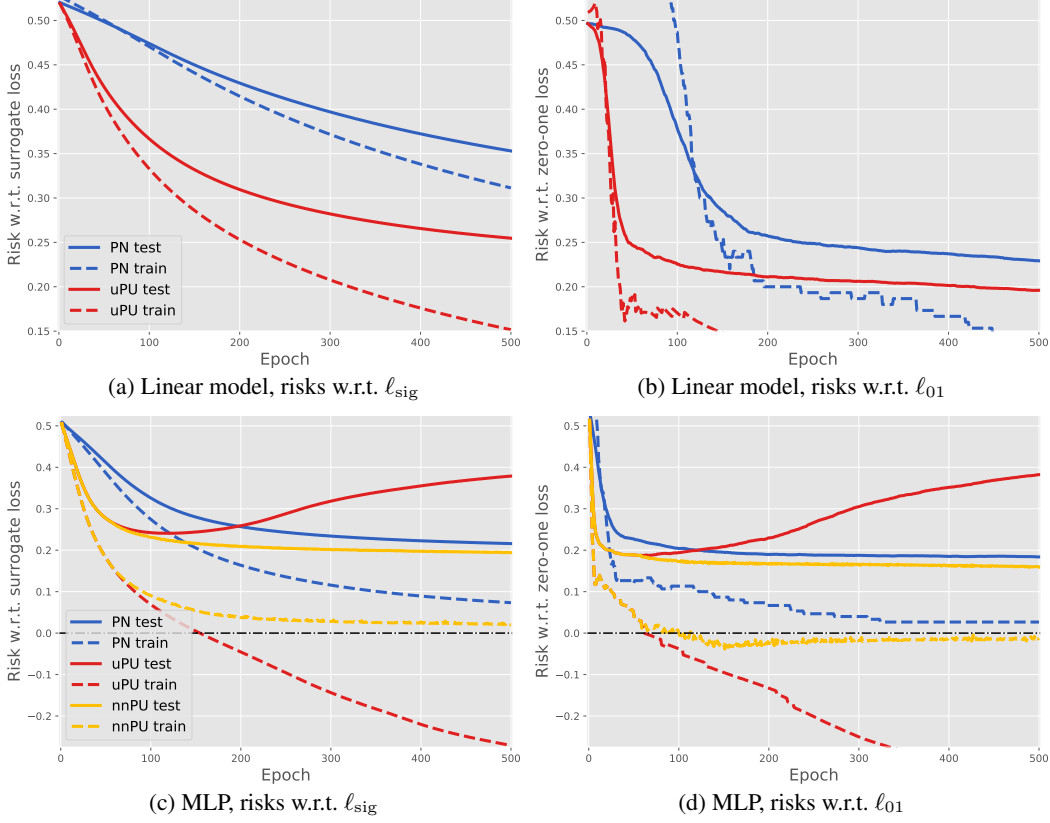


Figure 4: The extended version of Figure 1.

B Proofs

In this appendix, we prove all the theoretical results in Section 4.

B.1 Proof of Lemma 1

Let

$$p_{\text{p}}(\mathcal{X}_{\text{p}}) = p_{\text{p}}(x_1^{\text{p}}) \cdots p_{\text{p}}(x_{n_{\text{p}}}^{\text{p}}), \quad p(\mathcal{X}_{\text{u}}) = p(x_1^{\text{u}}) \cdots p(x_{n_{\text{u}}}^{\text{u}})$$

be the probability density functions of \mathcal{X}_{p} and \mathcal{X}_{u} . Then let $F_{\text{p}}(\mathcal{X}_{\text{p}})$ be the cumulative distribution function of \mathcal{X}_{p} , $F_{\text{u}}(\mathcal{X}_{\text{u}})$ be that of \mathcal{X}_{u} , and

$$F(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) = F_{\text{p}}(\mathcal{X}_{\text{p}}) \cdot F_{\text{u}}(\mathcal{X}_{\text{u}})$$

be the joint cumulative distribution function of $(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}})$. Given the above definitions, the measure of $\mathfrak{D}^-(g)$ is defined by

$$\Pr(\mathfrak{D}^-(g)) = \int_{(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \in \mathfrak{D}^-(g)} dF(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}),$$

where \Pr denotes the probability. Since $\tilde{R}_{\text{pu}}(g)$ is identical to $\hat{R}_{\text{pu}}(g)$ on $\mathfrak{D}^+(g)$ and different from $\hat{R}_{\text{pu}}(g)$ on $\mathfrak{D}^-(g)$, we have $\Pr(\mathfrak{D}^-(g)) = \Pr\{\tilde{R}_{\text{pu}}(g) \neq \hat{R}_{\text{pu}}(g)\}$. That is, the measure of $\mathfrak{D}^-(g)$ is non-zero if and only if $\tilde{R}_{\text{pu}}(g)$ differs from $\hat{R}_{\text{pu}}(g)$ with a non-zero probability.

Based on the facts that $\hat{R}_{\text{pu}}(g)$ is unbiased and $\tilde{R}_{\text{pu}}(g) - \hat{R}_{\text{pu}}(g) = 0$ on $\mathfrak{D}^+(g)$, we have

$$\begin{aligned} \mathbb{E}[\tilde{R}_{\text{pu}}(g)] - R(g) &= \mathbb{E}[\tilde{R}_{\text{pu}}(g) - \hat{R}_{\text{pu}}(g)] \\ &= \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^+(g)} \tilde{R}_{\text{pu}}(g) - \hat{R}_{\text{pu}}(g) dF(\mathcal{X}_p, \mathcal{X}_u) \\ &\quad + \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} \tilde{R}_{\text{pu}}(g) - \hat{R}_{\text{pu}}(g) dF(\mathcal{X}_p, \mathcal{X}_u) \\ &= \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} \tilde{R}_{\text{pu}}(g) - \hat{R}_{\text{pu}}(g) dF(\mathcal{X}_p, \mathcal{X}_u). \end{aligned}$$

As a result, $\mathbb{E}[\tilde{R}_{\text{pu}}(g)] - R(g) > 0$ if and only if $\int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} dF(\mathcal{X}_p, \mathcal{X}_u) > 0$ due to the fact $\tilde{R}_{\text{pu}}(g) - \hat{R}_{\text{pu}}(g) > 0$ on $\mathfrak{D}^-(g)$. That is, the bias of $\tilde{R}_{\text{pu}}(g)$ is positive if and only if the measure of $\mathfrak{D}^-(g)$ is non-zero.

We prove (7) by *the method of bounded differences*, for that

$$\mathbb{E}[\hat{R}_{\text{u}}^-(g) - \pi_p \hat{R}_{\text{p}}^-(g)] = R_{\text{u}}^-(g) - \pi_p R_{\text{p}}^-(g) = R_{\text{n}}^-(g) \geq \alpha.$$

We have assumed that $0 \leq \ell(t, \pm 1) \leq C_\ell$, and thus the change of $\hat{R}_{\text{p}}^-(g)$ will be no more than C_ℓ/n_p if some $x_i^p \in \mathcal{X}_p$ is replaced, or the change of $\hat{R}_{\text{u}}^-(g)$ will be no more than C_ℓ/n_u if some $x_i^u \in \mathcal{X}_u$ is replaced. Subsequently, *McDiarmid's inequality* [47] implies

$$\begin{aligned} \Pr\{R_{\text{n}}^-(g) - (\hat{R}_{\text{u}}^-(g) - \pi_p \hat{R}_{\text{p}}^-(g)) \geq \alpha\} &\leq \exp\left(-\frac{2\alpha^2}{n_p(C_\ell\pi_p/n_p)^2 + n_u(C_\ell/n_u)^2}\right) \\ &= \exp\left(-\frac{2\alpha^2/C_\ell^2}{\pi_p^2/n_p + 1/n_u}\right). \end{aligned}$$

Taking into account that

$$\begin{aligned} \Pr(\mathfrak{D}^-(g)) &= \Pr\{\hat{R}_{\text{u}}^-(g) - \pi_p \hat{R}_{\text{p}}^-(g) < 0\} \\ &\leq \Pr\{\hat{R}_{\text{u}}^-(g) - \pi_p \hat{R}_{\text{p}}^-(g) \leq R_{\text{n}}^-(g) - \alpha\} \\ &= \Pr\{R_{\text{n}}^-(g) - (\hat{R}_{\text{u}}^-(g) - \pi_p \hat{R}_{\text{p}}^-(g)) \geq \alpha\}, \end{aligned}$$

we complete the proof. \square

B.2 Proof of Theorem 2

It has been proven in Lemma 1 that

$$\mathbb{E}[\tilde{R}_{\text{pu}}(g)] - R(g) = \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} \tilde{R}_{\text{pu}}(g) - \hat{R}_{\text{pu}}(g) dF(\mathcal{X}_p, \mathcal{X}_u),$$

and thus the exponential decay of the bias in (8) is obtained via

$$\begin{aligned} \mathbb{E}[\tilde{R}_{\text{pu}}(g)] - R(g) &\leq \sup_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} (\tilde{R}_{\text{pu}}(g) - \hat{R}_{\text{pu}}(g)) \cdot \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} dF(\mathcal{X}_p, \mathcal{X}_u) \\ &= \sup_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} (\pi_p \hat{R}_{\text{p}}^-(g) - \hat{R}_{\text{u}}^-(g)) \cdot \Pr(\mathfrak{D}^-(g)) \\ &\leq C_\ell \pi_p \Delta_g. \end{aligned}$$

The deviation bound (9) is due to

$$\begin{aligned} |\tilde{R}_{\text{pu}}(g) - R(g)| &\leq |\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]| + |\mathbb{E}[\tilde{R}_{\text{pu}}(g)] - R(g)| \\ &\leq |\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]| + C_\ell \pi_p \Delta_g. \end{aligned}$$

The change of $\tilde{R}_{\text{pu}}(g)$ will be no more than $2C_\ell/n_p$ if some $x_i^p \in \mathcal{X}_p$ is replaced, or it will be no more than C_ℓ/n_u if some $x_i^u \in \mathcal{X}_u$ is replaced, and McDiarmid's inequality gives us

$$\Pr\{|\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]| \geq \epsilon\} \leq 2 \exp\left(-\frac{2\epsilon^2}{n_p(2C_\ell\pi_p/n_p)^2 + n_u(C_\ell/n_u)^2}\right),$$

or equivalently, with probability at least $1 - \delta$,

$$\begin{aligned} |\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]| &\leq \sqrt{\frac{\ln(2/\delta)C_\ell^2}{2} \left(\frac{4\pi_p^2}{n_p} + \frac{1}{n_u}\right)} \\ &\leq C_\delta \left(\frac{2\pi_p}{\sqrt{n_p}} + \frac{1}{\sqrt{n_u}}\right) \\ &= C_\delta \cdot \chi_{n_p, n_u}. \end{aligned}$$

On the other hand, the deviation bound (10) is due to

$$|\tilde{R}_{\text{pu}}(g) - R(g)| \leq |\tilde{R}_{\text{pu}}(g) - \hat{R}_{\text{pu}}(g)| + |\hat{R}_{\text{pu}}(g) - R(g)|,$$

where $|\tilde{R}_{\text{pu}}(g) - \hat{R}_{\text{pu}}(g)| > 0$ with probability at most Δ_g , and $|\hat{R}_{\text{pu}}(g) - R(g)|$ shares the same concentration inequality with $|\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]|$. \square

B.3 Proof of Theorem 3

For convenience, let $A = \pi_p \hat{R}_p^+(g)$ and $B = \hat{R}_u^-(g) - \pi_p \hat{R}_p^-(g)$, so that

$$R(g) = \mathbb{E}[A + B], \quad \hat{R}_{\text{pu}}(g) = A + B, \quad \tilde{R}_{\text{pu}}(g) = A + B_+,$$

where $B_+ = \max\{0, B\}$. Subsequently, let $R = R(g)$ for short, and then by definition,

$$\begin{aligned} \text{MSE}(\hat{R}_{\text{pu}}(g)) &= \mathbb{E}[(A + B - R)^2] \\ &= \mathbb{E}[(A + B)^2] - 2R \cdot \mathbb{E}[A + B] + R^2, \\ \text{MSE}(\tilde{R}_{\text{pu}}(g)) &= \mathbb{E}[(A + B_+ - R)^2] \\ &= \mathbb{E}[(A + B_+)^2] - 2R \cdot \mathbb{E}[A + B_+] + R^2. \end{aligned}$$

Hence,

$$\begin{aligned} \text{MSE}(\hat{R}_{\text{pu}}(g)) - \text{MSE}(\tilde{R}_{\text{pu}}(g)) &= \mathbb{E}[(A + B)^2] - \mathbb{E}[(A + B_+)^2] \\ &\quad - 2R \cdot (\mathbb{E}[A + B] - \mathbb{E}[A + B_+]). \end{aligned}$$

The first part $\mathbb{E}[(A + B)^2] - \mathbb{E}[(A + B_+)^2]$ can be rewritten as

$$\begin{aligned} \mathbb{E}[(A + B)^2] - \mathbb{E}[(A + B_+)^2] &= \mathbb{E}[2A(B - B_+) + B^2 - B_+^2] \\ &= \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^+(g)} 2A(B - B) + B^2 - B^2 dF(\mathcal{X}_p, \mathcal{X}_u) \\ &\quad + \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} 2A(B - 0) + B^2 - 0^2 dF(\mathcal{X}_p, \mathcal{X}_u) \\ &= \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} 2AB + B^2 dF(\mathcal{X}_p, \mathcal{X}_u). \end{aligned}$$

The second part $2R \cdot (\mathbb{E}[A + B] - \mathbb{E}[A + B_+])$ can be rewritten as

$$\begin{aligned} 2R \cdot (\mathbb{E}[A + B] - \mathbb{E}[A + B_+]) &= 2R \cdot \mathbb{E}[B - B_+] \\ &= 2R \cdot \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^+(g)} B - B dF(\mathcal{X}_p, \mathcal{X}_u) \\ &\quad + 2R \cdot \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} B - 0 dF(\mathcal{X}_p, \mathcal{X}_u) \\ &= \int_{(\mathcal{X}_p, \mathcal{X}_u) \in \mathfrak{D}^-(g)} 2RB dF(\mathcal{X}_p, \mathcal{X}_u). \end{aligned}$$

As a consequence,

$$\text{MSE}(\widehat{R}_{\text{pu}}(g)) - \text{MSE}(\widetilde{R}_{\text{pu}}(g)) = \int_{(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \in \mathfrak{D}^-(g)} (2A + B - 2R)B \, dF(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}),$$

which is exactly the left-hand side of (11) since $\widetilde{R}_{\text{pu}}(g) = A$ on $\mathfrak{D}^-(g)$.

In order to prove the rest, it suffices to show that $A - R \leq B$ on $\mathfrak{D}^-(g)$. By the assumption that ℓ satisfies (3),

$$\begin{aligned} A - R &= A - \mathbb{E}[A] - \mathbb{E}[B] \\ &= \pi_{\text{p}} \widehat{R}_{\text{p}}^+(g) - \pi_{\text{p}} R_{\text{p}}^+(g) - \mathbb{E}[B] \\ &= \pi_{\text{p}} R_{\text{p}}^-(g) - \pi_{\text{p}} \widehat{R}_{\text{p}}^-(g) - \mathbb{E}[B]. \end{aligned}$$

Thus, with probability one,

$$\begin{aligned} A - R &= \pi_{\text{p}} R_{\text{p}}^-(g) - \pi_{\text{p}} \widehat{R}_{\text{p}}^-(g) - \mathbb{E}[B] + (\widehat{R}_{\text{u}}^-(g) - \widehat{R}_{\text{u}}^-(g)) + (R_{\text{u}}^-(g) - R_{\text{u}}^-(g)) \\ &= (\widehat{R}_{\text{u}}^-(g) - \pi_{\text{p}} \widehat{R}_{\text{p}}^-(g)) - (R_{\text{u}}^-(g) - \pi_{\text{p}} R_{\text{p}}^-(g)) - \mathbb{E}[B] + (R_{\text{u}}^-(g) - \widehat{R}_{\text{u}}^-(g)) \\ &= B - 2\mathbb{E}[B] + (R_{\text{u}}^-(g) - \widehat{R}_{\text{u}}^-(g)) \\ &\leq B, \end{aligned}$$

where we used the assumptions that $\mathbb{E}[B] \geq \alpha$ and $R_{\text{u}}^-(g) - \widehat{R}_{\text{u}}^-(g) \leq 2\alpha$ almost surely on $\mathfrak{D}^-(g)$. To sum up, we have established that

$$\int_{(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \in \mathfrak{D}^-(g)} (2A + B - 2R)B \, dF(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \geq 3 \int_{(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \in \mathfrak{D}^-(g)} B^2 \, dF(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}).$$

Due to the fact that $B^2 > 0$ on $\mathfrak{D}^-(g)$ and the assumption that $\Pr(\mathfrak{D}^-(g)) > 0$, we know Eq. (11) is valid. Finally, for any $0 \leq \beta \leq C_{\ell} \pi_{\text{p}}$, it is clear that

$$\{(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \mid B < -\beta\} \subseteq \{(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \mid B < 0\} = \mathfrak{D}^-(g),$$

and $B < -\beta$ if and only if $\widetilde{R}_{\text{pu}}(g) - \widehat{R}_{\text{pu}}(g) > \beta$. These two facts imply that

$$\begin{aligned} \int_{(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \in \mathfrak{D}^-(g)} B^2 \, dF(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) &\geq \int_{(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \mid B < -\beta} B^2 \, dF(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \\ &\geq \beta^2 \int_{(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \mid B < -\beta} dF(\mathcal{X}_{\text{p}}, \mathcal{X}_{\text{u}}) \\ &= \beta^2 \Pr\{B < -\beta\} \\ &= \beta^2 \Pr\{\widetilde{R}_{\text{pu}}(g) - \widehat{R}_{\text{pu}}(g) > \beta\}, \end{aligned}$$

which proves (12) and the whole theorem. \square

B.4 Proof of Lemma 5

Preliminary An alternative definition of the Rademacher complexity will be used in the proof:

$$\mathfrak{R}'_{n,q}(\mathcal{G}) = \mathbb{E}_{\mathcal{X}} \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\sup_{g \in \mathcal{G}} \left| \frac{1}{n} \sum_{x_i \in \mathcal{X}} \sigma_i g(x_i) \right| \right].$$

For the sake of comparison, the one we have used in the statements of theoretical results is

$$\mathfrak{R}_{n,q}(\mathcal{G}) = \mathbb{E}_{\mathcal{X}} \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{x_i \in \mathcal{X}} \sigma_i g(x_i) \right].$$

This alternative version comes from [35, 36] of which authors are the pioneers of error bounds based on the Rademacher complexity. Without any composition, $\mathfrak{R}'_{n,q}(\mathcal{G}) \geq \mathfrak{R}_{n,q}(\mathcal{G})$ for arbitrary \mathcal{G} and $\mathfrak{R}'_{n,q}(\mathcal{G}) = \mathfrak{R}_{n,q}(\mathcal{G})$ if \mathcal{G} is closed under negation. However, with a composition

$$\ell \circ \mathcal{G} = \{\ell \circ g \mid g \in \mathcal{G}\}$$

where the loss ℓ is non-negative, the Rademacher complexity of the *composite function class* would generally not satisfy $\mathfrak{R}'_{n,q}(\ell \circ \mathcal{G}) = \mathfrak{R}_{n,q}(\ell \circ \mathcal{G})$ since $\ell \circ \mathcal{G}$ is generally not closed under negation. Furthermore, a vital disagreement arises when considering the contraction principle or property: if $\psi : \mathbb{R} \rightarrow \mathbb{R}$ is a Lipschitz continuous function with a Lipschitz constant L_ψ and satisfies $\psi(0) = 0$, we have

$$\begin{aligned}\mathfrak{R}_{n,q}(\psi \circ \mathcal{G}) &\leq L_\psi \mathfrak{R}_{n,q}(\mathcal{G}), \\ \mathfrak{R}'_{n,q}(\psi \circ \mathcal{G}) &\leq 2L_\psi \mathfrak{R}'_{n,q}(\mathcal{G}),\end{aligned}$$

according to *Talagrand's contraction lemma* [48] and its extension [28, 49]. Here, for $\mathfrak{R}_{n,q}(\psi \circ \mathcal{G})$ we can use Lemma 4.2 in [28] or Lemma 26.9 in [49] where $\psi(0) = 0$ is safely dropped, while for $\mathfrak{R}'_{n,q}(\psi \circ \mathcal{G})$ we have to use the original Theorem 4.12 in [48] where $\psi(0) = 0$ is required. In fact, the name of the lemma is after that ψ is a contraction if $\psi(0) = 0$ and $L_\psi = 1$.

Proof Firstly, we deal with the bias of $\tilde{R}_{\text{pu}}(g)$:

$$\begin{aligned}\sup_{g \in \mathcal{G}} |\tilde{R}_{\text{pu}}(g) - R(g)| &\leq \sup_{g \in \mathcal{G}} |\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]| + \sup_{g \in \mathcal{G}} |\mathbb{E}[\tilde{R}_{\text{pu}}(g)] - R(g)| \\ &\leq \sup_{g \in \mathcal{G}} |\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]| + C_\ell \pi_p \Delta,\end{aligned}\quad (16)$$

where we followed the assumption that $\inf_{g \in \mathcal{G}} R_{\text{n}}^-(g) \geq \alpha > 0$ and Theorem 2.

Secondly, we apply McDiarmid's inequality to the uniform deviation $\sup_{g \in \mathcal{G}} |\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]|$ to get that with probability at least $1 - \delta$,

$$\sup_{g \in \mathcal{G}} |\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]| - \mathbb{E}[\sup_{g \in \mathcal{G}} |\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]|] \leq C'_\delta \cdot \chi_{n_p, n_u}. \quad (17)$$

Notice that this concentration inequality is single-sided even though the uniform deviation itself is double-sided, which is different from the non-uniform deviation in Theorem 2.

Thirdly, we make *symmetrization* [50]. Suppose that $(\mathcal{X}'_p, \mathcal{X}'_u)$ is a *ghost sample*, then

$$\begin{aligned}\mathbb{E}[\sup_{g \in \mathcal{G}} |\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]|] &= \mathbb{E}_{(\mathcal{X}_p, \mathcal{X}_u)}[\sup_{g \in \mathcal{G}} |\tilde{R}_{\text{pu}}(g) - \mathbb{E}_{(\mathcal{X}'_p, \mathcal{X}'_u)}[\tilde{R}_{\text{pu}}(g)]|] \\ &\leq \mathbb{E}_{(\mathcal{X}_p, \mathcal{X}_u), (\mathcal{X}'_p, \mathcal{X}'_u)}[\sup_{g \in \mathcal{G}} |\tilde{R}_{\text{pu}}(g; \mathcal{X}_p, \mathcal{X}_u) - \tilde{R}_{\text{pu}}(g; \mathcal{X}'_p, \mathcal{X}'_u)|],\end{aligned}$$

where we applied *Jensen's inequality* twice since the absolute value and the supremum are convex. By decomposing the difference $|\tilde{R}_{\text{pu}}(g; \mathcal{X}_p, \mathcal{X}_u) - \tilde{R}_{\text{pu}}(g; \mathcal{X}'_p, \mathcal{X}'_u)|$, we can know that

$$\begin{aligned}&|\tilde{R}_{\text{pu}}(g; \mathcal{X}_p, \mathcal{X}_u) - \tilde{R}_{\text{pu}}(g; \mathcal{X}'_p, \mathcal{X}'_u)| \\ &= |\pi_p \hat{R}_p^+(g; \mathcal{X}_p) - \pi_p \hat{R}_p^+(g; \mathcal{X}'_p) \\ &\quad + \max\{0, \hat{R}_u^-(g; \mathcal{X}_u) - \pi_p \hat{R}_p^-(g; \mathcal{X}_p)\} - \max\{0, \hat{R}_u^-(g; \mathcal{X}'_u) - \pi_p \hat{R}_p^-(g; \mathcal{X}'_p)\}| \\ &\leq \pi_p |\hat{R}_p^+(g; \mathcal{X}_p) - \hat{R}_p^+(g; \mathcal{X}'_p)| + \pi_p |\hat{R}_p^-(g; \mathcal{X}_p) - \hat{R}_p^-(g; \mathcal{X}'_p)| + |\hat{R}_u^-(g; \mathcal{X}_u) - \hat{R}_u^-(g; \mathcal{X}'_u)|\end{aligned}$$

where we employed $|\max\{0, z\} - \max\{0, z'\}| \leq |z - z'|$. This decomposition results in

$$\begin{aligned}\mathbb{E}[\sup_{g \in \mathcal{G}} |\tilde{R}_{\text{pu}}(g) - \mathbb{E}[\tilde{R}_{\text{pu}}(g)]|] &\leq \pi_p \mathbb{E}_{\mathcal{X}_p, \mathcal{X}'_p}[\sup_{g \in \mathcal{G}} |\hat{R}_p^+(g; \mathcal{X}_p) - \hat{R}_p^+(g; \mathcal{X}'_p)|] \\ &\quad + \pi_p \mathbb{E}_{\mathcal{X}_p, \mathcal{X}'_p}[\sup_{g \in \mathcal{G}} |\hat{R}_p^-(g; \mathcal{X}_p) - \hat{R}_p^-(g; \mathcal{X}'_p)|] \\ &\quad + \mathbb{E}_{\mathcal{X}_u, \mathcal{X}'_u}[\sup_{g \in \mathcal{G}} |\hat{R}_u^-(g; \mathcal{X}_u) - \hat{R}_u^-(g; \mathcal{X}'_u)|].\end{aligned}\quad (18)$$

Fourthly, we relax those expectations in (18) to Rademacher complexities. The original ℓ may miss the origin, i.e., $\ell(0, y) \neq 0$, with which we need to cope. Let

$$\tilde{\ell}(t, y) = \ell(t, y) - \ell(0, y)$$

be a *shifted loss* so that $\tilde{\ell}(0, y) = 0$. Note that for all $t, t' \in \mathbb{R}$ and $y = \pm 1$,

$$\ell(t, y) - \ell(t', y) = \tilde{\ell}(t, y) - \tilde{\ell}(t', y).$$

Hence,

$$\begin{aligned}
\widehat{R}_p^+(g; \mathcal{X}_p) - \widehat{R}_p^+(g; \mathcal{X}'_p) &= (1/n_p) \sum_{x_i \in \mathcal{X}_p} \ell(g(x_i), +1) - (1/n_p) \sum_{x'_i \in \mathcal{X}'_p} \ell(g(x'_i), +1) \\
&= (1/n_p) \sum_{i=1}^{n_p} (\ell(g(x_i), +1) - \ell(g(x'_i), +1)) \\
&= (1/n_p) \sum_{i=1}^{n_p} (\tilde{\ell}(g(x_i), +1) - \tilde{\ell}(g(x'_i), +1)).
\end{aligned}$$

This is already a standard form where we can attach Rademacher variables to every $\tilde{\ell}(g(x_i), +1) - \tilde{\ell}(g(x'_i), +1)$, and it is a routine work to show that

$$\mathbb{E}_{\mathcal{X}_p, \mathcal{X}'_p} [\sup_{g \in \mathcal{G}} |\widehat{R}_p^+(g; \mathcal{X}_p) - \widehat{R}_p^+(g; \mathcal{X}'_p)|] \leq 2\mathfrak{R}_{n_p, p_p}(\tilde{\ell}(\cdot, +1) \circ \mathcal{G}).$$

The other two expectations can be handled analogously. As a result, (18) can be reduced to

$$\begin{aligned}
\mathbb{E}[\sup_{g \in \mathcal{G}} |\widetilde{R}_{pu}(g) - \mathbb{E}[\widetilde{R}_{pu}(g)]|] &\leq 2\pi_p \mathfrak{R}'_{n_p, p_p}(\tilde{\ell}(\cdot, +1) \circ \mathcal{G}) \\
&\quad + 2\pi_p \mathfrak{R}'_{n_p, p_p}(\tilde{\ell}(\cdot, -1) \circ \mathcal{G}) + 2\mathfrak{R}'_{n_u, p}(\tilde{\ell}(\cdot, -1) \circ \mathcal{G}). \quad (19)
\end{aligned}$$

Finally, we transform the Rademacher complexities of composite function classes in (19) to those of the original function class. It is obvious that $\tilde{\ell}$ shares the same Lipschitz constant L_ℓ with ℓ , and consequently

$$\begin{aligned}
\mathfrak{R}'_{n_p, p_p}(\tilde{\ell}(\cdot, +1) \circ \mathcal{G}) &\leq 2L_\ell \mathfrak{R}'_{n_p, p_p}(\mathcal{G}) = 2L_\ell \mathfrak{R}_{n_p, p_p}(\mathcal{G}) \\
\mathfrak{R}'_{n_p, p_p}(\tilde{\ell}(\cdot, -1) \circ \mathcal{G}) &\leq 2L_\ell \mathfrak{R}'_{n_p, p_p}(\mathcal{G}) = 2L_\ell \mathfrak{R}_{n_p, p_p}(\mathcal{G}) \\
\mathfrak{R}'_{n_u, p}(\tilde{\ell}(\cdot, -1) \circ \mathcal{G}) &\leq 2L_\ell \mathfrak{R}'_{n_u, p}(\mathcal{G}) = 2L_\ell \mathfrak{R}_{n_u, p}(\mathcal{G}),
\end{aligned} \quad (20)$$

where we used Talagrand's contraction lemma and the assumption that \mathcal{G} is closed under negation. Combining (16), (17), (19) and (20) finishes the proof of the uniform deviation bound (15). \square

B.5 Proof of Theorem 4

Based on Lemma 5, the estimation error bound (13) is proven through

$$\begin{aligned}
R(\tilde{g}_{pu}) - R(g^*) &= \left(\widetilde{R}_{pu}(\tilde{g}_{pu}) - \widetilde{R}_{pu}(g^*) \right) + \left(R(\tilde{g}_{pu}) - \widetilde{R}_{pu}(\tilde{g}_{pu}) \right) + \left(\widetilde{R}_{pu}(g^*) - R(g^*) \right) \\
&\leq 0 + 2 \sup_{g \in \mathcal{G}} |\widetilde{R}_{pu}(g) - R(g)| \\
&\leq 16L_\ell \pi_p \mathfrak{R}_{n_p, p_p}(\mathcal{G}) + 8L_\ell \mathfrak{R}_{n_u, p}(\mathcal{G}) + 2C'_\delta \cdot \chi_{n_p, n_u} + 2C_\ell \pi_p \Delta,
\end{aligned}$$

where $\widetilde{R}_{pu}(\tilde{g}_{pu}) \leq \widetilde{R}_{pu}(g^*)$ by the definition of \tilde{g}_{pu} . \square