

博士論文

Machine Learning with Weak Supervision from Risk Minimization Perspective

(リスク最小化の観点からの弱教師付き機械学習)

坂井 智哉

Contents

Abstract	v
Acknowledgements	vii
List of Figures	xi
List of Tables	xiv
Notation	xv
1 Introduction	1
1.1 Learning from Data	1
1.2 Machine Learning with Weak Supervision	2
1.2.1 Supervised Learning	2
1.2.2 Unsupervised Learning	3
1.2.3 Semi-Supervised Learning	4
1.2.4 Positive-Unlabeled Learning (PU Learning)	6
1.3 Contributions	7
1.3.1 Imbalanced Classification from Positive and Unlabeled Data . . .	7
1.3.2 Statistical Dependency Analysis from Positive and Unlabeled Data	7
1.3.3 Learning from Positive, Negative, and Unlabeled Data	8
1.4 Organization	8
2 Preliminaries and Related Work	11
2.1 Common Notation	11
2.2 Empirical Risk Minimization (ERM)	12
2.3 Classification in Statistical Machine Learning	13
2.3.1 Classifier	13
2.3.2 Loss Function	14
2.3.3 Performance Measure	15

2.3.4	Regularization	18
2.4	PN Classification (Supervised Classification)	19
2.5	PU Classification	20
2.5.1	Empirical Risk Minimization Approach	20
2.5.2	Theoretical Comparisons of PU Classification against Supervised Classification	24
2.6	Semi-Supervised Classification	26
2.6.1	Entropy Regularization	26
2.6.2	Manifold Regularization	26
2.6.3	Squared-Loss Mutual Information Regularization	27
2.6.4	Maximum-Margin Principle	28
2.6.5	Safe Approach	29
2.7	Class-Prior Estimation Methods	30
2.7.1	Class-Prior Shift	30
2.7.2	Class-Prior Estimation under Semi-Supervised Learning Setting .	31
2.7.3	Class-Prior Estimation under PU Learning Setting	32
2.8	Statistical Dependency	33
2.8.1	Mutual Information (MI)	33
2.8.2	Squared-Loss Mutual Information (SMI)	34
3	Imbalanced Classification from Positive and Unlabeled Data	39
3.1	Introduction	39
3.2	Problem Setting	40
3.3	PU-AUC Optimization	40
3.4	Practical Implementation	43
3.4.1	General Case	43
3.4.2	Analytical Solution for Squared Loss	44
3.4.3	Model Selection	45
3.5	Generalization Error Bounds	46
3.6	Experiments	47
3.7	Derivation of PU-AUC Risk Estimator	49
3.8	Proof of Theorem 5	51
4	Statistical Dependency Analysis from Positive and Unlabeled Data	57
4.1	Introduction	57
4.2	Problem Setting	58

4.3	SMI Estimation from Positive and Unlabeled Data	59
4.3.1	SMI with Positive and Unlabeled Data (PU-SMI)	59
4.3.2	PU-SMI Estimation	59
4.4	Practical Implementation	61
4.4.1	Neural Network and Linear-in-Parameter Models	61
4.4.2	Model Selection	62
4.5	Theoretical Analysis	63
4.6	Application based on PU-SMI	64
4.6.1	Dimension Reduction	64
4.6.2	Independence Test	65
4.7	Experiments	67
4.7.1	Experimental Analysis	67
4.7.2	Dimension Reduction	68
4.7.3	Independence Test	70
4.8	Proof of Theorem 9	71
5	Learning from Positive, Negative, and Unlabeled Data	75
5.1	Introduction	75
5.2	Problem Setting	76
5.3	Semi-Supervised Learning based on PU Learning	76
5.3.1	PUNU Learning	77
5.3.2	PNU Learning	78
5.3.3	PUNU vs. PNU Learning	80
5.4	Practical Implementation	81
5.4.1	General Case	82
5.4.2	PNU Classification	82
5.4.3	PNU-AUC Optimization	84
5.4.4	Model Selection	85
5.5	Theoretical Analyses	87
5.5.1	Generalization Error Bounds	87
5.5.2	Variance Reduction	90
5.6	Experiments on Balanced Classification	95
5.6.1	Experimental Analyses	95
5.6.2	Comparison with Existing Methods	98
5.6.3	Image Classification:	100

5.7	Experiments on Imbalanced Classification	103
5.7.1	Experimental Analysis	104
5.7.2	Comparison with Existing Methods	108
5.7.3	Text Classification	110
5.8	Proofs of Theorems	112
5.8.1	Proof of Theorem 13	113
5.8.2	Proof of Theorem 14	115
5.8.3	Proof of Theorem 15	117
5.8.4	Proofs of Theorems 16 and 17	119
5.8.5	Proof of Theorem 18	120
6	Conclusions and Future Work	123
6.1	Conclusions	123
6.2	Future Work	124
6.2.1	Semi-Supervised SMI Estimation	124
6.2.2	AUC Optimization from PU Data without Class-Priors	125
6.2.3	Multi-Class Extension of Semi-Supervised AUC Optimization . .	126
6.2.4	Essential Mechanism of Extracting Information from Data	127
	Bibliography	129

Abstract

Recent advances in the Internet have resulted in the generation of an enormous amount of data every day. By extracting useful rules and information from the data, we can utilize them for various purposes. However, since handling a large amount of data is a laborious task for humans, people want computers to conduct human-like information processing. To this end, *machine learning*, which is aimed at obtaining human-like intelligent systems, has been gathering growing attention and has achieved remarkable success in several tasks.

The success of machine learning mainly comes from supervised learning with a large amount of *labeled* data. However, a sufficient amount of labeled data is not always available in practice. To address this problem, researchers have explored a machine learning approach that can learn from *limited* or *weak supervision*. For example, semi-supervised learning, which utilizes abundant *unlabeled* data in addition to a small amount of labeled data, has been studied. However, most of the existing semi-supervised learning methods proposed so far rely on particular distributional assumptions, such as the *cluster assumption*, that are rarely satisfied in real-world applications. Once such distributional assumptions are violated, the performance of the existing methods degrades and sometimes becomes even worse than that of the plain supervised learning method using only a small amount of labeled data. Thus, it is desired to have semi-supervised learning methods that do not rely on strong distributional assumptions.

On the other hand, in real-world applications, it is conceivable that only positive and unlabeled data are available. Learning from such a situation is called *positive-unlabeled learning* (PU learning), which can also be regarded as a form of learning from weak supervision since labeled data is available only from the positive class. For example, in online ad click prediction problems, when collecting data, merely regarding an unclicked advertisement as negative data introduces a *bias*; users may have an interest in the advertisement, but may not click on the ad because of lack of time. PU learning provides a solution that cancels the bias. Since the PU learning situation often occurs in real-world applications, it is practically desirable to have PU learning methods for various machine learning tasks such as imbalanced classification and statistical dependency analysis.

To address the above issues, this dissertation is devoted to developing machine learning approaches under weak supervision. In particular, we focus on a risk minimization approach. That is, based on a learning criterion called a *risk* or *error*, machine learning methods learn decision rules and statistical dependency from data. The contributions of this dissertation are as follows: i) We propose a PU classification method for imbalanced data, where the number of positive and negative samples in underlying distributions is far from even. ii) We develop statistical dependency analysis tools for PU data. iii) We propose semi-supervised learning methods that do not rely on strong distributional assumptions. Below, we briefly explain these contributions.

Firstly, in addition to the PU learning scenario, we consider a class-imbalance situation. In this situation, maximizing the *area under the receiver operating characteristic curve* (AUC) is a standard approach rather than ordinary misclassification rate minimization. We propose a risk function based on AUC for the PU learning situation and prove that minimization of the risk estimator corresponds to minimization of the generalization error bound in AUC maximization. We experimentally demonstrate the practical effectiveness of the AUC maximization method from PU data.

Next, we develop statistical dependency analysis tools for PU data. As a statistical dependency measure, we employ *squared-loss mutual information* (SMI). We present an SMI estimation method from only PU data and prove its optimal convergence to true SMI. Then, we apply the SMI estimator to dimension reduction and independence test. Through numerical experiments, we investigate the behavior of the SMI estimator and demonstrate the effectiveness of the SMI-based statistical dependency analysis methods.

Then, we discuss a semi-supervised learning approach based on PU learning. Unlike most of the existing methods, our methods can leverage unlabeled data without the restrictive distributional assumptions such as the cluster assumption. We prove that the variance of the proposed risk estimators is smaller than that of their supervised counterparts, and generalization error bounds decrease with respect to the number of positive, negative, and unlabeled data with the optimal parametric rate without restrictive distributional assumptions. We experimentally show the effectiveness of the semi-supervised learning methods.

Finally, we summarize the conclusions of this dissertation and discuss future directions.

Acknowledgements

First of all, I am deeply grateful to my supervisor, Prof. Masashi Sugiyama. Without his support, my Ph.D. studies might have reached a dead end. Every time I was confronted with a difficulty or challenge, he encouraged me and provided valuable comments to help me overcome them. He patiently supported me as I progressed, and his unwavering confidence in me has helped relieve my anxiety while pursuing my education. I would also like to express my gratitude to Prof. Issei Sato and Prof. Junya Honda. They frankly communicated with me and helped provide refreshing ideas and advice. Moreover, I am very grateful to Prof. Masato Okada, Prof. Noboru Kunihiro, and Prof. Taiji Suzuki for reviewing and evaluating my thesis.

Next, I would like to express my great appreciation to Dr. Marthinus Christoffel du Plessis and Prof. Gang Niu. I greatly appreciated their advice, as it was based on their research experience and it opened my eyes to numerous essential points for conducting research and writing papers. Without them, my Ph.D. studies would have been almost impossible to complete.

Furthermore, my gratitude goes out all the members of Sugiyama-Sato-Honda Laboratory. In particular, I am grateful to Ikko Yamane. His insight helped not only my Ph.D. studies but also setting up all servers for numerical experiments from scratch. I am also grateful to Prof. Hiroaki Sasaki, Dr. Kiyoshi Irie, Prof. Hideko Kawakubo, Dr. Voot Tangkaratt, Inbal Horev, Hiroaki Shiino, Mina Ashizawa, Soma Yokoi, Takashi Ishida, Futoshi Futami, Weihua Hu, Ryuichi Kiryo, Masayoshi Hayashi, Ryosuke Kame-sawa, Han Bao, and Yusuke Konno. The secretaries, Ms. Ayako Tamai and Ms. Yuko Kawashima, kindly supported me. I deeply appreciate their help.

My research project was financially supported by the Japan Society for the Promotion of Science (JSPS) Fellowship Program (the grant number 15J09111). The financial support helped me in various aspects: preparing research environment, conducting research, and attending conferences. Without this support, the life of my Ph.D. studies would have been harder.

Last but not least, I would like to thank my parents and grandparents for their lasting support to me.

List of Figures

1.1	Illustration of a supervised learning situation for a binary class case. In a binary class case, one class is called <i>positive</i> class and the other is called <i>negative</i> class. Samples belonging positive and negative classes are respectively called positive and negative data, denoted by “○” and “×”.	3
1.2	Illustration of a unsupervised learning situation.	4
1.3	Illustration of a semi-supervised learning situation.	5
1.4	Illustration of a PU learning situation.	6
1.5	Organization of this Dissertation	9
2.1	Shape of Loss functions. The ramp loss ℓ_R , squared-loss ℓ_S , exponential loss ℓ_E , and double-hinge loss ℓ_{DH}	16
3.1	Average computation time on benchmark datasets when $\theta_P = 0.1$ over 50 trials. The computation time of the PU-AUC optimization method includes the class-prior estimation and the empirical risk minimization.	50
4.1	Average and standard error of the squared error of PU-SMI over 50 trials. (a) n_P is increased while $n_U = 400$ is fixed. (b) n_U is increased while $n_P = 200$ is fixed. The results show that both positive and unlabeled samples contribute to improving the estimation accuracy of SMI.	68
4.2	Illustration of linear dimension reduction. (a) Estimated subspace obtained by PUDR, where positive “○” and unlabeled samples “□” are used. (b) Estimated subspace obtained by LSDR, where positive “○” and negative samples “×” are used. The results show that from only PU data, the PUDR method gave a subspace similar to that obtained by its supervised counterpart.	69

4.3	Frequency of the type-II error (the number of times the null hypothesis is accepted when the alternative hypothesis is true) of the PU-SMI based independence test with the significance level 5% over 50 trials. (a) n_P is increased while $n_U = 400$ is fixed. (b) n_U is increased while $n_P = 100$ is fixed. In both cases, the frequency of the type-II error decreased as the number of positive/unlabeled samples increases.	71
5.1	Illustration of PUNU Learning	77
5.2	Illustration of PNU Learning	79
5.3	The average computation time over 30 trials for benchmark datasets. . .	96
5.4	Average and standard error of the ratio between the variance of empirical PNU risk and that of PN risk, $\text{Var}[\hat{R}_{\text{PNU}}^{\text{MR},\eta}(\hat{g}_{\text{PN}})] / \text{Var}[\hat{R}_{\text{PN}}^{\text{MR}}(\hat{g}_{\text{PN}})]$, as a function of the number of unlabeled samples over 100 trials. Although the variance reduction is proved for an infinite number of samples, it can be observed with a finite number of samples.	98
5.5	Average and standard error of the ratio between the misclassification rates of $\hat{g}_{\text{PN}}^{\text{PNU}}$ and $\hat{g}_{\text{PN}}^{\text{PN}}$ as a function of unlabeled samples over 1000 trials. In many cases, the ratio becomes less than 1, implying that the PNU risk is a promising alternative to the standard PN risk in validation if unlabeled data are available.	99
5.6	Average computation time over 50 trials for benchmark datasets when $n_L = 50$	100
5.7	Average computation time over 30 trials for the Places 205 dataset when $n_U = 10000$	104
5.8	Average with standard error of the ratio between the variance of the empirical PNU risk and that of the PN risk, $r = \text{Var}[\hat{R}_{\text{PNU}}^{\text{AUC},\eta}(\hat{g}_{\text{PN}})] / \text{Var}[\hat{R}_{\text{PN}}^{\text{AUC}}(\hat{g}_{\text{PN}})]$, as a function of the combination parameter η over 100 trials on the Banana dataset. The class-prior is $\theta_P = 0.1$ and the number of positive and negative samples varies as $(n_P, n_N) = (2, 8), (10, 10)$, and $(18, 2)$. Left: values of r as a function of η . Right: values for $\eta > 0$ are magnified.	105

- 5.9 Average with standard error of the ratio between the variance of the PNU-AUC risk and that of the PN-AUC risk, $r = \text{Var}[\hat{R}_{\text{PNU}}^{\text{AUC},\eta}(\hat{g}_{\text{PN}})] / \text{Var}[\hat{R}_{\text{PN}}^{\text{AUC}}(\hat{g}_{\text{PN}})]$, as a function of the combination parameter η over 100 trials on the Banana dataset. A class-prior varies as $\theta_{\text{P}} = 0.1, 0.2$, and 0.3 . Left: values of r as a function of η . Right: values for $\eta > 0$ are magnified. When $\theta_{\text{P}} = 0.1, 0.2$, the variance of the empirical PNU-AUC risk is smaller than that of the PN risk for $\eta > 0$. . . 106
- 5.10 Average with standard error of the AUC as a function of the noise ρ over 100 trials. The PNU-AUC optimization method used the noisy class-prior $\hat{\theta}_{\text{P}} = \theta_{\text{P}} + \rho$ in training. The plots show that when $\theta_{\text{P}} = 0.2$ and 0.3 , the performance of the PNU-AUC optimization method is stable even when the estimated class-prior has some noise. However, when $\theta_{\text{P}} = 0.1$, as the noise is close to $\rho = -0.09$, the performance largely decreases. Since the true class-prior is small, it is sensitive to the negative bias. 107
- 5.11 Average with standard error of the AUC as a function of the number of unlabeled data n_{U} over 20 trials. 108
- 5.12 Average computation time of each method on benchmark datasets when $n_{\text{L}} = 100$ and $\theta_{\text{P}} = 0.1$ over 50 trials. 110
- 5.13 Average computation time of each method on the text classification datasets. 112

List of Tables

- 3.1 Average and standard error of the estimated class-prior over 50 trials on benchmark datasets in PU learning setting. 48
- 3.2 Average and standard error of the AUC over 50 trials on benchmark datasets. The boldface denotes the best and comparable methods in terms of the average AUC according to the t-test at the significance level 5%. The bottom row shows the number of best/comparable cases of each method. 49
- 4.1 Mean absolute error (with standard error) between the estimated class-prior and the true value over 20 trials. KM means the class-prior estimation method based on kernel mean embedding, PCA is the principal component analysis. The boldface denotes the best and comparable approach in terms of the average absolute error according to the t-test at the significance level 5%. 70
- 5.1 Average and standard error of the misclassification rates of each method over 30 trials for benchmark datasets. The boldface denotes the best and comparable methods in terms of the average misclassification rate according to the t-test at the significance level 5%. 96
- 5.2 Average and standard error of the misclassification rates of each method over 50 trials for benchmark datasets. Boldface numbers denote the best and comparable methods in terms of average misclassifications rate according to a t-test at a significance level of 5%. The bottom row gives the number of best/comparable cases of each method. 101
- 5.3 The description of the dataset used in the image classification experiment. 102
- 5.4 Average and standard error of misclassification rates over 30 trials for the Places 205 dataset. Boldface numbers denote the best and comparable methods in terms of the average misclassification rate according to a t-test at a significance level of 5%. 103

5.5	Average and standard error of the estimated class-prior over 50 trials on benchmark datasets in semi-supervised learning setting.	109
5.6	Average and standard error of the AUC over 50 trials on benchmark datasets. The boldface denotes the best and comparable methods in terms of the average AUC according to the t-test at the significance level 5%. The bottom row shows the number of best/comparable cases of each method. SSRboost is an abbreviation for SSRankboost.	111
5.7	Average with standard error of the AUC over 20 trials on the text classification datasets. The boldface denotes the best and comparable methods in terms of the average AUC according to the t-test at the significance level 5%.	112

Notation

\mathbb{R}	The set of real numbers
$\mathbf{x} \in \mathbb{R}^d$	A d -dimensional vector of features
$\mathbf{x}^{(i)}$	An i -th dimension of a vector \mathbf{x}
y	An output or class label
$p(\mathbf{x}, y)$	The joint probability density function
$p(\mathbf{x})$	The marginal probability density of \mathbf{x}
$p(y)$	The marginal probability mass of y
$\mathcal{O}_p(\cdot)$	The probabilistic order
$\ \cdot\ _2$ or $\ \cdot\ $	The ℓ_2 norm or Euclidean norm
$\ \cdot\ _1$	The ℓ_1 norm
$I(c)$	The indicator function takes 1 if the condition c is true; otherwise 0
$\mathbf{1}_b$	The b -dimensional vector whose elements are all one
\mathbf{I}_b	The $b \times b$ identity matrix
$^\top$	Transpose of vectors and matrices
$\log(x)$	Natural logarithm of x

Chapter 1

Introduction

A recent success of machine learning such as (super) human-level performance in image classification (He et al., 2016) is based on supervised learning with a large amount of annotated data. However, such a large number of annotated data is not always available in real-world applications. This dissertation is devoted to developing machine learning algorithms without sufficient amount of supervision from humans.

1.1 Learning from Data

Since ancient Greece, human beings have wondered about their intelligence and the possibility of creating intelligence. In ancient mythology, the robot behaving like a human has already appeared. When mechanical machines made by gears emerged, some people thought that the machines could think/act like humans. Indeed, people have created intelligence with several gears, and a machine was applied to decrypting codes in Germany (Singh, 2000). Based on these studies, researchers are trying to implement human intelligence using computers composed of billions/trillions of semiconductor transistors instead of gears.

Since the 2000s, the Internet has been connecting computers worldwide via cables. Along with the popularization of the Internet, a vast amount of data is being created every day. These days, researchers are studying a way of acquiring intelligence from this data. *Artificial intelligence* is intended to mimic human intelligence; in particular, *machine learning* implements it based on probability and statistics, which enables computers to conduct human-like information processing through learning from data.

As machine learning has been successfully applied to various tasks, it has been gaining growing attention from several fields. One of the successful examples is image classification. The goal of classification is to identify category or class of data points. For example,

suppose that we are given pairs of a photo (input) and an animal name in the photo (output). With machine learning methods, computers can determine the name of an animal in a new photo. Furthermore, from the input-output data, we can analyze an input-output dependency to reveal the properties of the data by using machine learning methods. In the following section, we explain machine learning approaches to solve related problems.

1.2 Machine Learning with Weak Supervision

In this section, we introduce four machine learning approaches from the perspective of supervision.

1.2.1 Supervised Learning

In a mathematics lecture at school, we first learn the rules and method for solving a problem from a teacher. During the lecture, we may further practice solving problems with several exercises. Sometime after the lecture, we may take an exam where we are faced with new unseen problems that are related to previous exercises but not the same. If we want to obtain a high score on the exam, we need to learn the rule from the teacher properly.

From the viewpoint of machine learning, learning with full supervision (from humans), is called *supervised learning* (Hastie et al., 2009; Mohri et al., 2012; Shalev-Shwartz and Ben-David, 2014). Formally, in supervised learning, we are given pairs of input and output samples

$$\{(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n, \quad (1.1)$$

where \mathcal{X} and \mathcal{Y} are domains of the input and output, respectively (see Figure 1.1 for a binary case, i.e., $|\mathcal{Y}| = 2$). For example, the inputs are images taking animals and the outputs are names of animals in the images. By using supervised learning, a mapping function from inputs and outputs can be obtained by

$$g: \mathcal{X} \rightarrow \mathcal{Y}. \quad (1.2)$$

Then, we can answer a name of an animal on a new image by using a learned mapping function.



FIGURE 1.1: Illustration of a supervised learning situation for a binary class case. In a binary class case, one class is called *positive* class and the other is called *negative* class. Samples belonging positive and negative classes are respectively called positive and negative data, denoted by “○” and “×”.

Supervised learning is a strong approach to obtain accurate input-output relationships. Moreover, in classification problems, supervised learning was proved that it achieves optimal classification accuracy as the size of data increases. Roughly speaking, we can expect high classification accuracy if we have a large amount of labeled data. Recent advances in deep learning partly come from this fact, i.e., due to the progress of the Internet, the large size of the dataset available and it enables us to utilize for learning an accurate mapping function.

From the perspective of supervision, to execute supervised learning methods, we need to prepare annotated samples in advance, e.g., in the example of animal images classification, we need to assign a name of animals for each image. However, collecting labeled samples requires laborious annotation task and a large budget to hire specialists for annotation tasks. For each input, we need to give labels appropriately so that supervised learning methods can learn input-output relation correctly. Thus, to reduce labeling costs, researchers have studied alternative learning approach.

1.2.2 Unsupervised Learning

In contrast to supervised learning, supervision is not provided in unsupervised learning (Hastie et al., 2009; Mohri et al., 2012; Shalev-Shwartz and Ben-David, 2014). That is, the available information is inputs only as shown in Figure 1.2:

$$\{\mathbf{x}_i\}_{i=1}^n. \quad (1.3)$$

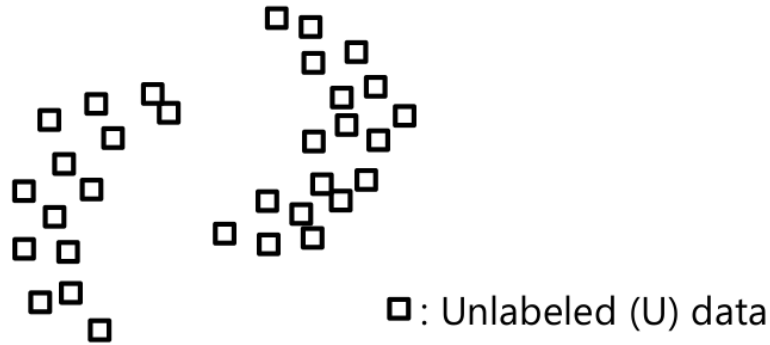


FIGURE 1.2: Illustration of a supervised learning situation.

The data without supervision are called *unlabeled data*. From obtained unlabeled data, unsupervised learning is designed to extract useful information.

For example, finding several groups from unlabeled data based on some characteristics is useful for people to understand properties of the data. This is called the *clustering*, and several algorithms were proposed. Another example is (unsupervised) *dimensionality reduction*; its goal is to obtain low-dimensional representation from (redundant) high-dimensional data for visualizing and analyzing the data without supervision.

Unsupervised learning does not require costly labeled data. In this sense, unsupervised learning is preferable to supervised learning. To extract information from unlabeled data, we often presume a particular assumption on the data, e.g., a sort of cluster structure is assumed for clustering methods. However, once the data distribution is far from the prespecified structure or assumptions, unsupervised learning methods do not work well as we expected. Moreover, since it is difficult to measure the quality of unsupervised learning, we often need some heuristic or quality measure based on a particular assumption to choose results obtained by unsupervised learning. If classification is the final goal, the labels would be precious information to achieve high performance.

1.2.3 Semi-Supervised Learning

Semi-supervised learning aims at obtaining high performance from both unlabeled and limited size of labeled data (see Figure 1.3), which can be regarded as a combination of supervised and unsupervised learning. From the perspective of supervision, supervision provided to semi-supervised learning is limited (weak) since the size of labeled data is

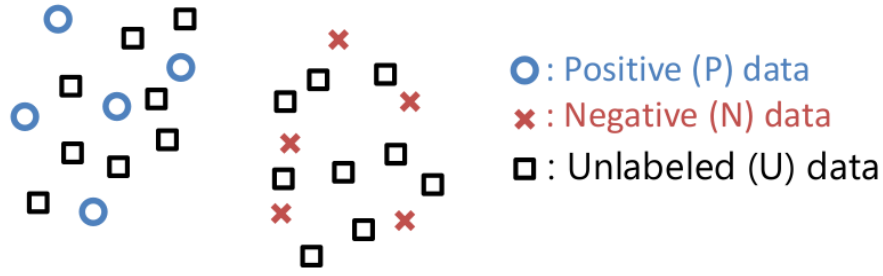


FIGURE 1.3: Illustration of a semi-supervised learning situation.

often assumed to be rather small compared with that of supervised learning. If semi-supervised learning methods can achieve high accuracy with a small amount of labeled data, we could manage both cost and accuracy.

In the beginning, most of the researchers tried to incorporate knowledge and techniques from unsupervised learning. For instance, [Chapelle et al. \(2002\)](#) proposed the *cluster kernel* based on the cluster assumption: “similar” data points tend to have the same label. Another example is the *manifold assumption*, which supposes that samples are distributed on a (near) low-dimensional manifold in the data space ([Belkin et al., 2006](#)). These assumptions can be summarized as the *low-density separation principle* that samples in different classes tend to be separated in a region with low data density ([Chapelle et al., 2006](#)). Although hundreds of studies have been devoted to developing semi-supervised learning methods based on the low-density separation principle, once the distributional assumptions are not satisfied, the performance of existing semi-supervised learning methods degrades or even worse than that of supervised learning methods with a limited amount of labeled data ([Cozman et al., 2003](#); [Sokolovska et al., 2008](#)).

To cope with this problem, recently several methods have been proposed to mitigate undesired performance degeneration of semi-supervised learning methods ([Li and Zhou, 2015](#); [Loog, 2016](#); [Krijthe and Loog, 2017](#)). For instance, a safe approach ([Li and Zhou, 2015](#)) aims at designing semi-supervised learning methods so that its performance does not become worse than the supervised learning methods. However, the method proposed by [Li and Zhou \(2015\)](#) still requires the distributional assumption for leveraging unlabeled data. Thus, if the distributional assumptions are not satisfied, the performance improvement is not expected even though the performance degeneration does not occur.

Therefore, it is practically desirable to have semi-supervised learning methods which do not strongly rely on the distributional assumptions.

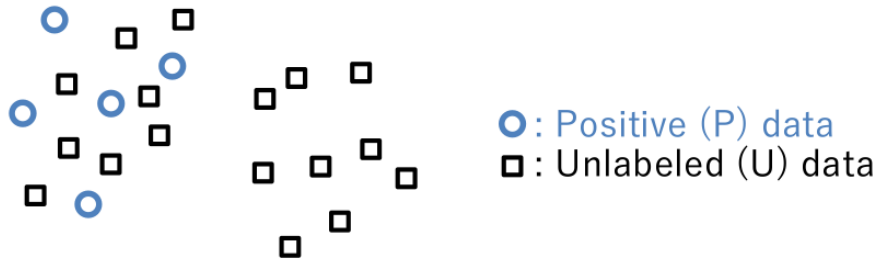


FIGURE 1.4: Illustration of a PU learning situation.

1.2.4 Positive-Unlabeled Learning (PU Learning)

In real-world applications, we are sometimes faced with the situation in which labeled data for the positive class are relatively easy to obtain but that for the negative class are difficult or almost impossible. In *positive-unlabeled learning* (PU learning) (Letouzey et al., 2000), we are only given positive and unlabeled data, but not given negative data (supervision is available only for the positive class), unlike ordinary supervised learning. Figure 1.4 illustrates a PU learning situation.

For example, suppose that we would like to classify cat photos out of animal photos. To this end, we assign cats and non-cats labels to some portion of photos. In this situation, assigning non-cats label would be more difficult than assigning cats label to the photos because we need to cover the several possibilities of non-cat photos such as dogs, lions, monkeys, ..., etc. This problem is known as the diversity of negative classes; the negative class contains several classes.

Another example is the case that labels for the negative class are hardly obtained. For instance, in recommendation system on social service, a service provider can observe a user's favorite articles through "like" button. In contrast, it is difficult to observe the user's disliked article because "dislike" button is usually not implemented. Moreover, whole articles that the user did not push the "like" button cannot be regarded merely as disliked articles since the user may not have read them yet: these articles are a mixture of liked and disliked articles. For predicting user's preference, classification problem in PU learning is needed to be solved.

One naive way of dealing with such a PU classification task is to regard unlabeled data as negative data and apply an ordinary supervised classification method to train a classifier. However, this approach causes bias because unlabeled data consist of both positive and negative data. To address this problem, several methods have been proposed to learn a decision rule from only PU data (Lee and Liu, 2003; Liu et al., 2003; Li and

Liu, 2005; Elkan and Noto, 2008). Recently, a PU learning method based on an empirical risk minimization framework was proposed (du Plessis et al., 2014, 2015) and shown that its risk function computed from only PU data is equivalent to the risk in supervised classification. Furthermore, a generalization property of PU learning was theoretically investigated, revealing that PU learning outperforms an ordinary supervised classification method under some conditions.

As the PU learning situation is perceived, it is practically desirable to have accurate classification and statistical data analysis methods for PU data.

1.3 Contributions

In this section, we describe the contributions of this dissertation.

1.3.1 Imbalanced Classification from Positive and Unlabeled Data

We develop a PU classification method for imbalanced classification, where the number of positive and negative samples in underlying distributions is far from even. For such a imbalanced classification problem, optimizing (maximizing) the *area under the receiver operating characteristic curve* (AUC) is a standard approach.

We derive the risk based on AUC optimization from only PU data (the PU-AUC risk) and reveal that the existing AUC optimization method uses a biased risk. We prove that minimization of the PU-AUC risk estimator corresponds to minimization of upper bound of the generalization error in AUC optimization. The result shows that PU data contribute to reducing the upper bound without any distributional assumptions.

We experimentally show the effectiveness of the method based on the PU-AUC risk on various datasets.

1.3.2 Statistical Dependency Analysis from Positive and Unlabeled Data

We present statistical dependency analysis tools for PU data.

To evaluate input-output dependency, we employ *squared-loss mutual information* (SMI) as a statistical dependency measure. With only PU data, we reformulate SMI (PU-SMI) that is equivalent to SMI from positive-negative data and show an estimation method. We establish a convergence rate of the PU-SMI estimator, revealing that PU data

help converge to true SMI. We then apply the PU-SMI estimator to two statistical data analysis tasks: dimension reduction and independence test. In particular, the SMI-based dimension reduction method can be executed without an estimate of class-prior that is a critical quantity for the existing PU classification method.

Through numerical experiments, we confirm that the SMI-based dimension reduction method gives a lower class-prior estimation error compared with other dimension reduction method. Additionally, we evaluate the performance of the SMI-based independence test, showing that the method can detect statistical dependency of input-output from only PU data.

1.3.3 Learning from Positive, Negative, and Unlabeled Data

We present a semi-supervised learning approach that does not rely on strong distributional assumptions.

An advantage of the semi-supervised learning approach is that it can employ any classification performance measures as long as a risk in PU learning is available, which is equivalent to that in PN learning. We show the classification methods based on the misclassification rate and AUC. Without strong distributional assumptions, we derive generalization error bounds, revealing that both labeled and unlabeled contribute to reducing the upper bound of the generalization error. Furthermore, we show the variance of the risk estimators based on PU learning becomes smaller than that of the plain supervised learning.

Through extensive experiments, we investigate the behavior of the semi-supervised learning methods and confirm its effectiveness.

1.4 Organization

This dissertation is organized into six chapters shown in Figure 1.5.

Chapter 1 introduces machine learning approach from the perspective of supervision and describes the contributions of this dissertation.

Chapter 2 provides the background of this dissertation. First of all, we review basic concepts of classification in machine learning. Based on the knowledge, we review existing methods in supervised learning, positive-unlabeled learning, and semi-supervised learning methods. Then, we introduce class-prior estimation methods used in experiments

and statistical dependency measures. The notations used in Sections 2.3 and Section 2.5 are frequently used in subsequent chapters.

In Chapter 3, we introduce a method for imbalanced classification in the PU learning setting. We show a learning criterion with a measure for imbalanced classification in a PU learning scenario. Then, we report theoretical and empirical properties of the imbalanced classification method for PU data.

Next, we discuss statistical dependency analysis for PU data in Chapter 4. We present an SMI estimation method from only PU data and apply the estimators to statistical dependency analysis tasks. We confirm the effectiveness of the SMI-based statistical dependency analysis methods.

In Chapter 5, we introduce a semi-supervised learning framework based on positive-unlabeled learning. The framework can be applied to both balanced and imbalanced classification. We investigate the theoretical property and practical effectiveness of the semi-supervised learning approach.

Finally, we present conclusions and discuss future perspective in Chapter 6.

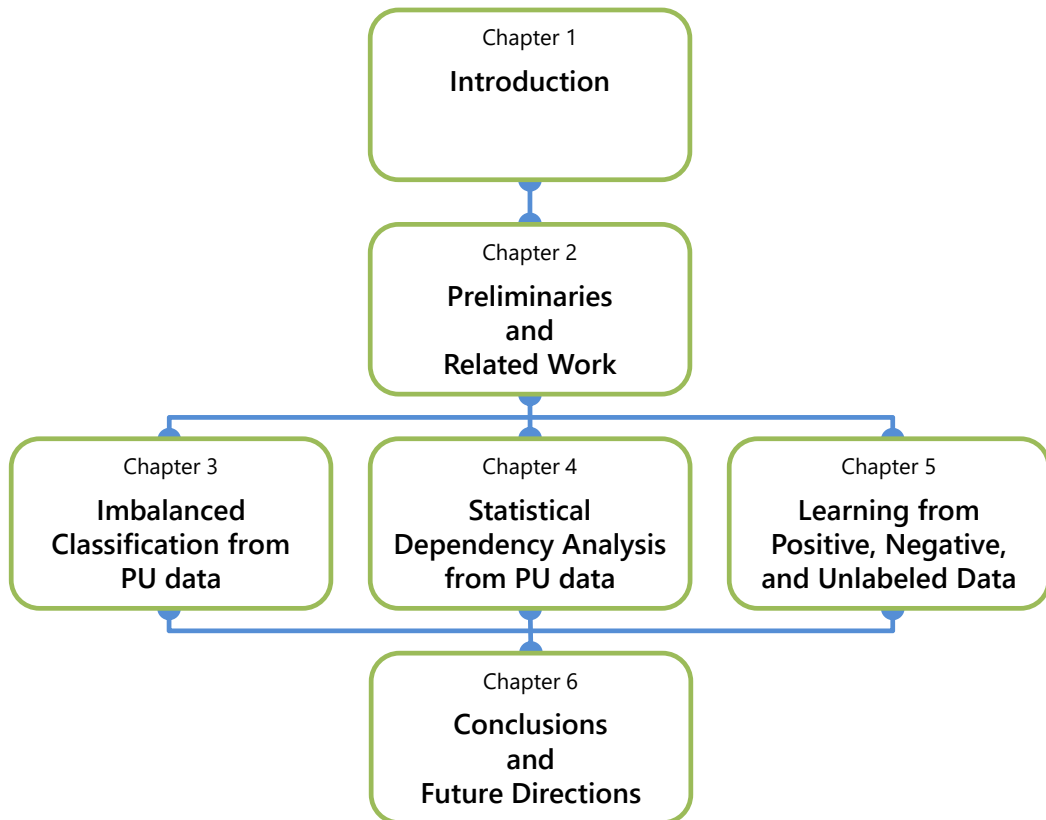


FIGURE 1.5: Organization of this Dissertation

Chapter 2

Preliminaries and Related Work

In this chapter, we review basic concepts of machine learning and related work of this dissertation. First, we describe the common notation in Section 2.1, and explain empirical risk minimization in Section 2.2. In Section 2.3, 2.4, and 2.5, we review basic components of classification, supervised classification, and PU classification. The notation introduced in these sections are frequently used in the subsequent chapters. Then, we review existing semi-supervised learning methods and class-prior estimation methods in Section 2.6 and 2.7, respectively, which are used in numerical experiments. In Section 2.8, we review statistical dependency measures.

2.1 Common Notation

Unless otherwise noted, in this chapter, we focus on binary classification tasks. Let random variables $\mathbf{x} \in \mathbb{R}^d$ and $y \in \{+1, -1\}$ be equipped with probability density $p(\mathbf{x}, y)$, where d is a positive integer.

Let us define positive (P), negative (N), and unlabeled (U) data as

$$\{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = +1), \quad (2.1)$$

$$\{\mathbf{x}_j^N\}_{j=1}^{n_N} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = -1), \quad (2.2)$$

$$\{\mathbf{x}_k^U\}_{k=1}^{n_U} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}) = \theta_P p(\mathbf{x} \mid y = +1) + \theta_N p(\mathbf{x} \mid y = -1), \quad (2.3)$$

where $\theta_P := p(y = +1)$ and $\theta_N := p(y = -1)$ are the class-prior probabilities. Moreover, let E_P , E_N , and E_U be the expectations over $p(\mathbf{x} \mid y = +1)$, $p(\mathbf{x} \mid y = -1)$, and $p(\mathbf{x})$, respectively.

We use this notation to review the basic concepts of machine learning and existing methods in the subsequent sections.

2.2 Empirical Risk Minimization (ERM)

Here, we explain a learning approach based on minimizing a *risk* or *error*, called *risk minimization*. In classification tasks, a possible choice for the risk is the ratio between the number of data points misclassified by a model and the total number of data points, where the model is assumed for approximating a target function such as a discrimination function.¹ A small risk means that a model is somehow close to the *true* or *optimal* target function which achieves the minimum of the risk.

Let $g: \mathbb{R}^d \rightarrow \mathbb{R}$ be a specific model and $R: \mathbb{R} \rightarrow \mathbb{R}$ be a risk function in a certain task.² In the risk function, the quality of the model is evaluated as averaged or expected performance, the expectation is computed by the underlying distribution. Then, the learned model is obtained by minimizing the risk function:

$$g^* := \operatorname{argmin}_g R(g), \quad (2.4)$$

where g^* is called the *true* or *optimal* model. However, since the underlying distribution is not accessible in practice, we approximate the risk function by using samples collected from the distribution. Let $\hat{R}(g)$ be the *empirical* risk function that is an approximation of the true risk $R(g)$. Practically, we obtain the learned model by minimizing the empirical risk:

$$\hat{g} := \operatorname{argmin}_g \hat{R}(g), \quad (2.5)$$

where \hat{g} is the estimated model. This procedure is called *empirical risk minimization* (ERM).

In this dissertation, we design a risk function for a specific machine learning task such as classification from PU data, and approximate the risk function by using obtained samples. In the PU learning scenario, the risk function can be approximated by P and U data, and in the semi-supervised learning scenario, it can be approximated by P, N, and U data.

In the next section, we see how to apply ERM to classification tasks.

¹ Another example is a sum of squared errors of outputs between the model and the target function.

² For the sake of simplicity, we focus on $g: \mathbb{R}^d \rightarrow \mathbb{R}$ and $R: \mathbb{R} \rightarrow \mathbb{R}$, but the definition can be extended to more general settings, e.g., the domain of g can be a set and the range of g can be categorical.

2.3 Classification in Statistical Machine Learning

In this section, we review basic components in classification. There are three main components: a function to classify a data point, a measure to evaluate the quality of the function, and a criterion to train the function.

2.3.1 Classifier

A classifier plays a central role in classification. Mathematically, a classifier is a function mapping d -dimensional vector $\mathbf{x} \in \mathbb{R}^d$ (input) to a scalar (output):

$$g: \mathbb{R}^d \rightarrow \mathbb{R}. \quad (2.6)$$

Under a specific rule, we classify data points into each class by outputs of the classifier. In binary classification, we often use a sign of the classifier $\text{sign}(g(\mathbf{x}))$ for the decision rule. Below, we introduce typical classifiers used in machine learning.

The linear-in-input classifier is a simple and widely used classifier defined as

$$g(\mathbf{x}) = \sum_{\ell=1}^d v_{\ell} x^{(\ell)} + v_0 = \mathbf{v}^{\top} \bar{\mathbf{x}}, \quad (2.7)$$

where $\mathbf{v} := (v_0, v_1, \dots, v_d)^{\top}$ and $\bar{\mathbf{x}} := (1, x^{(1)}, \dots, x^{(d)})^{\top}$. The linear-in-input classifier is assumed that a hyperplane can separate most of data points. Thus, it is useful when a decision boundary of data points is almost linear, which is likely to happen in high-dimensional data. However, the data points are not always separated linearly. In such a case, applying a non-linear transformation to data points, the data can be classified more correctly in transformed space.

Let $\phi(\mathbf{x})$ be a basis function, a specific non-linear transformation for inputs. The linear-in-parameter classifier is defined as

$$g(\mathbf{x}) = \sum_{\ell=1}^b w_{\ell} \phi_{\ell}(\mathbf{x}) = \mathbf{w}^{\top} \boldsymbol{\phi}(\mathbf{x}), \quad (2.8)$$

where $\mathbf{w} := (w_1, \dots, w_b)^{\top}$ is the vector of the parameters, $\boldsymbol{\phi}(\mathbf{x}) := (\phi_1(\mathbf{x}), \dots, \phi_b(\mathbf{x}))^{\top}$ is the vector of the basis functions, and b is the number of basis functions. By setting $b = n$, where n is the number of training samples, and using the *kernel* function $K(\mathbf{x}, \mathbf{x}')$,

the linear-in-parameter classifier can be regarded as the *kernel classifier*:

$$g(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i), \quad (2.9)$$

where $\{\alpha_i\}_{i=1}^n$ is a set of parameters. For the kernel function, a popular choice is to use the *Gaussian kernel* function

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right), \quad (2.10)$$

where $\|\cdot\|$ denotes the Euclidean distance and σ is the Gaussian width. The Gaussian kernel can be interpreted as that it maps data points to infinite dimensional space ([Schölkopf and Smola, 2001](#)). In the infinite dimensional space, transformed data points can be separated easily.

2.3.2 Loss Function

To evaluate the quality of a classifier, a loss function also plays a vital role in machine learning. So far, many loss functions have been proposed and studied from several aspects such as tasks and algorithms. Here, we briefly review loss functions used in classification.

In binary classification, if a pattern \mathbf{x} with label $y \in \{\pm 1\}$ is correctly classified by a classifier g , it satisfies $yg(\mathbf{x}) > 0$; otherwise $yg(\mathbf{x}) \leq 0$. To evaluate the quality of the classifier, the loss functions map $yg(\mathbf{x})$ to some value

$$\ell: \mathbb{R} \rightarrow \mathbb{R}. \quad (2.11)$$

The *zero-one* loss function is the most basic loss function in classification tasks and can be defined as

$$\ell_{0-1}(m) := \begin{cases} 0 & (m > 0), \\ 1 & (m \leq 0). \end{cases} \quad (2.12)$$

If a pattern is correctly classified, i.e., $yg(\mathbf{x}) > 0$, the zero-one loss function does not give any loss; otherwise, it gives 1 as the loss. This corresponds to checking whether a pattern is correctly classified or not. Thus, an accurate classifier can be found by searching a classifier that can correctly classify a pattern in terms of the zero-one loss function.

Although obtaining a classifier in terms of the zero-one loss function is our ultimate goal, due to the discrete nature of the zero-one loss function, we use *surrogate* loss functions, in practice. So far, several surrogate loss functions have been proposed for classification. One famous loss function is the *hinge loss* function, which is used for the *support vector machine/classifier* (Vapnik, 1995; Cortes and Vapnik, 1995), defined as

$$\ell_H(m) := \max(0, 1 - m). \quad (2.13)$$

In statistics, the *logistic loss* function,

$$\ell_L(m) := \log(1 + \exp(-m)), \quad (2.14)$$

which is used for the *logistic regression* (Hastie et al., 2009), has been used for classification. The squared loss is also often used:

$$\ell_S(m) := (1 - m)^2. \quad (2.15)$$

In addition to the above, there are many other loss functions listed below:

- Squared hinge loss: $\ell_{SH}(m) := (\max(0, 1 - m))^2$,
- Ramp loss: $\ell_R(m) := \min(2, \max(0, 1 - m))$,
- Absolute loss: $\ell_A(m) := |1 - m|$,
- Exponential loss: $\ell_E(m) := \exp(-m)$,
- Double hinge loss: $\ell_{DH}(m) := \max(-2m, \max(0, 1 - m))$,
- Sigmoid loss: $\ell_{SG}(m) := \frac{1}{1 + \exp(-m)}$.

These loss functions can be chosen from the viewpoint of tasks, algorithms, and a property of loss functions.

2.3.3 Performance Measure

Now, we have classifiers and loss functions. The final component is performance measure for classification.

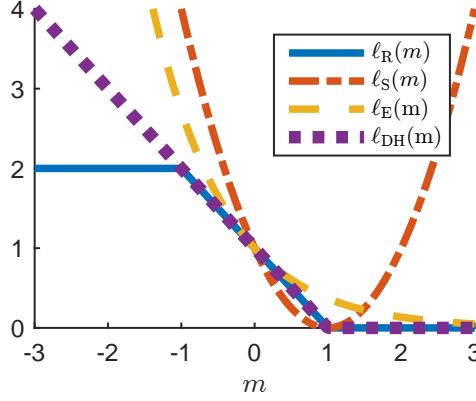


FIGURE 2.1: Shape of Loss functions. The ramp loss ℓ_R , squared-loss ℓ_S , exponential loss ℓ_E , and double-hinge loss ℓ_{DH} .

Misclassification Rate: The misclassification rate is defined as

$$\text{MR}(g) := \sum_{y=\pm 1} \int I(yg(\mathbf{x}) \leq 0) p(\mathbf{x}, y) d\mathbf{x} = \mathbb{E}_{p(\mathbf{x}, y)}[I(yg(\mathbf{x}) \leq 0)], \quad (2.16)$$

where $\mathbb{E}_{p(\mathbf{x}, y)}$ is the expectation over joint distribution, and $I(c)$ is the *indicator* function taking 1 if c is true; otherwise 0. As its name indicates, the misclassification rate computes the rate of the number of misclassified samples. To obtain an accurate classifier regarding the misclassification rate, we train a classifier by minimizing the following risk (the MR risk):

$$R^{\text{MR}}(g) = \mathbb{E}_{p(\mathbf{x}, y)}[\ell(yg(\mathbf{x}))]. \quad (2.17)$$

If we adopt the zero-one loss function for evaluating the MR risk, we recover the original misclassification rate. However, as mentioned, the surrogate loss functions are used for ease of optimization.

Since the underlying *true* distribution $p(\mathbf{x}, y)$ is unknown in practice, with *independently and identically distributed* (i.i.d.) *paired* samples drawn from the joint distribution

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y), \quad (2.18)$$

where n is the number of samples, we compute an *empirical* MR risk given by

$$\hat{R}^{\text{MR}}(g) := \frac{1}{n} \sum_{i=1}^n \ell(y_i g(\mathbf{x}_i)). \quad (2.19)$$

Based on empirical risk minimization (ERM), we find a classifier that minimizes the empirical risk.

Let

$$I^{\text{MR}}(g) := \mathbb{E}_{p(\mathbf{x}, y)}[\ell_{0-1}(yg(\mathbf{x}))] \quad (2.20)$$

be the MR risk with the zero-one loss function over the joint distribution, called the *generalization error* in misclassification rate minimization. Our ultimate goal is to obtain the classifier that minimizes the generalization error, but there is a gap since we use the surrogate loss functions in practice. Although the shape of loss functions is different, [Bartlett et al. \(2006\)](#) elucidated that minimization of the MR risk with a loss function belonging the *classification-calibrated loss* is proved to achieve that with the zero-one loss function asymptotically. For example, the hinge loss is the classification calibrated. Thus, minimization of the MR risk with the hinge loss asymptotically achieves minimization of the generalization error even though the hinge loss is a surrogate of the zero-one loss function.

Area under the Receiver Operating Characteristic Curve (AUC): The misclassification rate has been widely used and studied so far. However, the misclassification rate is not a general metric for classification performance. When the two classes are imbalanced, where the class-prior θ_P is far from 0.5, the use of the *area under the receiver operating characteristic curve* (AUC) ([Hanley and McNeil, 1982](#)) is a preferable approach ([Cortes and Mohri, 2003](#)). AUC is defined as

$$\text{AUC}(g) := \mathbb{E}_P[\mathbb{E}_N[I(g(\mathbf{x}^P) \geq g(\mathbf{x}^N))]], \quad (2.21)$$

where $\mathbb{E}_P[g(\mathbf{x})] := \int g(\mathbf{x})p(\mathbf{x} \mid y = +1)d\mathbf{x}$, $\mathbb{E}_N[g(\mathbf{x})] := \int g(\mathbf{x})p(\mathbf{x} \mid y = -1)d\mathbf{x}$, and \mathbf{x}^P and \mathbf{x}^N are random variables drawn from $p(\mathbf{x} \mid y = +1)$ and $p(\mathbf{x} \mid y = -1)$, respectively.

Unlike the misclassification rate, AUC takes into account the order of estimates of a classifier. This property is, in turn, the important for imbalanced classification. For example, when the class-prior is 0.1,³ the classifier can achieve the higher classification

³ 10% of samples drawn from the distribution $p(\mathbf{x}, y)$ tends to be positive samples and the rest 90% tends to be negative samples.

accuracy 90% by simply classifying all test samples into negative class. However, the classifier cannot necessarily achieve higher AUC.

Based on the above principle, a classifier trained by maximizing AUC has been explored. Firstly, AUC can be transformed as follows:

$$\begin{aligned} \text{AUC}(g) &= \mathbb{E}_P[\mathbb{E}_N[I(g(\mathbf{x}^P) \geq g(\mathbf{x}^N))]] \\ &= 1 - \mathbb{E}_P[\mathbb{E}_N[I(g(\mathbf{x}^P) < g(\mathbf{x}^N))]]. \end{aligned} \quad (2.22)$$

Then, the risk for AUC maximization (the AUC risk) is defined as

$$R^{\text{AUC}}(g) := \mathbb{E}_P[\mathbb{E}_N[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^N))]]. \quad (2.23)$$

If the zero-one loss is used, the risk coincides with $1 - \text{AUC}(g)$. Maximization of AUC corresponds to minimization of the AUC risk R^{AUC} . As well as misclassification rate minimization, a classifier can be trained based on AUC through the ERM framework.

Similarly to the classification-calibrated loss (Bartlett et al., 2006), the consistency of ERM with the surrogate loss in AUC optimization has been studied (Gao and Zhou, 2015; Gao et al., 2016), i.e., whether the AUC risk minimization with a specific surrogate loss function is asymptotically consistent with that with the zero-one loss. It revealed that the AUC risk minimization with the squared loss, exponential loss, and logistic loss functions shown to be consistent, while the hinge loss and absolute loss functions are *not* consistent.

2.3.4 Regularization

As discussed in the previous section, the empirical risk minimization with an appropriate loss function is asymptotically equivalent to the true risk minimization. However, the size of data is finite in practice, applying only ERM leads to a classifier fitted to the observed data completely, but it does not work well on unseen test data, which is known as *overfitting*. To avoid overfitting, several *regularization* techniques has been investigated (Hastie et al., 2009; Shalev-Shwartz and Ben-David, 2014).

The ℓ_2 -regularizer is a standard choice. For a linear-in-parameter model $g(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$, the ℓ_2 -regularizer is defined as

$$W_2(\mathbf{w}) := \sum_{\ell=1}^b w_\ell^2 = \|\mathbf{w}\|_2^2. \quad (2.24)$$

The ℓ_2 -regularizer controls complexity of the model in terms of the ℓ_2 -norm of the parameter vector.

The ℓ_1 -regularizer is also widely used in machine learning, which forces most of the elements of parameter vector to be zero, and leads to a *sparse* parameter vector. The ℓ_1 -regularizer is defined as

$$W_1(\mathbf{w}) := \sum_{\ell=1}^b |w_\ell| = \|\mathbf{w}\|_1. \quad (2.25)$$

With the regularizational functional $W(g)$, we obtain the trained classifier by solving the following optimization problem:

$$\hat{g} := \operatorname{argmin}_g \left[\hat{R}(g) + W(g) \right]. \quad (2.26)$$

2.4 PN Classification (Supervised Classification)

To solve binary or positive-negative (PN) classification problems, supervised classification can be applied with positive and negative (labeled) samples.

Recall that the MR risk is defined as

$$R^{\text{MR}}(g) = \mathbb{E}_{p(\mathbf{x}, y)} [\ell(yg(\mathbf{x}))]. \quad (2.27)$$

Let us define the MR risk in PN classification (the PN-MR risk) by decomposing the MR risk in terms of class:

$$R^{\text{MR}}(g) = \theta_P \mathbb{E}_P[\ell(g(\mathbf{x}))] + \theta_N \mathbb{E}_N[\ell(-g(\mathbf{x}))] := R_{\text{PN}}^{\text{MR}}(g). \quad (2.28)$$

The goal of PN classification is to obtain a classifier that minimizes the PN-MR risk as small as it can. In ERM, we try to accomplish such a classifier by using samples collected from the distributions.

Suppose we are given positive and negative samples drawn from class-conditional distributions:

$$\{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{i.i.d.}{\sim} p(\mathbf{x} \mid y = +1), \quad (2.29)$$

$$\{\mathbf{x}_j^N\}_{j=1}^{n_N} \stackrel{i.i.d.}{\sim} p(\mathbf{x} \mid y = -1). \quad (2.30)$$

With these samples, the empirical PN-MR risk can be obtained by

$$\hat{R}_{\text{PN}}^{\text{MR}}(g) := \frac{\theta_{\text{P}}}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \ell(g(\mathbf{x}_i^{\text{P}})) + \frac{\theta_{\text{N}}}{n_{\text{N}}} \sum_{j=1}^{n_{\text{N}}} \ell(-g(\mathbf{x}_j^{\text{N}})). \quad (2.31)$$

We obtain a classifier by minimizing the empirical PN risk with, e.g, the ℓ_2 -regularizer.

2.5 PU Classification

In this section, we review an existing PU learning method based on empirical risk minimization and theoretical comparisons of PU classification against PN classification.

2.5.1 Empirical Risk Minimization Approach

Recently, [du Plessis et al. \(2014, 2015\)](#) proposed methods based on an empirical risk minimization framework for PU classification. Their idea is to derive a risk equivalent to the PN-MR risk from only PU data without any distributional assumptions on the data. Note that they focused on the misclassification rate as a classification performance measure.

Suppose that we are given *positive* (P) and *unlabeled* (U) data:

$$\{\mathbf{x}_i^{\text{P}}\}_{i=1}^{n_{\text{P}}} \stackrel{i.i.d.}{\sim} p(\mathbf{x} \mid y = +1), \quad (2.32)$$

$$\{\mathbf{x}_k^{\text{U}}\}_{k=1}^{n_{\text{U}}} \stackrel{i.i.d.}{\sim} p(\mathbf{x}) := \theta_{\text{P}} p(\mathbf{x} \mid y = +1) + \theta_{\text{N}} p(\mathbf{x} \mid y = -1), \quad (2.33)$$

where $\theta_{\text{P}} := p(y = +1)$ and $\theta_{\text{N}} := p(y = -1)$ are the class-prior probabilities. By definition of the marginal density, we obtain

$$\mathbb{E}_{\text{U}}[\ell(-g(\mathbf{x}))] = \theta_{\text{P}} \mathbb{E}_{\text{P}}[\ell(-g(\mathbf{x}))] + \theta_{\text{N}} \mathbb{E}_{\text{N}}[\ell(-g(\mathbf{x}))]. \quad (2.34)$$

Equivalently,

$$\theta_{\text{N}} \mathbb{E}_{\text{N}}[\ell(-g(\mathbf{x}))] = \mathbb{E}_{\text{U}}[\ell(-g(\mathbf{x}))] - \theta_{\text{P}} \mathbb{E}_{\text{P}}[\ell(-g(\mathbf{x}))]. \quad (2.35)$$

Recall the PN-MR risk

$$R_{\text{PN}}^{\text{MR}}(g) = \theta_{\text{P}} \mathbb{E}_{\text{P}}[\ell(g(\mathbf{x}))] + \theta_{\text{N}} \mathbb{E}_{\text{N}}[\ell(-g(\mathbf{x}))]. \quad (2.36)$$

By plugging Eq. (2.35) into the second term of the PN-MR risk, we obtain the risk in PU classification (the PU-MR risk):

$$\begin{aligned} R_{\text{PN}}^{\text{MR}}(g) &= \theta_{\text{P}} \mathbb{E}_{\text{P}}[\ell(g(\mathbf{x})) - \ell(-g(\mathbf{x}))] + \mathbb{E}_{\text{U}}[\ell(-g(\mathbf{x}))] \\ &= \theta_{\text{P}} \mathbb{E}_{\text{P}}[\tilde{\ell}(g(\mathbf{x}))] + \mathbb{E}_{\text{U}}[\ell(-g(\mathbf{x}))] := R_{\text{PU}}^{\text{MR}}(g), \end{aligned} \quad (2.37)$$

where

$$\tilde{\ell}(m) := \ell(m) - \ell(-m) \quad (2.38)$$

is a composite loss function. Thus, without distributional assumptions such as the cluster assumption, we have the following lemma:

Lemma 1. *The MR, PN-MR, and PU-MR risks are equivalent:*

$$R^{\text{MR}}(g) = R_{\text{PN}}^{\text{MR}}(g) = R_{\text{PU}}^{\text{MR}}(g). \quad (2.39)$$

The above lemma also shows that the minimizers of the PN-MR and PU-MR risks are equivalent even though their expressions are different. In contrast to the PN-MR risk, the PU-MR risk can be computed by the expectations over P and U data. Therefore, we can train a classifier by using only PU data. In practice, we use the empirical PU-MR risk:

$$\hat{R}_{\text{PU}}^{\text{MR}}(g) := \frac{\theta_{\text{P}}}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \tilde{\ell}(g(\mathbf{x}_i^{\text{P}})) + \frac{1}{n_{\text{U}}} \sum_{k=1}^{n_{\text{U}}} \ell(-g(\mathbf{x}_k^{\text{U}})), \quad (2.40)$$

where the expectations are replaced with the corresponding sample averages. The empirical PU-MR risk converges to the PU-MR risk, i.e., the MR risk in the order of $\mathcal{O}_p(1/\sqrt{n_{\text{P}}} + 1/\sqrt{n_{\text{U}}})$, which is the optimal convergence rate without any additional distributional assumptions (Vapnik, 1998).

Non-Convex Approach: If the loss function satisfies

$$\ell(m) + \ell(-m) = 1, \quad (2.41)$$

the composite loss function becomes $\tilde{\ell}(m) = 2\ell(m) - 1$. We thus obtain the *non-convex* PU-MR risk and the empirical risk as

$$R_{\text{N-PU}}^{\text{MR}}(g) := 2\theta_{\text{P}} \mathbb{E}_{\text{P}}[\ell(g(\mathbf{x}))] + \mathbb{E}_{\text{U}}[\ell(-g(\mathbf{x}))] - \theta_{\text{P}}, \quad (2.42)$$

$$\hat{R}_{\text{N-PU}}^{\text{MR}}(g) := \frac{2\theta_{\text{P}}}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \ell(g(\mathbf{x}_i^{\text{P}})) + \frac{1}{n_{\text{U}}} \sum_{k=1}^{n_{\text{U}}} \ell(-g(\mathbf{x}_k^{\text{U}})) - \theta_{\text{P}}. \quad (2.43)$$

This formulation can be seen as cost-sensitive classification of P and U data with weight $2\theta_{\text{P}}$ (du Plessis et al., 2014).

The ramp loss used in the robust support vector machine (Collobert et al., 2006),

$$\ell_{\text{R}}(m) := \frac{1}{2} \max(0, \min(2, 1 - m)), \quad (2.44)$$

satisfies the condition (2.41). However, the use of the ramp loss (and any other losses that satisfy the condition (2.41)) yields a non-convex optimization problem, which may be solved locally by the *concave-convex procedure* (CCCP) (Yuille and Rangarajan, 2002; Collobert et al., 2006; du Plessis et al., 2014).

Convex Approach: If a convex surrogate loss function satisfies

$$\ell(m) - \ell(-m) = -m, \quad (2.45)$$

the composite loss function becomes a linear function $\tilde{\ell}(m) = -m$ (du Plessis et al., 2015). For example, the logistic loss $\ell_{\text{L}}(m) = \log(1 + \exp(-m))$, double-hinge loss $\ell_{\text{DH}}(m) = \max(-2m, \max(0, 1 - m))$, scaled squared loss $(m - 1)^2/4$ satisfy the condition. We thus obtain the *convex* PU-MR risk and the empirical risk as

$$R_{\text{C-PU}}^{\text{MR}}(g) := \theta_{\text{P}} \mathbb{E}_{\text{P}}[-g(\mathbf{x})] + \mathbb{E}_{\text{U}}[\ell(-g(\mathbf{x}))], \quad (2.46)$$

$$\hat{R}_{\text{C-PU}}^{\text{MR}}(g) := -\frac{\theta_{\text{P}}}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} g(\mathbf{x}_i^{\text{P}}) + \frac{1}{n_{\text{U}}} \sum_{k=1}^{n_{\text{U}}} \ell(-g(\mathbf{x}_k^{\text{U}})). \quad (2.47)$$

This formulation yields the convex optimization problem that can be solved efficiently.

NU Classification: As a mirror of PU classification, we can consider NU classification. The risk in NU classification (the NU-MR risk) and its empirical version are given by

$$R_{\text{NU}}^{\text{MR}}(g) := \theta_{\text{N}} \mathbb{E}_{\text{N}}[\tilde{\ell}(-g(\mathbf{x}))] + \mathbb{E}_{\text{U}}[\ell(g(\mathbf{x}))], \quad (2.48)$$

$$\hat{R}_{\text{NU}}^{\text{MR}}(g) := \frac{\theta_{\text{N}}}{n_{\text{N}}} \sum_{j=1}^{n_{\text{N}}} \tilde{\ell}(-g(\mathbf{x}_j^{\text{N}})) + \frac{1}{n_{\text{U}}} \sum_{k=1}^{n_{\text{U}}} \ell(g(\mathbf{x}_k^{\text{U}})). \quad (2.49)$$

NU classification is used in theoretical analysis in the next section and semi-supervised learning methods in Chapter 5.

Here, without distributional assumptions, we have the following lemma:

Lemma 2. *The MR, PN-MR, PU-MR, and NU-MR risks are equivalent:*

$$R^{\text{MR}}(g) = R_{\text{PN}}^{\text{MR}}(g) = R_{\text{PN}}^{\text{PU}}(g) = R_{\text{PN}}^{\text{NU}}(g). \quad (2.50)$$

In contrast to the PN-MR risk, the NU-MR risk can be computed by the expectations over N and U data. Therefore, we can train a classifier by using only NU data. Additionally, Lemma 2 shows that the PU-MR and NU-MR risks are equivalent. However, the practical behavior is different. Indeed, in a theoretical analysis in the subsequent section, we compare PN, PU, and NU classification methods and see that theoretically expected performances are different depending on the number of P, N, and U samples and the class-prior θ_{P} .

Similarly to PU classification, the non-convex and convex NU-MR risks are expressed as

$$R_{\text{N-NU}}^{\text{MR}}(g) := 2\theta_{\text{N}} \mathbb{E}_{\text{N}}[\ell(-g(\mathbf{x}))] + \mathbb{E}_{\text{U}}[\ell(g(\mathbf{x}))] - \theta_{\text{N}}, \quad (2.51)$$

$$R_{\text{C-NU}}^{\text{MR}}(g) := \theta_{\text{N}} \mathbb{E}_{\text{N}}[g(\mathbf{x})] + \mathbb{E}_{\text{U}}[\ell(g(\mathbf{x}))]. \quad (2.52)$$

In practice, the empirical versions are used for training a classifier:

$$\hat{R}_{\text{N-NU}}^{\text{MR}}(g) := \frac{2\theta_{\text{N}}}{n_{\text{N}}} \sum_{j=1}^{n_{\text{N}}} \ell(-g(\mathbf{x}_j^{\text{N}})) + \frac{1}{n_{\text{U}}} \sum_{k=1}^{n_{\text{U}}} \ell(g(\mathbf{x}_k^{\text{U}})) - \theta_{\text{N}}, \quad (2.53)$$

$$\hat{R}_{\text{C-NU}}^{\text{MR}}(g) := \frac{\theta_{\text{N}}}{n_{\text{N}}} \sum_{j=1}^{n_{\text{N}}} g(\mathbf{x}_j^{\text{N}}) + \frac{1}{n_{\text{U}}} \sum_{k=1}^{n_{\text{U}}} \ell(g(\mathbf{x}_k^{\text{U}})). \quad (2.54)$$

2.5.2 Theoretical Comparisons of PU Classification against Supervised Classification

du Plessis et al. (2014, 2015) proposed the MR risk in PU classification equivalent to the MR risk in PN (supervised) classification. Even though those two risks are equivalent, from a qualitative point of view, those risks utilize different data. One may ask a question that which learning approach, PU or PN classification, is better? To answer this question, Niu et al. (2016) theoretically compared the classification performance between PU and PN classification regarding estimation error bounds.

We assume that \mathcal{G} is a set of functions for the classifiers defined as

$$\mathcal{G} := \{g(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}) \mid \|\mathbf{w}\| \leq C_w; \forall \mathbf{x}: \|\boldsymbol{\phi}(\mathbf{x})\| \leq C_\phi\}, \quad (2.55)$$

where $C_w > 0$ and $C_\phi > 0$ are certain positive constants. This assumption is reasonable because the ℓ_2 -regularizer included in training and the use of bounded basis functions, e.g., the Gaussian kernel basis, ensures that the minimizer of the empirical risk belongs to such the function class \mathcal{G} . Let us define the PN, PU, and NU classifiers as

$$\hat{g}_{\text{PN}} := \operatorname{argmin}_{g \in \mathcal{G}} \hat{R}_{\text{PN}}^{\text{MR}}(g), \quad (2.56)$$

$$\hat{g}_{\text{PU}} := \operatorname{argmin}_{g \in \mathcal{G}} \hat{R}_{\text{PU}}^{\text{MR}}(g), \quad (2.57)$$

$$\hat{g}_{\text{NU}} := \operatorname{argmin}_{g \in \mathcal{G}} \hat{R}_{\text{NU}}^{\text{MR}}(g), \quad (2.58)$$

respectively.

The estimation error bounds can be obtained separately with probability at least $1 - \delta$:

$$R^{\text{MR}}(\hat{g}_{\text{PN}}) - R^{\text{MR}}(g^*) \leq h^{\text{MR}}(\delta) \left(\frac{\theta_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{\theta_{\text{N}}}{\sqrt{n_{\text{N}}}} \right), \quad (2.59)$$

$$R^{\text{MR}}(\hat{g}_{\text{PU}}) - R^{\text{MR}}(g^*) \leq h^{\text{MR}}(\delta) \left(\frac{2\theta_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}} \right), \quad (2.60)$$

$$R^{\text{MR}}(\hat{g}_{\text{NU}}) - R^{\text{MR}}(g^*) \leq h^{\text{MR}}(\delta) \left(\frac{2\theta_{\text{N}}}{\sqrt{n_{\text{N}}}} + \frac{1}{\sqrt{n_{\text{U}}}} \right), \quad (2.61)$$

where $g^* := \operatorname{argmin}_{g \in \mathcal{G}} R^{\text{MR}}(g)$ is the optimal decision function and $h^{\text{MR}}(\delta) := 4L_\ell C_w C_\phi + \sqrt{2 \log(4/\delta)}$.

We then compare the bounds between PN and PU (NU) classification. Since $h(\delta)\theta_{\text{P}}/\sqrt{n_{\text{P}}}$ (resp., $h(\delta)\theta_{\text{N}}/\sqrt{n_{\text{N}}}$) appears both of Eq. (2.59) and Eq. (2.60) (resp.,

Eq. 2.61), we subtract the term from both bounds and compute the ratio of the two bounds:

$$\rho_{\text{PU},\text{PN}} := \frac{\theta_{\text{P}}/\sqrt{n_{\text{P}}} + 1/\sqrt{n_{\text{U}}}}{\theta_{\text{N}}/\sqrt{n_{\text{N}}}}, \quad (2.62)$$

$$\rho_{\text{NU},\text{PN}} := \frac{\theta_{\text{N}}/\sqrt{n_{\text{N}}} + 1/\sqrt{n_{\text{U}}}}{\theta_{\text{P}}/\sqrt{n_{\text{P}}}}. \quad (2.63)$$

When $\rho_{\text{PU},\text{PN}} < 1$, PU classification is more promising than PN classification; on the other hand, $\rho_{\text{PU},\text{PN}} > 1$ means that PN classification is more promising than PU classification. Similarly, we can compare NU and PN classification by $\rho_{\text{NU},\text{PN}}$.

Furthermore, when vast amounts of unlabeled data are available ($n_{\text{U}} \rightarrow \infty$ which is faster than $n_{\text{P}} \rightarrow \infty$ and $n_{\text{N}} \rightarrow \infty$), we have the results of *asymptotic comparison* (Niu et al., 2016):

$$\bar{\rho}_{\text{PU},\text{PN}} := \lim_{n_{\text{U}} \rightarrow \infty} \rho_{\text{PU},\text{PN}} = \frac{\theta_{\text{P}}/\sqrt{n_{\text{P}}}}{\theta_{\text{N}}/\sqrt{n_{\text{N}}}}, \quad (2.64)$$

$$\bar{\rho}_{\text{NU},\text{PN}} := \lim_{n_{\text{U}} \rightarrow \infty} \rho_{\text{NU},\text{PN}} = \frac{\theta_{\text{N}}/\sqrt{n_{\text{N}}}}{\theta_{\text{P}}/\sqrt{n_{\text{P}}}}. \quad (2.65)$$

Since

$$\bar{\rho}_{\text{PU},\text{PN}} \times \bar{\rho}_{\text{NU},\text{PN}} = 1, \quad (2.66)$$

there are only three cases needed to be considered:

$$\bar{\rho}_{\text{PU},\text{PN}} < 1 \text{ and } \bar{\rho}_{\text{NU},\text{PN}} > 1, \quad (2.67)$$

$$\bar{\rho}_{\text{PU},\text{PN}} > 1 \text{ and } \bar{\rho}_{\text{NU},\text{PN}} < 1, \quad (2.68)$$

$$\bar{\rho}_{\text{PU},\text{PN}} = 1 \text{ and } \bar{\rho}_{\text{NU},\text{PN}} = 1. \quad (2.69)$$

The first and second cases indicate that the classifier based on either PU or NU classification is more promising than that based on PN classification although PU or NU classification can not access negative or positive data, respectively. The exception is the third case which happens when $n_{\text{P}}/n_{\text{N}} = \theta_{\text{P}}^2/\theta_{\text{N}}^2$, but the probability of occurring the third case is close to zero because n_{P} and n_{N} are independently determined by θ_{P} and θ_{N} , in practice.

2.6 Semi-Supervised Classification

In this section, we review the existing semi-supervised learning algorithms, which are compared with our methods in experiments.

To simplify the notation, we introduce the followings:

$$\{(\mathbf{x}_i^L, y_i)\}_{i=1}^{n_L} = \{(\mathbf{x}_i^P, +1)\}_{i=1}^{n_P} \cup \{(\mathbf{x}_j^N, -1)\}_{j=1}^{n_N},$$

where $n_L = n_P + n_N$ and $\{c_i\}_{i=1}^{n_L}$ is a set of weights taking $c_i = \theta_P/n_P$ if $y_i = +1$ and θ_N/n_N if $y_i = -1$.

2.6.1 Entropy Regularization

Based on the *entropy minimization principle*, *entropy regularization* (ER) (Grandvalet and Bengio, 2004) was proposed. From unlabeled samples, ER computes the conditional entropy of class labels conditioned on the inputs.

Specifically, let $q_w(y | \mathbf{x})$ be a model of posterior distribution $p(y | \mathbf{x})$. The objective function consists of the likelihood function and the entropy term:

$$\max_{\mathbf{w} \in \mathbb{R}^b} \left[\sum_{i=1}^{n_L} c_i \log q_w(y_i | \mathbf{x}_i^L) + \lambda_E \sum_{k=1}^{n_U} \sum_{y=\pm 1} q_w(y | \mathbf{x}_k^U) \log q_w(y | \mathbf{x}_k^U) \right], \quad (2.70)$$

where $\lambda_E \geq 0$ is the regularization parameter. The first term can be regarded as (weighted) conditional likelihood, and the second term can be regarded as negative entropy. The entropy regularization is known as that it favors a decision boundary lying on a low-density region (Chapelle et al., 2006). Therefore, if the low-density separation principle does not hold on data points, the regularization scheme does not work well.

From the optimization viewpoint, the objective function in Eq. (2.70) is non-convex. Hence, the solution which we obtain is usually local optima. Following to Grandvalet and Bengio (2006), we first minimize the objective without entropy regularization, i.e., ordinary logistic regression, and then minimize the objective in Eq. (2.70) using the first solution as an initial solution for a quasi-newton method (Nocedal and Wright, 2006).

2.6.2 Manifold Regularization

The *manifold assumption* presumes that data points are supported on a (near) low-dimensional manifold (Belkin et al., 2006). Based on the manifold assumption, Belkin

et al. (2006) proposed the *manifold regularization* approach, encoding a manifold structure in the data to a regularizer.

Laplacian support vector machines (LapSVM) (Belkin et al., 2006) is the *support vector machine* (SVM) (Cortes and Vapnik, 1995) with the manifold regularization. Let us define a kernel model as $g(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i)$, where α is dual variable, $K(\mathbf{x}, \mathbf{x}')$ is the kernel function. The problem of LapSVM (in primal) (Melacci and Belkin, 2011) can be expressed as

$$\min_{\alpha \in \mathbb{R}^n} \left[\sum_{i=1}^{n_L} c_i \ell_{\text{SH}}(1 - y_i \mathbf{k}_i^\top \alpha) + C_A \alpha^\top \mathbf{K} \alpha + C_I \alpha^\top \mathbf{K} \mathbf{L} \mathbf{K} \alpha \right], \quad (2.71)$$

where C_A is the weight of the norm of the function (or *ambient* norm), C_I is the weight of the norm of function in the low dimensional manifold (or *intrinsic* norm), $\mathbf{L} = \mathbf{D} - \mathbf{W}$, \mathbf{W} is the adjacency matrix of the data graph, \mathbf{D} is the diagonal matrix whose element is $D_{i,i} = \sum_{j=1}^n W_{i,j}$, $k_i^{(j)} = K(\mathbf{x}_i^L, \mathbf{x}_j)$, and $K_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$.

Compared with the originally proposed LapSVM objective function (Belkin et al., 2006), the above objective employs the squared hinge loss ℓ_{SH} for labeled samples so that its gradient can be computed, unlike the hinge loss. Melacci and Belkin (2011) showed that solving the above objective function by preconditioned conjugate gradient along with early stopping significantly reduces the computation time.

2.6.3 Squared-Loss Mutual Information Regularization

Squared-Loss Mutual Information Regularization (SMIR) (Niu et al., 2013) is based on the *information maximization principle* (Linsker, 1988). The idea of SMIR is to learn the class-posterior probability $p(y \mid \mathbf{x})$ while maximizing *squared-loss mutual information* (SMI) (Suzuki et al., 2009) defined as

$$\text{SMI} := \frac{1}{2} \sum_{y=\pm 1} \int p(\mathbf{x}) p(y) \left(\frac{p(\mathbf{x}, y)}{p(\mathbf{x}) p(y)} - 1 \right)^2 d\mathbf{x}. \quad (2.72)$$

Unlike maximization of *mutual information* (Cover and Thomas, 2006), the optimization problem of SMIR is strictly convex under mild conditions. It thus can be solved efficiently and achieves the unique globally optimal solution.

Specifically, let $q_\alpha(y \mid \mathbf{x}) := \alpha_y^\top \mathbf{D}^{-\frac{1}{2}} \mathbf{K}^{-\frac{1}{2}} \phi_n(\mathbf{x})$ be the model of the class-posterior probability, where $K_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$, $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$, and $d_i = \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j)$.

The optimization problem of SMIR can be formulated as

$$\min_{\mathbf{A} \in \mathbb{R}^{n \times 2}} \left[\widehat{\Delta}_2(p, q_{\alpha}) - \gamma_S \widehat{\text{SMI}} + \frac{\lambda_S}{2} \text{tr}(\mathbf{A}^\top \mathbf{A}) \right], \quad (2.73)$$

where $\mathbf{A} := (\alpha_{+1}, \alpha_{-1})$ is the matrix representation of model parameters, γ_S and λ_S are parameters for trading off between the SMI regularization and ℓ_2 -regularization. $\Delta_2(p, q)$ is the squared difference between p and q defined as

$$\Delta_2(p, q) := \frac{1}{2} \sum_{y=\pm 1} \int (p(y | \mathbf{x}) - q(y | \mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x}. \quad (2.74)$$

As explained, if the condition $\lambda_S > \eta / (n \min_{y=\pm 1} p(y))$ is satisfied, the optimization problem can be strictly convex with respect to \mathbf{A} .

2.6.4 Maximum-Margin Principle

SVM has been widely used in many applications and known as obtaining a large margin decision boundary. Based on the *maximum-margin principle*, *semi-supervised SVM* (S3VM) aims at learning a classifier by minimizing the risk computed from labeled data while maximizing margin on unlabeled data (Chapelle et al., 2008).

While the original S3VM solves the non-convex optimization problem, weakly labeled support vector machines (WellSVM) proposed by Li et al. (2013) solve the tightly relaxed convex optimization problem, which can be solved efficiently.

The optimization problem of WellSVM can be expressed as

$$\min_{\{\mu_t\} \in \mathcal{M}} \max_{\alpha \in \mathcal{A}} \left[\mathbf{1}^\top \alpha - \frac{1}{2} \alpha^\top \left[\sum_{t: \hat{\mathbf{y}}_t \in \mathcal{B}} \mu_t \left(\mathbf{K} \odot (\hat{\mathbf{y}}_t \hat{\mathbf{y}}_t^\top) \right) \right] \alpha \right], \quad (2.75)$$

where

$$\mathcal{A} = \{\alpha \mid 0 \leq \alpha_i \leq C_L c_i, 0 \leq \alpha_j \leq C_U, i \in \mathcal{L}, j \in \mathcal{U}\}, \quad (2.76)$$

$$\mathcal{B} = \{\hat{\mathbf{y}} \mid \hat{y}_i = y_i, \hat{y}_j \in \{\pm 1\}, \frac{1}{n_U} \sum_j I(\hat{y}_j = +1) = \theta_p, i \in \mathcal{L}, j \in \mathcal{U}\}, \quad (2.77)$$

\odot denotes the element-wise multiplication (a.k.a. *Hadamard* product), C_L and C_U are the regularization parameter on the losses for labeled and unlabeled data, \mathcal{L} and \mathcal{U} are

the index sets for labeled and unlabeled data, and $I(\cdot)$ is the indicator function. To solve Eq.(2.75), the authors utilized the *cutting plane algorithm* (Kelley, 1960) by iteratively generating label. More specifically, we first initialize a label vector $\hat{\mathbf{y}}$ and the working set \mathcal{C} to $\hat{\mathbf{y}}$. We then solve Eq.(2.75) with respect to α by standard supervised learning methods. After that, we generate a violated vector $\hat{\mathbf{y}}$ (Li et al., 2013), and add it into \mathcal{C} . We repeat the above procedure until the decrease of objective function is smaller than a threshold.

2.6.5 Safe Approach

As discussed in Section 1.2.3, once the distributional assumptions are not satisfied, the performance of the existing methods becomes even worse than the method based on supervised learning. To mitigate this problem, a *safe* semi-supervised learning approach was proposed to avoid undesirable performance degeneration (Li and Zhou, 2015).

Based on the safe approach, *safe S3VM* (S4VM) (Li and Zhou, 2015) was proposed so that its performance is not worse than the ordinary supervised SVM. Specifically, let $\mathbf{y}^* \in \mathbb{R}^{n_U}$ be the ground-truth label assignment and $\mathbf{y}^{\text{svm}} \in \mathbb{R}^{n_U}$ be the predictive labels of SVM on unlabeled instances. The $\text{gain}(\mathbf{y}, \mathbf{y}^*, \mathbf{y}^{\text{svm}})$ and $\text{loss}(\mathbf{y}, \mathbf{y}^*, \mathbf{y}^{\text{svm}})$ respectively measure the gained and lost accuracies compared to SVM, where $\mathbf{y} := (y_1, \dots, y_{n_U})^\top$. The goal of S4VM is to maximize the improvement against SVM:

$$\hat{\mathbf{y}} = \underset{\mathbf{y} \in \{\pm 1\}^{n_U}}{\operatorname{argmax}} \min_{\bar{\mathbf{y}} \in \mathcal{M}} \left[\text{gain}(\mathbf{y}, \bar{\mathbf{y}}, \mathbf{y}^{\text{svm}}) - \lambda_{\text{S4VM}} \cdot \text{loss}(\mathbf{y}, \bar{\mathbf{y}}, \mathbf{y}^{\text{svm}}) \right], \quad (2.78)$$

where

$$\text{gain}(\mathbf{y}, \bar{\mathbf{y}}, \mathbf{y}^{\text{svm}}) = \sum_{j=n_L+1}^{n_L+n_U} \left(c_P \frac{1+y_j}{2} + c_N \frac{1-y_j}{2} \right) \frac{1+y_j \bar{y}_j}{2} \frac{1-y_j^{\text{svm}} \bar{y}_j}{2}, \quad (2.79)$$

$$\text{loss}(\mathbf{y}, \bar{\mathbf{y}}, \mathbf{y}^{\text{svm}}) = \sum_{j=n_L+1}^{n_L+n_U} \left(c_P \frac{1+y_j}{2} + c_N \frac{1-y_j}{2} \right) \frac{1-y_j \bar{y}_j}{2} \frac{1+y_j^{\text{svm}} \bar{y}_j}{2}, \quad (2.80)$$

$c_P = \theta_P/n_P$, $c_N = \theta_N/n_N$, and λ_{S4VM} is a parameter for controlling the trade-off between gain and loss functions. S4VM assumes that the ground-truth is realized by a low-density separator, i.e., $\mathbf{y}^* \in \mathcal{M}$, where $\mathcal{M} := \{\bar{\mathbf{y}}_t\}_{t=1}^T$ is a pool of low-density separators obtained by, e.g., S3VM. If the class-prior of unlabeled data is θ_P different from that of labeled data,

we restrict that the assigned labels \mathbf{y} are in $\mathcal{B} = \{\mathbf{y} \in \{\pm 1\}^{n_U} \mid -\beta \leq \frac{1}{n_U} \sum_{i=1}^{n_U} y_i - \theta_P \leq \beta\}$, where β is a small constant controlling the inconsistency of class prior.

Although S4VM can avoid unexpected performance degeneration, its core mechanism for utilizing unlabeled data still relies on the low-density separation principle. Thus, the performance degeneration may not happen, but S4VM cannot leverage information of unlabeled data.

2.7 Class-Prior Estimation Methods

In this section, we review class-prior estimation methods used in numerical experiments in this dissertation.

2.7.1 Class-Prior Shift

In practice, we are often faced with *dataset shift* or *distribution change*, where distributions between training data and test data are different:

$$p_{\text{tr}}(\mathbf{x}, y) \neq p_{\text{te}}(\mathbf{x}, y), \quad (2.81)$$

where p_{tr} and p_{te} are training and test densities, respectively. When class-conditional densities are the same but class-priors are different between training and test distributions, the situation is called *class-prior shift* (Quiñonero Candela et al., 2009).⁴ Formally, class-prior shift can be expressed as

$$p_{\text{tr}}(\mathbf{x} \mid y) = p_{\text{te}}(\mathbf{x} \mid y), \quad (2.82)$$

$$p_{\text{tr}}(y) \neq p_{\text{te}}(y). \quad (2.83)$$

Class-prior shift occurs in many real-world applications. For example, suppose that we would like to obtain a classifier for discriminating a face image into male or female. The classifier can be trained with, for instance, face images collected from a lab. However, the face images often contains mostly face images from males and smaller size of face images from females. Thus, if we simply apply an ordinary classification method, a learned classifier will be reflected by the proportion of male and female faces in the training data. On the other hand, in real-world, the proportion of male and female is almost equal. Thus, the learned classifier on training data will perform poorly on test distribution.

⁴ There are other situations studied such as *covariate shift* and *target shift*.

To cope with the class-prior shift situation, we train a classifier with training samples weighted by class-priors of test data. For example, if we have domain knowledge about class-priors of test distribution, we can train a classifier with the class-priors, e.g., the use of a weighted SVM with prespecified weights. However, it is not always the case; we need an estimate of class-prior of test data.

For estimating a class-prior, several class-prior estimation methods have been proposed for both PU and semi-supervised learning situations, and shown its usefulness in experiments (Kawakubo et al., 2016; Ramaswamy et al., 2016; du Plessis et al., 2017).

2.7.2 Class-Prior Estimation under Semi-Supervised Learning Setting

Here, we review a class-prior estimation method based on *energy distance* minimization (Kawakubo et al., 2016), which employs the energy distance (Székely and Rizzo, 2013) for measuring the difference between two distributions.

The energy distance (ED) between densities p and q is defined as

$$\text{ED}(p, q) := \int \|\psi_p(\mathbf{t}) - \psi_q(\mathbf{t})\|^2 \left(\frac{\pi^{\frac{d+1}{2}}}{\Gamma(\frac{d+1}{2})} \|\mathbf{t}\|^{d+1} \right) d\mathbf{t}, \quad (2.84)$$

where $\|\cdot\|$ denotes the Euclidean distance, Γ is the *gamma* function, and ψ_p and ψ_q are the characteristic functions of p and q , respectively. An advantage of the energy distance is that it can be expressed as expectations of Euclid distance:

$$\text{ED}(p, q) = 2 \mathbb{E}_{\mathbf{x} \sim p, \bar{\mathbf{x}} \sim q} [\|\mathbf{x} - \bar{\mathbf{x}}\|] - \mathbb{E}_{\mathbf{x}, \bar{\mathbf{x}} \sim p} [\|\mathbf{x} - \bar{\mathbf{x}}\|] - \mathbb{E}_{\mathbf{x}, \bar{\mathbf{x}} \sim q} [\|\mathbf{x} - \bar{\mathbf{x}}\|]. \quad (2.85)$$

Let us define a distribution model parameterized by class-priors as

$$q_\theta(\mathbf{x}) = \theta_P p(\mathbf{x} \mid y = +1) + \theta_N p(\mathbf{x} \mid y = -1). \quad (2.86)$$

To obtain estimated class-priors, we minimize the energy distance between a marginal density p and the distribution model q_{θ_P} . The estimated class-prior can be obtained by

minimizing the energy distance with respect to θ_P and θ_N :

$$\hat{\theta}_P, \hat{\theta}_N = \operatorname{argmin}_{\theta=\{\theta_P, \theta_N\}} \operatorname{ED}(p, q_\theta), \quad (2.87)$$

$$\begin{aligned} \operatorname{ED}(p, q_\theta) &= 2\theta_P \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}), \bar{\mathbf{x}} \sim p(\mathbf{x}|y=+1)}[\|\mathbf{x} - \bar{\mathbf{x}}\|] + 2\theta_N \mathbb{E}_{\mathbf{x} \sim p, \bar{\mathbf{x}} \sim p(\mathbf{x}|y=-1)}[\|\mathbf{x} - \bar{\mathbf{x}}\|] \\ &\quad - \mathbb{E}_{\mathbf{x}, \bar{\mathbf{x}} \sim p}[\|\mathbf{x} - \bar{\mathbf{x}}\|] - \theta_P^2 \mathbb{E}_{\mathbf{x}, \bar{\mathbf{x}} \sim p(\mathbf{x}|y=+1)}[\|\mathbf{x} - \bar{\mathbf{x}}\|] \\ &\quad - \theta_N^2 \mathbb{E}_{\mathbf{x}, \bar{\mathbf{x}} \sim p(\mathbf{x}|y=-1)}[\|\mathbf{x} - \bar{\mathbf{x}}\|] \\ &\quad + \theta_P \theta_N \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|y=+1), \bar{\mathbf{x}} \sim p(\mathbf{x}|y=-1)}[\|\mathbf{x} - \bar{\mathbf{x}}\|]. \end{aligned} \quad (2.88)$$

2.7.3 Class-Prior Estimation under PU Learning Setting

In this section, we review class-prior estimation methods for the PU learning situation (du Plessis and Sugiyama, 2014; du Plessis et al., 2017). To estimate a class-prior from PU data, du Plessis and Sugiyama (2014) proposed a *partial distribution matching* approach.

Let us define a *partial model* as

$$q'_{\theta_P}(\mathbf{x}) := \theta_P p(\mathbf{x} \mid y = +1). \quad (2.89)$$

Then, we obtain an estimate of a class-prior by minimizing a discrepancy measure, such as the *Kullback-Leibler* (KL) divergence, between the partial model and the true marginal density with respect to θ_P :

$$\hat{\theta}_P = \operatorname{argmin}_{0 \leq \theta_P \leq 1} \operatorname{Div}_f[p, q'_{\theta_P}], = \operatorname{argmin}_{0 \leq \theta_P \leq 1} \int p(\mathbf{x}) f\left(\frac{q'_{\theta_P}(\mathbf{x})}{p(\mathbf{x})}\right) d\mathbf{x}, \quad (2.90)$$

where Div_f is the f -divergence (Ali and Slivey, 1966; Csizár, 1967) and $f(t)$ is a convex function with $f(1) = 0$.

Although the partial distribution matching via f -divergences had been shown its effectiveness on various datasets experimentally, it was known that the partial matching via f -divergences overestimates the true class-prior unless distributions has no overlap. To cope with this problem, a partial distribution matching via *penalized* f -divergences was proposed. In penalized f -divergences, we replace the function f in Eq. (2.90) with the function satisfying the following condition:

$$\tilde{f}(t) = \begin{cases} f(t) & 0 \leq t \leq 1, \\ \infty & \text{otherwise.} \end{cases} \quad (2.91)$$

The penalized f -divergence is restricting the range of the density-ratio of $q'_{\theta_P}/p(\mathbf{x})$ so as to be $[0, 1]$. We then obtain an estimate of the class-prior by solving the following optimization problem:

$$\hat{\theta}_P = \operatorname{argmin}_{0 \leq \theta_P \leq 1} \int p(\mathbf{x}) \tilde{f}\left(\frac{q'_{\theta_P}(\mathbf{x})}{p(\mathbf{x})}\right) d\mathbf{x}. \quad (2.92)$$

du Plessis et al. (2017) prove that consistency, stability, and an estimation error of their class-prior estimation method. Furthermore, they showed an estimate of class-prior can be computed efficiently with the partial matching via the penalized L_1 -distance.

2.8 Statistical Dependency

In this section, we review two statistical dependency measures: *mutual information* (MI) (Cover and Thomas, 2006) and *squared-loss mutual information* (SMI) (Suzuki et al., 2009). Note that in this dissertation, since our focus is a binary output, the definition of these statistical dependency measures will be given by the binary classification case.

2.8.1 Mutual Information (MI)

Mutual information (MI), originally developed in *information theory* (Shannon, 1948), is a vital tools for statistical dependency analysis. MI is defined as

$$\text{MI} := \sum_{y=\pm 1} \int p(\mathbf{x}, y) \log \left(\frac{p(\mathbf{x}, y)}{p(\mathbf{x})p(y)} \right) d\mathbf{x}. \quad (2.93)$$

MI can be regarded as the *Kullback-Leibler divergence* from $p(\mathbf{x}, y)$ to $p(\mathbf{x})p(y)$. Thus, MI is non-negative and takes zero if and only if $p(\mathbf{x}, y) = p(\mathbf{x})p(y)$, i.e., \mathbf{x} and y are statistically independent. This property allows us to measure the dependency between \mathbf{x} and y .

Based on the *information maximization principle*, empirical estimators of MI have been employed for solving various machine learning tasks (Torkkola, 2003; Krause et al., 2010). However, since such MI estimators are known to be sensitive to outliers (Basu et al., 1998), its practical usefulness is limited. Thus, more practical dependency measures are desirable.

2.8.2 Squared-Loss Mutual Information (SMI)

To cope with the sensitivity issue of MI, the squared-loss variant of MI has been proposed and its supervised estimator was developed (Suzuki et al., 2009).

Squared-loss MI (SMI) is defined as

$$\text{SMI} := \frac{1}{2} \sum_{y=\pm 1} \int p(\mathbf{x})p(y) \left(\frac{p(\mathbf{x}, y)}{p(\mathbf{x})p(y)} - 1 \right)^2 d\mathbf{x}. \quad (2.94)$$

SMI can be regarded as the *Pearson divergence* (Pearson, 1990) from $p(\mathbf{x}, y)$ to $p(\mathbf{x})p(y)$. That is, SMI is also non-negative and takes zero if and only if $p(\mathbf{x}, y) = p(\mathbf{x})p(y)$.

Unlike MI, SMI is more robust to outliers and has superior numerical properties (Kanamori et al., 2012). So far, various SMI-based machine learning algorithms have been proposed and demonstrated their effectiveness: independence test (Sugiyama and Suzuki, 2011), dimension reduction (Suzuki and Sugiyama, 2013), clustering (Sugiyama et al., 2014), object matching (Yamada et al., 2015), and image registration (Sakai et al., 2015).

SMI Estimation: In practice, since $p(\mathbf{x}, y)$ is unknown, SMI cannot be computed directly. Here, we review an SMI estimation method from samples. Suppose that we are given a set of *independent and identically distributed* (i.i.d.) paired samples from underlying joint density:

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y). \quad (2.95)$$

A naive approach to estimate SMI is as follows: first densities $p(\mathbf{x}, y)$, $p(\mathbf{x})$, and $p(y)$ are estimated from the samples, and then the estimated densities $\hat{p}(\mathbf{x}, y)$, $\hat{p}(\mathbf{x})$, and $\hat{p}(y)$ are plugged into the definition of SMI. However, such a two-step approach is not reliable since estimation error of densities can be magnified, in particular, when division by estimated densities is performed in $\frac{\hat{p}(\mathbf{x}, y)}{\hat{p}(\mathbf{x})\hat{p}(y)}$. To cope with this problem, the SMI estimation method via the direct density-ratio estimation was proposed (Suzuki et al., 2009; Sakai and Sugiyama, 2014).

Let

$$r(\mathbf{x}, y) := \frac{p(\mathbf{x}, y)}{p(\mathbf{x})p(y)} \quad (2.96)$$

be the density-ratio function and $q(\mathbf{x}, y)$ be a model of the density-ratio. The model is learned so that it minimizes the squared error (the PN-SMI risk) defined as

$$\begin{aligned} R_{\text{PN}}^{\text{SMI}}(q) &= \frac{1}{2} \sum_{y=\pm 1} \int (q(\mathbf{x}, y) - r(\mathbf{x}, y))^2 p(\mathbf{x}) p(y) d\mathbf{x} - C \\ &= \frac{1}{2} \sum_{y=\pm 1} \int q^2(\mathbf{x}, y) p(\mathbf{x}) p(y) d\mathbf{x} - \sum_{y=\pm 1} \int q(\mathbf{x}, y) p(\mathbf{x}, y) d\mathbf{x}, \end{aligned} \quad (2.97)$$

where we used $r(\mathbf{x}, y)p(\mathbf{x})p(y) = p(\mathbf{x}, y)$ to obtain the second term and $C = \frac{1}{2} \sum_{y=\pm 1} \int r^2(\mathbf{x}, y) p(\mathbf{x}) p(y) d\mathbf{x}$. By replacing the expectations with corresponding sample averages, we obtain an empirical squared error (the empirical PN-SMI risk) as

$$\widehat{R}_{\text{PN}}^{\text{SMI}}(q) = \sum_{y=\pm 1} \frac{n_y}{2n^2} \sum_{i=1}^n q^2(\mathbf{x}_i, y) - \frac{1}{n} \sum_{i=1}^n q(\mathbf{x}_i, y_i), \quad (2.98)$$

where n_y is the number of samples for each class in $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$. The learned model \widehat{q} can be obtained by minimizing the regularized empirical PN-SMI risk:

$$\widehat{q} := \underset{q}{\operatorname{argmin}} \left[\widehat{R}_{\text{PN}}^{\text{SMI}}(q) + \lambda W(q) \right], \quad (2.99)$$

where W is a regularizational functional and $\lambda \geq 0$ is the regularization parameter. With the estimated parameter, an SMI approximator can be obtained by

$$\widehat{\text{SMI}} := \frac{1}{n} \sum_{i=1}^n \widehat{q}(\mathbf{x}_i, y_i) - \sum_{y=\pm 1} \frac{n_y}{2n^2} \sum_{i=1}^n \widehat{q}^2(\mathbf{x}_i, y) - \frac{1}{2}, \quad (2.100)$$

since SMI can be expressed as

$$\text{SMI} = \frac{1}{2} \sum_{y=\pm 1} \int r^2(\mathbf{x}, y) p(\mathbf{x}) p(y) d\mathbf{x} - \frac{1}{2} \quad (2.101)$$

$$= \sum_{y=\pm 1} \int r(\mathbf{x}, y) p(\mathbf{x}, y) d\mathbf{x} - \frac{1}{2} \sum_{y=\pm 1} \int r^2(\mathbf{x}, y) p(\mathbf{x}) p(y) d\mathbf{x} - \frac{1}{2}, \quad (2.102)$$

where the first equation can be obtained by expanding the definition of SMI with the relation $r(\mathbf{x}, y)p(\mathbf{x}, y) = p(\mathbf{x})p(y)$, and the equality between the first and second equations can be confirmed by using $r(\mathbf{x}, y)p(\mathbf{x}, y) = p(\mathbf{x})p(y)$ for the second term in the second equation.

Note that the above SMI approximator can be obtained by *Legendre-Fenchel* convex

duality (Boyd and Vandenberghe, 2004) as in the f -divergence estimation (Keziou, 2003; Nguyen et al., 2007):

$$\text{SMI} = -\inf_g R_{\text{PN}}^{\text{SMI}}(g) - \frac{1}{2}. \quad (2.103)$$

This implies that minimizing the PN-SMI risk corresponds to maximizing lower bound of SMI.

Practical Implementation: In practice, the model g can be either parametric or non-parametric. A simple choice is to use a linear-in-parameter model:

$$q(\mathbf{x}, y) := \sum_{\ell=1}^b \alpha_{\ell} \psi_{\ell}(\mathbf{x}, y) = \boldsymbol{\alpha}^{\top} \boldsymbol{\psi}(\mathbf{x}, y), \quad (2.104)$$

where $\boldsymbol{\alpha} := (\alpha_1, \dots, \alpha_b)^{\top}$ is the parameter vector, $\boldsymbol{\psi}(\mathbf{x}, y) := (\psi_1(\mathbf{x}, y), \dots, \psi_b(\mathbf{x}, y))$ is the basis function vector, b is the number of basis functions, and $^{\top}$ denotes the transpose of vectors and matrices. Then, with the ℓ_2 -regularizer $W(q) = \boldsymbol{\alpha}^{\top} \boldsymbol{\alpha} / 2$, the optimization problem yields

$$\hat{\boldsymbol{\alpha}} := \underset{\boldsymbol{\alpha} \in \mathbb{R}^b}{\text{argmin}} \left[\frac{1}{2} \boldsymbol{\alpha}^{\top} \widehat{\mathbf{H}} \boldsymbol{\alpha} - \boldsymbol{\alpha}^{\top} \widehat{\mathbf{h}} + \frac{\lambda}{2} \boldsymbol{\alpha}^{\top} \boldsymbol{\alpha} \right], \quad (2.105)$$

where

$$\widehat{H}_{\ell, \ell'} := \sum_{y=\pm 1} \frac{n_y}{n^2} \sum_{i=1}^n \psi_{\ell}(\mathbf{x}_i, y) \psi_{\ell'}(\mathbf{x}_i, y)^{\top}, \quad (2.106)$$

$$\widehat{h}_{\ell} := \frac{1}{n} \sum_{i=1}^n \psi_{\ell}(\mathbf{x}_i, y_i). \quad (2.107)$$

The solution can be obtained analytically by differentiating the objective function with respect to $\boldsymbol{\alpha}$ and set it to zero:

$$\hat{\boldsymbol{\alpha}} = (\widehat{\mathbf{H}} + \lambda \mathbf{I}_b)^{-1} \widehat{\mathbf{h}}, \quad (2.108)$$

where \mathbf{I}_b is the $b \times b$ identity matrix. Then, the density-ratio can be estimated by

$$\widehat{q}(\mathbf{x}, y) = \widehat{\boldsymbol{\alpha}}^{\top} \boldsymbol{\psi}(\mathbf{x}, y). \quad (2.109)$$

With the estimated density-ratio, the SMI approximator can be expressed as

$$\widehat{\text{SMI}} = \widehat{\boldsymbol{\alpha}}^\top \widehat{\boldsymbol{h}} - \frac{1}{2} \widehat{\boldsymbol{\alpha}}^\top \widehat{\boldsymbol{H}} \widehat{\boldsymbol{\alpha}} - \frac{1}{2}. \quad (2.110)$$

Chapter 3

Imbalanced Classification from Positive and Unlabeled Data

In this chapter, we discuss learning from positive and unlabeled data in an imbalanced classification scenario.

3.1 Introduction

Supervised learning is widely used and studied in many fields. However, in practice, positive data can be easily accessed, but negative data are not due to diverse representations of negative class or additional expensive costs to collect negative labels. For instance, a friend in a social network service is easy to observe through users' declaration, while a non-friend is difficult unless we conduct a user study. To address this issue, learning from positive and unlabeled data (PU learning) has been studied and revealed useful results.

In addition to the PU learning scenario, we are sometimes faced with the situation where the number of positive data is significantly smaller than that of negative data, i.e., the class-ratio is *imbalanced*. In such an imbalanced classification problem, maximizing the *area under the receiver operating characteristic curve* (AUC) (Hanley and McNeil, 1982) is a standard approach (Cortes and Mohri, 2003). While the misclassification rate relies on the sign of the score of a single sample, AUC is governed by the ranking of the scores of *two* samples. Based on this principle, various supervised methods for directly optimizing AUC have been developed so far and demonstrated to be useful (Herschtal and Raskutti, 2004; Zhao et al., 2011; Rakhlin et al., 2012; Kotlowski et al., 2011; Ying et al., 2016). Previously, a *pairwise ranking* method for PU data has been developed (Sundararajan et al., 2011), which can be regarded as an AUC optimization method for PU data. However, it merely regards unlabeled data as negative data, and thus the obtained classifier is biased.

In this chapter, we present a PU learning method for imbalanced classification. In contrast to [Sundararajan et al. \(2011\)](#), the method is based on an unbiased risk estimator. Furthermore, we prove the generalization error bound of the method and show that unlabeled data contribute to reducing an upper bound on the generalization error with the optimal parametric convergence rate. Experiments confirm the effectiveness of the PU learning method for imbalanced classification.

3.2 Problem Setting

Let random variables $\mathbf{x} \in \mathbb{R}^d$ and $y \in \{+1, -1\}$ be equipped with probability density $p(\mathbf{x}, y)$, where d is a positive integer. Let us consider a binary classification problem from \mathbf{x} to y , given two sets of samples called the *positive* (P) and *unlabeled* (U) data:

$$\mathcal{X}_P := \{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = +1), \quad (3.1)$$

$$\mathcal{X}_U := \{\mathbf{x}_k^U\}_{k=1}^{n_U} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}) = \theta_P p(\mathbf{x} \mid y = +1) + \theta_N p(\mathbf{x} \mid y = -1), \quad (3.2)$$

where

$$\theta_P := p(y = +1) \text{ and } \theta_N := p(y = -1).$$

Let E_P , E_N , and E_U be the expectations over $p(\mathbf{x} \mid y = +1)$, $p(\mathbf{x} \mid y = -1)$, and $p(\mathbf{x})$, respectively.

Moreover, let $g: \mathbb{R}^d \rightarrow \mathbb{R}$ be an arbitrary real-valued decision function for binary classification, and classification is performed based on its sign. Let $\ell: \mathbb{R} \rightarrow \mathbb{R}$ be a loss function. The goal is to obtain a better classifier in the imbalanced classification scenario.

3.3 PU-AUC Optimization

In this section, we explain the AUC optimization method for the PU learning setting.

Similarly to [du Plessis et al. \(2014, 2015\)](#), the risk function is derived for AUC optimization in PU learning setting. That is, the risk function is equivalent to the risk of its supervised counterpart, and it depends only on positive and unlabeled data distributions. In the derivation and theoretical analysis, the class-priors θ_P and θ_N are assumed to be known. In practice, they are replaced with their estimates obtained, e.g., by [Ramaswamy et al. \(2016\)](#) and [du Plessis et al. \(2017\)](#), and references therein.

Recall the PN-AUC risk in Eq. (2.23):

$$R_{\text{PN}}^{\text{AUC}}(g) := \mathbb{E}_{\text{P}}[\mathbb{E}_{\text{N}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\mathbf{x}^{\text{N}}))]]. \quad (3.3)$$

By definition of the marginal density $p(\mathbf{x}) = \theta_{\text{P}}p(\mathbf{x} \mid y = +1) + \theta_{\text{N}}p(\mathbf{x} \mid y = -1)$, we have

$$\mathbb{E}_{\text{U}}[\cdot] = \theta_{\text{P}} \mathbb{E}_{\text{P}}[\cdot] + \theta_{\text{N}} \mathbb{E}_{\text{N}}[\cdot]. \quad (3.4)$$

Thus, we obtain the following relation:

$$\begin{aligned} & \mathbb{E}_{\text{P}}[\mathbb{E}_{\text{U}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\mathbf{x}^{\text{U}}))]] \\ &= \theta_{\text{P}} \mathbb{E}_{\text{P}}[\mathbb{E}_{\bar{\text{P}}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\bar{\mathbf{x}}^{\text{P}}))]] + \theta_{\text{N}} \mathbb{E}_{\text{P}}[\mathbb{E}_{\text{N}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\mathbf{x}^{\text{N}}))]] \\ &= \theta_{\text{P}} \mathbb{E}_{\text{P}}[\mathbb{E}_{\bar{\text{P}}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\bar{\mathbf{x}}^{\text{P}}))]] + \theta_{\text{N}} R_{\text{PN}}^{\text{AUC}}(g), \end{aligned} \quad (3.5)$$

where $\mathbb{E}_{\bar{\text{P}}}$ denotes the expectation over $p(\bar{\mathbf{x}}^{\text{P}} \mid y = +1)$. Dividing the above equation by θ_{N} and rearranging it, we can express the PN-AUC risk based on PU data (the PU-AUC risk) as

$$\begin{aligned} R_{\text{PN}}^{\text{AUC}}(g) &= \frac{1}{\theta_{\text{N}}} \mathbb{E}_{\text{P}}[\mathbb{E}_{\text{U}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\mathbf{x}^{\text{U}}))]] \\ &\quad - \frac{\theta_{\text{P}}}{\theta_{\text{N}}} \mathbb{E}_{\text{P}}[\mathbb{E}_{\bar{\text{P}}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\bar{\mathbf{x}}^{\text{P}}))]] := R_{\text{PU}}^{\text{AUC}}(g). \end{aligned} \quad (3.6)$$

Thus, we immediately obtain the following lemma.

Lemma 3. *The AUC, PN-AUC, and PU-AUC risks are equivalent:*

$$R^{\text{AUC}}(g) = R_{\text{PN}}^{\text{AUC}}(g) = R_{\text{PU}}^{\text{AUC}}(g). \quad (3.7)$$

In contrast to the PN-AUC risk, the PU-AUC risk can be computed by the expectations over P and U data, but the minimizers are equivalent, meaning that the performance of the classifier obtained by using the PU-AUC risk is the same as the one obtained by using the PN-AUC risk, i.e., supervised learning. We refer to the method minimizing the PU-AUC risk as *PU-AUC optimization*. We explain a practical implementation in Section 3.4 and investigate the superiority of $R_{\text{PU}}^{\text{AUC}}$ in Section 3.5.

As discussed in Section 3.1, [Sundararajan et al. \(2011\)](#) proposed a pairwise ranking method for PU data, which can be regarded as an AUC optimization method for PU data.

Their approach simply regards unlabeled data as negative data, and the ranking SVM (Joachims, 2002) is applied to PU data so that the score of positive data tends to be higher than that of unlabeled data. Although this approach is simple and shown computationally efficient in experiments, the obtained classifier is biased. From the mathematical viewpoint, the existing method ignores the second term in Eq. (3.6) and maximizes only the first term with the hinge loss function. However, the effect of ignoring the second term is not negligible when the class prior, θ_P , is not sufficiently small. In contrast, the PU-AUC risk includes the second term so that it is equivalent to its supervised counterpart, i.e., the PN-AUC risk.

For the next chapter, we will present the symmetry of the PU-AUC risk, that is, the risk in AUC optimization from negative and unlabeled data (the NU-AUC risk). From the definition of the marginal density, we have

$$\begin{aligned} & \mathbb{E}_U[\mathbb{E}_N[\ell(g(\mathbf{x}^U) - g(\mathbf{x}^N))]] \\ &= \theta_P \mathbb{E}_P[\mathbb{E}_N[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^N))]] + \theta_N \mathbb{E}_N[\mathbb{E}_{\bar{N}}[\ell(g(\mathbf{x}^N) - g(\bar{\mathbf{x}}^N))]] \\ &= \theta_P R_{PN}^{\text{AUC}}(g) + \theta_N \mathbb{E}_N[\mathbb{E}_{\bar{N}}[\ell(g(\mathbf{x}^N) - g(\bar{\mathbf{x}}^N))]], \end{aligned} \quad (3.8)$$

where $\mathbb{E}_{\bar{N}}$ denotes the expectation over $p(\bar{\mathbf{x}}^N \mid y = -1)$. Rearranging the above equation, we can obtain the PN-AUC risk in Eq. (2.23) based on negative and unlabeled data (the NU-AUC risk):

$$\begin{aligned} R_{PN}^{\text{AUC}}(g) &= \frac{1}{\theta_P} \mathbb{E}_U[\mathbb{E}_N[\ell(g(\mathbf{x}^U) - g(\mathbf{x}^N))]] \\ &\quad - \frac{\theta_N}{\theta_P} \mathbb{E}_N[\mathbb{E}_{\bar{N}}[\ell(g(\mathbf{x}^N) - g(\bar{\mathbf{x}}^N))]] := R_{NU}^{\text{AUC}}(g). \end{aligned} \quad (3.9)$$

We refer to the method minimizing the NU-AUC risk as *NU-AUC optimization*. Finally, we have the following lemma.

Lemma 4. *The AUC, PN-AUC, PU-AUC, and NU-AUC risks are equivalent:*

$$R^{\text{AUC}}(g) = R_{PN}^{\text{AUC}}(g) = R_{PU}^{\text{AUC}}(g) = R_{NU}^{\text{AUC}}(g). \quad (3.10)$$

This property of equivalence will be utilized for developing semi-supervised learning methods in Chapter 5.

3.4 Practical Implementation

In this section, we present the implementation details of the PU-AUC optimization method.

3.4.1 General Case

In practice, the AUC risk R^{AUC} is replaced with its empirical version \hat{R}^{AUC} , i.e., the expectations in R^{AUC} are replaced with the corresponding sample averages.

Here, we focus on the linear-in-parameter model given by

$$g(\mathbf{x}) = \sum_{\ell=1}^b w_{\ell} \phi_{\ell}(\mathbf{x}) = \mathbf{w}^{\top} \boldsymbol{\phi}(\mathbf{x}), \quad (3.11)$$

where $^{\top}$ denotes the transpose of vectors and matrices, b is the number of basis functions, $\mathbf{w} = (w_1, \dots, w_b)^{\top}$ is a parameter vector, and $\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_b(\mathbf{x}))^{\top}$ is a basis function vector. The linear-in-parameter model allows us to express the difference of two classifiers also in a linear-in-parameter form as

$$g(\mathbf{x}) - g(\mathbf{x}') = \mathbf{w}^{\top} (\boldsymbol{\phi}(\mathbf{x}) - \boldsymbol{\phi}(\mathbf{x}')) = \mathbf{w}^{\top} \bar{\boldsymbol{\phi}}(\mathbf{x}, \mathbf{x}'), \quad (3.12)$$

where

$$\bar{\boldsymbol{\phi}}(\mathbf{x}, \mathbf{x}') := \boldsymbol{\phi}(\mathbf{x}) - \boldsymbol{\phi}(\mathbf{x}') \quad (3.13)$$

is a composite basis function vector. This property makes the implementation simple with certain loss functions. Then, we train the classifier by minimizing the ℓ_2 -regularized empirical AUC risk:

$$\min_{\mathbf{w}} \left[\hat{R}^{\text{AUC}}(g) + \lambda \|\mathbf{w}\|^2 \right], \quad (3.14)$$

where $\lambda \geq 0$ is the regularization parameter.

3.4.2 Analytical Solution for Squared Loss

For the squared loss $\ell_S(m) := (1 - m)^2$, the empirical PU-AUC risk (See Section 3.7 for the derivation) can be expressed as

$$\begin{aligned}\widehat{R}_{\text{PU}}^{\text{AUC}}(g) &= \frac{1}{\theta_N n_P n_U} \sum_{i=1}^{n_P} \sum_{k=1}^{n_U} \ell_S(\mathbf{w}^\top \bar{\boldsymbol{\phi}}(\mathbf{x}_i^P, \mathbf{x}_k^U)) \\ &\quad - \frac{\theta_P}{\theta_N n_P (n_P - 1)} \sum_{i=1}^{n_P} \sum_{i'=1}^{n_P} \ell_S(\mathbf{w}^\top \bar{\boldsymbol{\phi}}(\mathbf{x}_i^P, \mathbf{x}_{i'}^P)) + \frac{\theta_P}{\theta_N (n_P - 1)} \\ &= 1 - 2\mathbf{w}^\top \widehat{\mathbf{h}}_{\text{PU}}^{\text{AUC}} + \mathbf{w}^\top \widehat{\mathbf{H}}_{\text{PU}}^{\text{AUC}} \mathbf{w} - \mathbf{w}^\top \widehat{\mathbf{H}}_{\text{PP}}^{\text{AUC}} \mathbf{w},\end{aligned}\quad (3.15)$$

where

$$\widehat{\mathbf{h}}_{\text{PU}}^{\text{AUC}} := \frac{1}{\theta_N n_P} \boldsymbol{\Phi}_P^\top \mathbf{1}_{n_P} - \frac{1}{\theta_N n_U} \boldsymbol{\Phi}_U^\top \mathbf{1}_{n_U}, \quad (3.16)$$

$$\begin{aligned}\widehat{\mathbf{H}}_{\text{PU}}^{\text{AUC}} &:= \frac{1}{\theta_N n_P} \boldsymbol{\Phi}_P^\top \boldsymbol{\Phi}_P - \frac{1}{\theta_N n_P n_U} \boldsymbol{\Phi}_U^\top \mathbf{1}_{n_U} \mathbf{1}_{n_P}^\top \boldsymbol{\Phi}_P \\ &\quad - \frac{1}{\theta_N n_P n_U} \boldsymbol{\Phi}_P^\top \mathbf{1}_{n_P} \mathbf{1}_{n_U}^\top \boldsymbol{\Phi}_U + \frac{1}{\theta_N n_U} \boldsymbol{\Phi}_U^\top \boldsymbol{\Phi}_U,\end{aligned}\quad (3.17)$$

$$\widehat{\mathbf{H}}_{\text{PP}}^{\text{AUC}} := \frac{2\theta_P}{\theta_N (n_P - 1)} \boldsymbol{\Phi}_P^\top \boldsymbol{\Phi}_P - \frac{2\theta_P}{\theta_N n_P (n_P - 1)} \boldsymbol{\Phi}_P^\top \mathbf{1}_{n_P} \mathbf{1}_{n_P}^\top \boldsymbol{\Phi}_P, \quad (3.18)$$

$$\boldsymbol{\Phi}_P := (\boldsymbol{\phi}(\mathbf{x}_1^P), \dots, \boldsymbol{\phi}(\mathbf{x}_{n_P}^P))^\top, \quad (3.19)$$

$$\boldsymbol{\Phi}_U := (\boldsymbol{\phi}(\mathbf{x}_1^U), \dots, \boldsymbol{\phi}(\mathbf{x}_{n_U}^U))^\top, \quad (3.20)$$

and $\mathbf{1}_b$ is the b -dimensional vector whose elements are all one. With the ℓ_2 -regularizer, we can analytically obtain the solution¹ by

$$\widehat{\mathbf{w}}_{\text{PU}}^{\text{AUC}} := (\widehat{\mathbf{H}}_{\text{PU}}^{\text{AUC}} - \widehat{\mathbf{H}}_{\text{PP}}^{\text{AUC}} + \lambda \mathbf{I}_b)^{-1} \widehat{\mathbf{h}}_{\text{PU}}^{\text{AUC}}, \quad (3.21)$$

where \mathbf{I}_b is the b -dimensional identity matrix.

The computational complexity of computing $\widehat{\mathbf{h}}_{\text{PU}}^{\text{AUC}}$, $\widehat{\mathbf{H}}_{\text{PU}}^{\text{AUC}}$, and $\widehat{\mathbf{H}}_{\text{PP}}^{\text{AUC}}$ are $\mathcal{O}((n_P + n_U)b)$, $\mathcal{O}((n_P + n_U)b^2)$, and $\mathcal{O}(n_P b^2)$, respectively. Then, solving a system of linear equations to obtain the solution $\widehat{\mathbf{w}}_{\text{PU}}^{\text{AUC}}$ requires the computational complexity of $\mathcal{O}(b^3)$.

¹To be precise, the obtained solution is a critical point (stationary point) of the optimization problem. However, with the squared loss, since the first term in Eq. (3.6) and ℓ_2 -regularizer are strongly convex functions, the optimization problem will become convex if the regularization parameter λ is sufficiently large.

In total, the computational complexity of this PU-AUC optimization method is $\mathcal{O}((n_P + n_U)b^2 + b^3)$.

From the viewpoint of computational complexity, the squared loss and the exponential loss are more efficient than the logistic loss because these loss functions reduce the nested summations to individual ones. More specifically, for example, in the PU-AUC optimization method, the logistic loss requires $\mathcal{O}(n_P n_U)$ operations for evaluating the first term in the PU-AUC risk, i.e., the loss over positive and unlabeled samples. In contrast, the squared loss and exponential loss reduce the number of operations for loss evaluation to $\mathcal{O}(n_P + n_U)$. For example, the exponential loss over positive and unlabeled data can be computed as follows:

$$\begin{aligned} \sum_{i=1}^{n_P} \sum_{k=1}^{n_U} \ell_E(g(\mathbf{x}_i^P) - g(\mathbf{x}_k^U)) &= \sum_{i=1}^{n_P} \sum_{k=1}^{n_U} \exp(-g(\mathbf{x}_i^P) + g(\mathbf{x}_k^U)) \\ &= \sum_{i=1}^{n_P} \exp(-g(\mathbf{x}_i^P)) \sum_{k=1}^{n_U} \exp(g(\mathbf{x}_k^U)). \end{aligned} \quad (3.22)$$

Thus, the number of operations for loss evaluation is reduced to $n_P + n_U + 1$ rather than $n_P n_U$. This property is beneficial especially when we handle large scale datasets.

3.4.3 Model Selection

In practice, the performance of the PU-AUC optimization method depends on hyperparameters. To tune the hyperparameters, cross-validation is available. Although labeled data for all the classes are not available unlike supervised learning, we can conduct cross-validation by only using positive and unlabeled data.

Specifically, we first split the positive and unlabeled samples, \mathcal{S}^P and \mathcal{S}^U , into K disjoint subsets, $\{\mathcal{S}_k^P\}_{k=1}^K$ and $\{\mathcal{S}_k^U\}_{k=1}^K$ of almost the same size; we construct K disjoint subsets $\{\mathcal{S}_k := \mathcal{S}_k^P \cup \mathcal{S}_k^U\}_{k=1}^K$. Then, we obtain the classifier \hat{g}_k with training samples $\mathcal{S} \setminus \mathcal{S}_k$. The approximation error is computed by the hold-out samples \mathcal{S}_k . We repeat this procedure for $k = 1, \dots, K$, and compute its mean \hat{J}_{PU}^{CV} . Finally, we choose the best model that minimizes \hat{J}_{PU}^{CV} from all model candidates.

3.5 Generalization Error Bounds

As the classifier g , we assume the linear-in-parameter model. Let \mathcal{G} be a function class of bounded hyperplanes:

$$\mathcal{G} := \{g(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}) \mid \|\mathbf{w}\| \leq C_w; \forall \mathbf{x}: \|\boldsymbol{\phi}(\mathbf{x})\| \leq C_\phi\}, \quad (3.23)$$

where $C_w > 0$ and $C_\phi > 0$ are certain positive constants. This assumption is reasonable because the ℓ_2 -regularizer included in training and the use of bounded basis functions, e.g., the Gaussian kernel basis, ensure that the minimizer of the empirical AUC risk belongs to such function class \mathcal{G} . Furthermore, we assume that a surrogate loss is bounded from above by C_ℓ and denote the Lipschitz constant by L . For simplicity, we focus on a surrogate loss satisfying $\ell_{0.1}(m) \leq \ell(m)$, but the theoretical analysis can be easily extended to the loss satisfying $\ell_{0.1}(m) \leq M\ell(m)$ with a certain $M > 0$. For example, the squared loss and the exponential loss satisfy the condition $M = 1$. Moreover, these losses are bounded in our setting, since the input to $\ell(m)$, i.e., g is bounded.

For the sake of simplicity, we assume that another set of positive samples $\{\bar{\mathbf{x}}_{i'}^P\}_{i'=1}^{n_P}$ is available for computing the nested expectations over $p(\mathbf{x}^P \mid y = +1)$ and $p(\bar{\mathbf{x}}^P \mid y = +1)$ in the second term of Eq. (3.6), unlike the practical implementation. On the other hand, for approximating the PU-AUC risk function, we can split positive samples $\{\mathbf{x}_i^P\}_{i=1}^{n_P}$ into two parts, e.g., $\{\mathbf{x}_i^P\}_{i=1}^{n_P/2}$ and $\{\bar{\mathbf{x}}_{i'}^P\}_{i'=1}^{n_P/2} = \{\mathbf{x}_i^P\}_{i=n_P/2+1}^{n_P}$ provided that n_P is even. Approximating the risk function in the latter way loosens the upper bound compared with the one obtained by the former approach. In our implementation (see Sections 3.4 and 3.7), we approximate the nested expectations without largely reducing the number of samples, unlike the latter approach. Thus, the practical performance of our risk estimator might be in between the former and the later approaches.

Let

$$I^{\text{AUC}}(g) = \mathbb{E}_P[\mathbb{E}_N[\ell_{0.1}(g(\mathbf{x}^P) - g(\mathbf{x}^N))]] \quad (3.24)$$

be the generalization error of g in AUC optimization. For convenience, we define

$$c^{\text{AUC}}(\delta) := 2\sqrt{2}LC_\ell C_w C_\phi + \frac{3}{2}\sqrt{2\log(2/\delta)}. \quad (3.25)$$

The following generalization error bounds are proved for the PU-AUC/NU-AUC risks (its proof is available in Section 3.8):

Theorem 5. For any $\delta > 0$, the following inequalities hold separately with probability at least $1 - \delta$ for all $g \in \mathcal{G}$:

$$I^{\text{AUC}}(g) \leq \hat{R}_{\text{PU}}^{\text{AUC}}(g) + c^{\text{AUC}}(\delta/2) \left(\frac{1}{\theta_{\text{N}} \sqrt{\min(n_{\text{P}}, n_{\text{U}})}} + \frac{\theta_{\text{P}}}{\theta_{\text{N}} \sqrt{n_{\text{P}}}} \right), \quad (3.26)$$

$$I^{\text{AUC}}(g) \leq \hat{R}_{\text{NU}}^{\text{AUC}}(g) + c^{\text{AUC}}(\delta/2) \left(\frac{1}{\theta_{\text{P}} \sqrt{\min(n_{\text{N}}, n_{\text{U}})}} + \frac{\theta_{\text{N}}}{\theta_{\text{P}} \sqrt{n_{\text{N}}}} \right), \quad (3.27)$$

where $\hat{R}_{\text{PU}}^{\text{AUC}}$ and $\hat{R}_{\text{NU}}^{\text{AUC}}$ are unbiased empirical risk estimators corresponding to $R_{\text{PU}}^{\text{AUC}}$ and $R_{\text{NU}}^{\text{AUC}}$, respectively.

Theorem 5 guarantees that $I^{\text{AUC}}(g)$ can be bounded from above by the empirical risk, $\hat{R}^{\text{AUC}}(g)$, plus the confidence terms of order

$$\mathcal{O}_p\left(\frac{1}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}}\right) \text{ and } \mathcal{O}_p\left(\frac{1}{\sqrt{n_{\text{N}}}} + \frac{1}{\sqrt{n_{\text{U}}}}\right).$$

Since n_{P} (n_{N}) and n_{U} can increase independently in our setting, this is the optimal convergence rate without any additional assumptions (Vapnik, 1998; Mendelson, 2008).

3.6 Experiments

In this section, we compare the PU-AUC optimization method against the existing AUC optimization method based on the ranking SVM (PU-RSVM) (Sundararajan et al., 2011). We used 8 benchmark datasets from the *IDA Benchmark Repository* (Rätsch et al., 2001), the *LIBSVM* (Chang and Lin, 2011), and the *UCI Machine Learning Repository* (Lichman, 2013).

As the classifier, we used the linear-in-parameter model with the Gaussian kernel basis function:

$$\phi_{\ell}(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_{\ell}\|^2}{2\sigma^2}\right),$$

where $\sigma > 0$ is the Gaussian bandwidth, $\{\mathbf{x}_{\ell}\}_{\ell=1}^b$ are the samples randomly selected from unlabeled samples $\{\mathbf{x}_i^{\text{U}}\}_{i=1}^{n_{\text{U}}}$. The number of basis functions was set at $b = \min(n_{\text{U}}, 200)$. The candidates of the Gaussian bandwidth were

$$\text{median}(\{\|\mathbf{x}_i - \mathbf{x}_j\|\}_{i,j=1}^n) \times \left\{ \frac{1}{8}, \frac{1}{4}, \frac{1}{2}, 1, 2 \right\} \quad (3.28)$$

TABLE 3.1: Average and standard error of the estimated class-prior over 50 trials on benchmark datasets in PU learning setting.

Dataset	d	$\theta_P = 0.1$	$\theta_P = 0.2$
Banana	2	0.15 (0.01)	0.24 (0.01)
skin_nonskin	3	0.10 (0.00)	0.20 (0.01)
cod-rna	8	0.32 (0.01)	0.40 (0.02)
Magic	10	0.34 (0.01)	0.39 (0.01)
Image	18	0.13 (0.01)	0.24 (0.01)
Twonorm	20	0.27 (0.00)	0.34 (0.00)
Waveform	21	0.31 (0.00)	0.38 (0.01)
mushrooms	112	0.19 (0.00)	0.28 (0.00)

and those of the regularization parameter were $\{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$. The bandwidth and regularization parameter were determined by five-fold cross-validation. We trained a classifier with samples of size $n_P = 100$ and $n_U = 1000$ under the different class-priors $\theta_P = 0.1$ and 0.2 . For the PU-AUC optimization method, the squared loss function was used and the class-prior was estimated by the distribution matching method (du Plessis et al., 2017). The results of the estimated class-prior are summarized in Table 3.1.

Table 3.2 lists the average with standard error of the AUC over 50 trials, showing that the PU-AUC optimization method achieves larger AUC than PU-RSVM. In particular, when $\theta_P = 0.2$, the difference between PU-RSVM and our method becomes larger compared with the difference when $\theta_P = 0.1$. This is because PU-RSVM can be interpreted as regarding unlabeled data as negative; it causes larger bias when θ_P increases.

Figure 3.1 summarizes the average computation time over 50 trials. The computation time of the PU-AUC optimization method includes both the class-prior estimation and the empirical risk minimization. The results show that the PU-AUC optimization method requires almost twice computation time as that of PU-RSVM, but it would be acceptable in practice to obtain better performance.

TABLE 3.2: Average and standard error of the AUC over 50 trials on benchmark datasets. The boldface denotes the best and comparable methods in terms of the average AUC according to the t-test at the significance level 5%. The bottom row shows the number of best/comparable cases of each method.

Dataset	d	$\theta_P = 0.1$		$\theta_P = 0.2$	
		PU-AUC	PU-RSVM	PU-AUC	PU-RSVM
Banana	2	95.4 (0.1)	88.5 (0.2)	95.0 (0.1)	85.9 (0.2)
skin_nonskin	3	99.9 (0.0)	86.2 (0.3)	99.8 (0.0)	73.8 (0.5)
cod-rna	8	98.1 (0.1)	62.8 (0.3)	97.6 (0.1)	59.4 (0.4)
Magic	10	87.4 (0.2)	82.8 (0.2)	86.4 (0.2)	81.9 (0.1)
Image	18	97.5 (0.1)	82.2 (0.2)	96.7 (0.2)	77.0 (0.4)
Twonorm	20	99.6 (0.0)	85.5 (0.2)	99.5 (0.0)	79.2 (0.4)
Waveform	21	96.6 (0.1)	73.9 (0.6)	96.3 (0.1)	59.5 (1.0)
mushrooms	112	99.8 (0.0)	93.9 (0.3)	99.6 (0.1)	84.7 (0.3)
#Best/Comp.		8	0	8	0

3.7 Derivation of PU-AUC Risk Estimator

Here, we present the derivation of the PU-AUC risk estimator. Recall that the PU-AUC risk in Eq. (3.6) is defined as

$$R_{\text{PU}}^{\text{AUC}}(g) = \frac{1}{\theta_N} \mathbb{E}_P[\mathbb{E}_U[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^U))]] - \frac{\theta_P}{\theta_N} \mathbb{E}_P[\mathbb{E}_{\bar{P}}[\ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P))]]. \quad (3.29)$$

If one additional set of positive samples $\{\bar{\mathbf{x}}_i^P\}_{i=1}^{n_P}$ is available, we obtain the unbiased PU-AUC risk estimator by

$$\hat{R}_{\text{PU}}^{\text{AUC}}(g) = \frac{1}{\theta_N} \sum_{i=1}^{n_P} \sum_{k=1}^{n_U} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_k^U)) - \frac{\theta_P}{\theta_N n_P^2} \sum_{i=1}^{n_P} \sum_{i'=1}^{n_P} \ell(g(\mathbf{x}_i^P) - g(\bar{\mathbf{x}}_{i'}^P)). \quad (3.30)$$

We used this estimator in our theoretical analyses because learning is not involved. However, obtaining one additional set of samples is not always possible in practice. Thus,

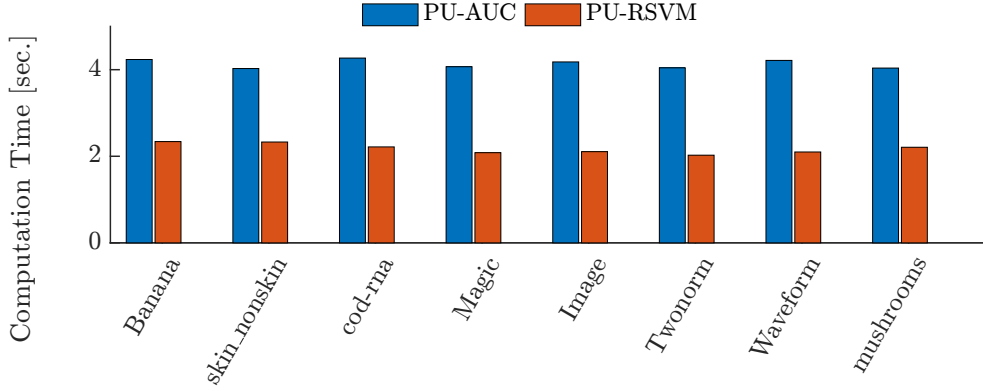


FIGURE 3.1: Average computation time on benchmark datasets when $\theta_P = 0.1$ over 50 trials. The computation time of the PU-AUC optimization method includes the class-prior estimation and the empirical risk minimization.

instead of the above risk estimator, we use the following risk estimator in our implementation:

$$\begin{aligned} \hat{R}_{\text{PU}}^{\text{AUC}}(g) &= \frac{1}{\theta_N} \sum_{i=1}^{n_P} \sum_{k=1}^{n_U} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_k^U)) \\ &\quad - \frac{\theta_P}{\theta_N} \left(\frac{1}{n_P(n_P - 1)} \sum_{i=1}^{n_P} \sum_{i'=1}^{n_P} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_{i'}^P)) - \frac{\ell(0)}{n_P - 1} \right). \end{aligned} \quad (3.31)$$

This estimator is also unbiased. To show unbiasedness of this estimator, let us rewrite the second term of the PU-AUC risk without coefficient in Eq. (3.6) as

$$\mathbb{E}_P[\mathbb{E}_{\bar{P}}[\ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P))]] = \mathbb{E}_{\mathbf{x}^P, \bar{\mathbf{x}}^P}[\ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P))]. \quad (3.32)$$

The unbiased estimator can be expressed as

$$\frac{1}{n_P(n_P - 1)} \sum_{i=1}^{n_P} \sum_{i'=1}^{n_P} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_{i'}^P)) - \frac{\ell(0)}{n_P - 1}, \quad (3.33)$$

because the expectation of the above estimator can be computed as follows:

$$\begin{aligned}
& \mathbb{E}_{\mathbf{x}_1^P, \dots, \mathbf{x}_{n_P}^P} \left[\frac{1}{n_P(n_P - 1)} \sum_{i=1}^{n_P} \sum_{i'=1}^{n_P} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_{i'}^P)) - \frac{\ell(0)}{n_P - 1} \right] \\
&= \mathbb{E}_{\mathbf{x}_1^P, \dots, \mathbf{x}_{n_P}^P} \left[\frac{1}{n_P(n_P - 1)} \left(\sum_{i=1}^{n_P} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_i^P)) + \sum_{i=1}^{n_P} \sum_{i' \neq i}^{n_P} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_{i'}^P)) \right) \right. \\
&\quad \left. - \frac{\ell(0)}{n_P - 1} \right] \\
&= \mathbb{E}_{\mathbf{x}_1^P, \dots, \mathbf{x}_{n_P}^P} \left[\frac{1}{n_P(n_P - 1)} \left(\sum_{i=1}^{n_P} \ell(0) + \sum_{i=1}^{n_P} \sum_{i' \neq i}^{n_P} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_{i'}^P)) \right) - \frac{\ell(0)}{n_P - 1} \right] \\
&= \frac{1}{n_P(n_P - 1)} \sum_{i=1}^{n_P} \sum_{i' \neq i}^{n_P} \mathbb{E}_{\mathbf{x}_i^P, \mathbf{x}_{i'}^P} \left[\ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_{i'}^P)) \right] \\
&= \frac{1}{n_P(n_P - 1)} \sum_{i=1}^{n_P} \sum_{i' \neq i}^{n_P} \mathbb{E}_{\mathbf{x}^P, \bar{\mathbf{x}}^P} \left[\ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P)) \right] \\
&= \mathbb{E}_{\mathbf{x}^P, \bar{\mathbf{x}}^P} \left[\ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P)) \right]. \tag{3.34}
\end{aligned}$$

If the squared loss function $\ell(m) = (1 - m)^2$ is used, $\ell(0) = 1$ (cf. the implementation with the squared loss in Section 3.4.2). Therefore, the PU-AUC risk estimator is unbiased.

3.8 Proof of Theorem 5

Finally, we show the proof of generalization error bounds in Section 3.5. The proof is based on [Usunier et al. \(2006\)](#).

Let

$$\{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} q(\mathbf{x}) \quad \text{and} \tag{3.35}$$

$$\{\mathbf{x}'_j\}_{j=1}^{n'} \stackrel{\text{i.i.d.}}{\sim} q'(\mathbf{x}) \tag{3.36}$$

be two sets of samples drawn from the distribution equipped with densities $q(\mathbf{x})$ and $q'(\mathbf{x})$, respectively. Let \mathcal{G} be a function class of bounded hyperplanes:

$$\mathcal{G} := \{g(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) \mid \|\mathbf{w}\| \leq C_w; \forall \mathbf{x}: \|\phi(\mathbf{x})\| \leq C_\phi\}, \tag{3.37}$$

where $C_w > 0$ and $C_\phi > 0$ are certain positive constants. The AUC risk over distributions q and q' and its empirical version can be expressed as

$$R^{\text{AUC}}(g) := \mathbb{E}_{\mathbf{x} \sim q}[\mathbb{E}_{\mathbf{x}' \sim q'}[\ell(g(\mathbf{x}) - g(\mathbf{x}'))]], \quad (3.38)$$

$$\hat{R}^{\text{AUC}}(g) := \frac{1}{nn'} \sum_{i=1}^n \sum_{j=1}^{n'} \ell(g(\mathbf{x}_i) - g(\mathbf{x}'_j)). \quad (3.39)$$

For convenience, let us define

$$c^{\text{AUC}}(\delta) := 2\sqrt{2}LC_\ell C_w C_\phi + \frac{3}{2}\sqrt{2\log(2/\delta)}. \quad (3.40)$$

Firstly, we have the following theorem:

Theorem 6. *For any $\delta > 0$, the following inequality holds with probability at least $1 - \delta$ for any $g \in \mathcal{G}$:*

$$R^{\text{AUC}}(g) - \hat{R}^{\text{AUC}}(g) \leq c^{\text{AUC}}(\delta) \frac{1}{\sqrt{\min(n, n')}}. \quad (3.41)$$

Proof. By slightly modifying Theorem 7 in [Usunier et al. \(2006\)](#) to fit our setting, for any $\delta > 0$, with probability at least $1 - \delta$ for any $g \in \mathcal{G}$, we have

$$\begin{aligned} R^{\text{AUC}}(g) - \hat{R}^{\text{AUC}}(g) &\leq \frac{2LC_\ell C_w \sqrt{\max(n, n')}}{nn'} \sqrt{\sum_{i=1}^n \sum_{j=1}^{n'} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}'_j)\|^2} \\ &\quad + 3\sqrt{\frac{\log(2/\delta)}{2\min(n, n')}}. \end{aligned} \quad (3.42)$$

Then, applying the triangle inequality to the first term in Eq. (3.42), we have

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^{n'} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}'_j)\|^2 &\leq n' \sum_{i=1}^n \|\phi(\mathbf{x}_i)\|^2 + n \sum_{j=1}^{n'} \|\phi(\mathbf{x}'_j)\|^2 \\ &\leq 2nn' C_\phi^2. \end{aligned} \quad (3.43)$$

Plugging the above into Eq. (3.42), we conclude the proof of the theorem. \square

Next, we prove the risk bounds of the PU-AUC and NU-AUC risks. For the sake of simplicity, we assume that the expectations of \mathbb{E}_P (\mathbb{E}_N) and $\mathbb{E}_{\bar{P}}$ ($\mathbb{E}_{\bar{N}}$) in the PU-AUC (NU-AUC) risks are estimated by using the same number of samples n_P (n_N), i.e., $\{\mathbf{x}_i^P\}_{i=1}^{n_P}$

and $\{\bar{\mathbf{x}}_{i'}^{\mathbf{P}}\}_{i'=1}^{n_{\mathbf{P}}}$ are used for approximating $\mathbb{E}_{\mathbf{P}}$ and $\mathbb{E}_{\bar{\mathbf{P}}}$.² By using Theorem 6, we have the following lemma:

Lemma 7. *For any $\delta > 0$, the following inequalities hold separately with probability at least $1 - \delta$ for any $g \in \mathcal{G}$:*

$$R_{\mathbf{PU}}^{\text{AUC}}(g) - \widehat{R}_{\mathbf{PU}}^{\text{AUC}}(g) \leq c^{\text{AUC}}(\delta/2) \left(\frac{1}{\theta_{\mathbf{N}} \sqrt{\min(n_{\mathbf{P}}, n_{\mathbf{U}})}} + \frac{\theta_{\mathbf{P}}}{\theta_{\mathbf{N}} \sqrt{n_{\mathbf{P}}}} \right), \quad (3.44)$$

$$R_{\mathbf{NU}}^{\text{AUC}}(g) - \widehat{R}_{\mathbf{NU}}^{\text{AUC}}(g) \leq c^{\text{AUC}}(\delta/2) \left(\frac{1}{\theta_{\mathbf{P}} \sqrt{\min(n_{\mathbf{N}}, n_{\mathbf{U}})}} + \frac{\theta_{\mathbf{N}}}{\theta_{\mathbf{P}} \sqrt{n_{\mathbf{N}}}} \right). \quad (3.45)$$

Proof. Recall that the PU-AUC and NU-AUC risks are expressed as

$$R_{\mathbf{PU}}^{\text{AUC}}(g) = \frac{1}{\theta_{\mathbf{N}}} \mathbb{E}_{\mathbf{P}}[\mathbb{E}_{\mathbf{U}}[\ell(g(\mathbf{x}^{\mathbf{P}}) - g(\mathbf{x}^{\mathbf{U}}))]] - \frac{\theta_{\mathbf{P}}}{\theta_{\mathbf{N}}} \mathbb{E}_{\mathbf{P}}[\mathbb{E}_{\bar{\mathbf{P}}}[\ell(g(\mathbf{x}^{\mathbf{P}}) - g(\bar{\mathbf{x}}^{\mathbf{P}}))]], \quad (3.46)$$

$$R_{\mathbf{NU}}^{\text{AUC}}(g) = \frac{1}{\theta_{\mathbf{P}}} \mathbb{E}_{\mathbf{U}}[\mathbb{E}_{\mathbf{N}}[\ell(g(\mathbf{x}^{\mathbf{U}}) - g(\mathbf{x}^{\mathbf{N}}))]] - \frac{\theta_{\mathbf{N}}}{\theta_{\mathbf{P}}} \mathbb{E}_{\mathbf{N}}[\mathbb{E}_{\bar{\mathbf{N}}}[\ell(g(\mathbf{x}^{\mathbf{N}}) - g(\bar{\mathbf{x}}^{\mathbf{N}}))]]. \quad (3.47)$$

² Otherwise, we can consider the obtained samples are divided into two parts, e.g., $n_{\mathbf{P}}/2$, and the bound in Eq. (3.50) are loosen by two times.

From Theorem 6, for any $\delta > 0$, we have these uniform deviation bounds with probability at least $1 - \delta/2$:

$$\sup_{g \in \mathcal{G}} \left(\mathbb{E}_P[\mathbb{E}_U[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^U))]] - \frac{1}{n_P n_U} \sum_{i=1}^{n_P} \sum_{k=1}^{n_U} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_k^U)) \right) \leq \frac{c^{\text{AUC}}(\delta/2)}{\sqrt{\min(n_P, n_U)}}, \quad (3.48)$$

$$\sup_{g \in \mathcal{G}} \left(\mathbb{E}_U[\mathbb{E}_N[\ell(g(\mathbf{x}^U) - g(\mathbf{x}^N))]] - \frac{1}{n_N n_U} \sum_{k=1}^{n_U} \sum_{j=1}^{n_N} \ell(g(\mathbf{x}_k^U) - g(\mathbf{x}_j^N)) \right) \leq \frac{c^{\text{AUC}}(\delta/2)}{\sqrt{\min(n_N, n_U)}}, \quad (3.49)$$

$$\sup_{g \in \mathcal{G}} \left(\mathbb{E}_P[\mathbb{E}_{\bar{P}}[\ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P))]] - \frac{1}{n_P^2} \sum_{i=1}^{n_P} \sum_{i'=1}^{n_P} \ell(g(\mathbf{x}_i^P) - g(\bar{\mathbf{x}}_{i'}^P)) \right) \leq \frac{c^{\text{AUC}}(\delta/2)}{\sqrt{n_P}}, \quad (3.50)$$

$$\sup_{g \in \mathcal{G}} \left(\mathbb{E}_N[\mathbb{E}_{\bar{N}}[\ell(g(\mathbf{x}^N) - g(\bar{\mathbf{x}}^N))]] - \frac{1}{n_N^2} \sum_{j=1}^{n_N} \sum_{j'=1}^{n_N} \ell(g(\mathbf{x}_j^N) - g(\bar{\mathbf{x}}_{j'}^N)) \right) \leq \frac{c^{\text{AUC}}(\delta/2)}{\sqrt{n_N}}. \quad (3.51)$$

Then, simple calculation showed that for any $\delta > 0$, with probability $1 - \delta$, we have

$$\begin{aligned} & \sup_{g \in \mathcal{G}} \left(R_{\text{PU}}^{\text{AUC}}(g) - \hat{R}_{\text{PU}}^{\text{AUC}}(g) \right) \\ & \leq \frac{1}{\theta_N} \sup_{g \in \mathcal{G}} \left(\mathbb{E}_P[\mathbb{E}_U[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^U))]] - \frac{1}{n_P n_U} \sum_{i=1}^{n_P} \sum_{k=1}^{n_U} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_k^U)) \right) \\ & \quad + \frac{\theta_P}{\theta_N} \sup_{g \in \mathcal{G}} \left(\mathbb{E}_P[\mathbb{E}_{\bar{P}}[\ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P))]] - \frac{1}{n_P^2} \sum_{i=1}^{n_P} \sum_{i'=1}^{n_P} \ell(g(\mathbf{x}_i^P) - g(\bar{\mathbf{x}}_{i'}^P)) \right) \\ & \leq c^{\text{AUC}}(\delta/2) \left(\frac{1}{\theta_N \sqrt{\min(n_P, n_U)}} + \frac{\theta_P}{\theta_N \sqrt{n_P}} \right), \end{aligned} \quad (3.52)$$

where we used

$$\sup(x + y) \leq \sup(x) + \sup(y), \quad (3.53)$$

$$R_{\text{PU}}^{\text{AUC}}(g) \leq \frac{1}{\theta_N} \mathbb{E}_P[\mathbb{E}_U[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^U))]] + \frac{\theta_P}{\theta_N} \mathbb{E}_P[\mathbb{E}_{\bar{P}}[\ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P))]]. \quad (3.54)$$

Similarly, for the NU-AUC risk, we have

$$\sup_{g \in \mathcal{G}} \left(R_{\text{NU}}^{\text{AUC}}(g) - \widehat{R}_{\text{NU}}^{\text{AUC}}(g) \right) \leq c^{\text{AUC}}(\delta/2) \left(\frac{1}{\theta_{\text{P}} \sqrt{\min(n_{\text{N}}, n_{\text{U}})}} + \frac{\theta_{\text{N}}}{\theta_{\text{P}} \sqrt{n_{\text{N}}}} \right). \quad (3.55)$$

Eqs. (3.52) and (3.55) conclude the lemma. \square

Finally, we give the proof of Theorem 5.

Proof. Assume the loss satisfying $\ell_{0.1}(m) \leq M\ell(m)$. Then, we have $I^{\text{AUC}}(g) \leq M \cdot R^{\text{AUC}}(g)$. When $M = 1$ such as the squared loss $\ell_{\text{S}}(m)$ and the exponential loss $\ell_{\text{E}}(m)$, $I^{\text{AUC}}(g) \leq R^{\text{AUC}}(g)$ holds. This observation yields Theorem 5. \square

Chapter 4

Statistical Dependency Analysis from Positive and Unlabeled Data

In this chapter, we develop statistical dependency analysis tools for positive and unlabeled data.

4.1 Introduction

Capturing relation between input and output is an important task in statistical data analysis. To this end, a statistical dependency measure plays a central role. Among statistical dependency measures, based on the *information maximization principle*, *mutual information* (MI) (Cover and Thomas, 2006), in particular, its empirical estimator has been employed to solve various machine learning tasks (Torkkola, 2003; Krause et al., 2010). However, since the MI estimator is known to be sensitive to outliers (Basu et al., 1998), it limits the range of applications of the MI-based methods. To cope with this problem, *squared-loss MI* (SMI) and its empirical estimator were proposed (Suzuki et al., 2009). Unlike MI, the SMI estimator is robust to noise and outliers, and having superior numerical properties (Kanamori et al., 2009), allowing us to apply the estimator to various statistical dependency analysis (Sugiyama, 2013).

On the other hand, as discussed throughout this dissertation, it is conceivable that only positive and unlabeled data are available (the PU learning situation) in real-world applications. So far, attention in PU learning has mostly been paid to classification (Elkan and Noto, 2008; du Plessis et al., 2015), and there are few studies on statistical dependency analysis for PU data. Recently, a PU-based independence test was proposed (Sechidis et al., 2014), but its range of applications is limited since the method can only be applied

to a categorical univariate input. Furthermore, there seems no other statistical dependency analysis, such as dimension reduction, for PU data. To analyze PU data from many aspects, it is practically desirable to have statistical dependency analysis for PU data.

In this chapter, we discuss the problem of estimating SMI and SMI-based statistical data analyses under the PU learning setting. We first derive an SMI expression over PU data (PU-SMI), which is equivalent to its supervised counterpart, i.e., an SMI expression over PN data. Then, we develop an estimation method for PU-SMI and show that the PU-SMI estimator converges to true SMI at the optimal parametric rate. We apply the PU-SMI estimator to dimension reduction and independence test. In particular, dimension reduction can be executed without estimating the class-prior, which is indispensable for most of the existing PU classification methods. However, the class-prior estimation from PU data is known to be challenging (du Plessis and Sugiyama, 2014; Ramaswamy et al., 2016), especially for high-dimensional data. Thus, our dimension reduction can mitigate the issue by providing informative low-dimensional representations.

4.2 Problem Setting

Let $\mathbf{x} \in \mathbb{R}^d$ be an input pattern, $y \in \{\pm 1\}$ be a corresponding class label, and $p(\mathbf{x}, y)$ be the underlying joint density, where d is a positive integer. In PU learning, we cannot access negatively labeled data, but we can utilize unlabeled data. Let us define *positive* (P) and *unlabeled* (U) data as

$$\mathcal{X}_P := \{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = +1), \quad (4.1)$$

$$\mathcal{X}_U := \{\mathbf{x}_k^U\}_{k=1}^{n_U} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}) = \theta_P p(\mathbf{x} \mid y = +1) + \theta_N p(\mathbf{x} \mid y = -1), \quad (4.2)$$

where

$$\theta_P := p(y = +1) \text{ and } \theta_N := p(y = -1).$$

The goal is to develop statistical data analysis tools for PU data. To this end, we first show an SMI estimation method from PU data and apply the estimator to data analysis tasks.

4.3 SMI Estimation from Positive and Unlabeled Data

In this section, we first reformulate SMI with only positive and unlabeled data, which is equivalent to SMI with positive and negative data, and then describe an SMI estimation method.

4.3.1 SMI with Positive and Unlabeled Data (PU-SMI)

Let us express SMI in Eq. (2.94) as

$$\text{SMI} = \frac{\theta_P}{2} \int \left(\frac{p(\mathbf{x} | y = +1)}{p(\mathbf{x})} - 1 \right)^2 p(\mathbf{x}) d\mathbf{x} + \frac{\theta_N}{2} \int \left(\frac{p(\mathbf{x} | y = -1)}{p(\mathbf{x})} - 1 \right)^2 p(\mathbf{x}) d\mathbf{x}. \quad (4.3)$$

From the definition of the marginal density $p(\mathbf{x})$, we have

$$\begin{aligned} \theta_N \frac{p(\mathbf{x} | y = -1)}{p(\mathbf{x})} &= 1 - \theta_P \frac{p(\mathbf{x} | y = +1)}{p(\mathbf{x})}, \\ \theta_N \left(\frac{p(\mathbf{x} | y = -1)}{p(\mathbf{x})} - 1 \right) &= \theta_P \left(1 - \frac{p(\mathbf{x} | y = +1)}{p(\mathbf{x})} \right), \\ \left(\frac{p(\mathbf{x} | y = -1)}{p(\mathbf{x})} - 1 \right)^2 &= \frac{\theta_P^2}{\theta_N^2} \left(\frac{p(\mathbf{x} | y = +1)}{p(\mathbf{x})} - 1 \right)^2, \end{aligned} \quad (4.4)$$

where the equality between the first and second equations can be confirmed by using $\theta_P + \theta_N = 1$. Plugging the last equation into the second term of Eq. (4.3), we then obtain an expression of SMI only with positive and unlabeled data (PU-SMI) as

$$\text{SMI} = \frac{\theta_P}{2\theta_N} \int \left(\frac{p(\mathbf{x} | y = +1)}{p(\mathbf{x})} - 1 \right)^2 p(\mathbf{x}) d\mathbf{x} := \text{PU-SMI}. \quad (4.5)$$

This immediately leads to the following lemma.

Lemma 8. *SMI and PU-SMI are equivalent.*

Compared with the original definition of SMI, PU-SMI can be expressed as only the probability densities of positive and unlabeled data. In the next section, we explain the estimation method for PU-SMI.

4.3.2 PU-SMI Estimation

Next, we show an SMI estimation method from only positive and unlabeled data.

Let

$$s(\mathbf{x}) := \frac{p(\mathbf{x} \mid y = +1)}{p(\mathbf{x})} \quad (4.6)$$

be the density-ratio function and g be a model of s , and we learn s so that the squared error (the PU-SMI risk) between model g and true ratio s is minimized:

$$R_{\text{PU}}^{\text{SMI}}(g) := \frac{1}{2} \int \left(g(\mathbf{x}) - \frac{p(\mathbf{x} \mid y = +1)}{p(\mathbf{x})} \right)^2 p(\mathbf{x}) d\mathbf{x} - C_{\text{PU}} \quad (4.7)$$

$$= \frac{1}{2} \int g^2(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} - \int g(\mathbf{x}) p(\mathbf{x} \mid y = +1) d\mathbf{x}, \quad (4.8)$$

where $C_{\text{PU}} := \frac{1}{2} \int \frac{p^2(\mathbf{x} \mid y = +1)}{p(\mathbf{x})} d\mathbf{x}$. An empirical PU-SMI risk can be obtained by

$$\hat{R}_{\text{PU}}^{\text{SMI}}(g) := \frac{1}{2n_{\text{U}}} \sum_{k=1}^{n_{\text{U}}} g^2(\mathbf{x}_k^{\text{U}}) - \frac{1}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} g(\mathbf{x}_i^{\text{P}}). \quad (4.9)$$

Similarly to the existing SMI estimation method, we formulate our optimization problem as

$$\hat{g} := \underset{g}{\operatorname{argmin}} \left[\hat{R}_{\text{PU}}^{\text{SMI}}(g) + \lambda W(g) \right], \quad (4.10)$$

where $\lambda \geq 0$ is the regularization parameter.

With the obtained density-ratio estimator, a PU-SMI estimator can be given by

$$\widehat{\text{PU-SMI}} = \frac{\theta_{\text{P}}}{\theta_{\text{N}}} \left(\frac{1}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \hat{g}(\mathbf{x}_i^{\text{P}}) - \frac{1}{2n_{\text{U}}} \sum_{k=1}^{n_{\text{U}}} \hat{g}^2(\mathbf{x}_k^{\text{U}}) - \frac{1}{2} \right), \quad (4.11)$$

since PU-SMI can be expressed as

$$\begin{aligned} \text{PU-SMI} &= \frac{\theta_{\text{P}}}{\theta_{\text{N}}} \left(\frac{1}{2} \int s(\mathbf{x}) p(\mathbf{x} \mid y = +1) d\mathbf{x} - \frac{1}{2} \right) \\ &= \frac{\theta_{\text{P}}}{\theta_{\text{N}}} \left(\int s(\mathbf{x}) p(\mathbf{x} \mid y = +1) d\mathbf{x} - \frac{1}{2} \int s^2(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} - \frac{1}{2} \right), \end{aligned} \quad (4.12)$$

where the first equation can be obtained by expanding Eq. (4.5) and using $s(\mathbf{x})p(\mathbf{x}) = p(\mathbf{x} \mid y = +1)$, and the equality of the first and second equations can be confirmed by plugging $s(\mathbf{x})p(\mathbf{x}) = p(\mathbf{x} \mid y = +1)$ into the second term of the second equation.

Similarly to [Suzuki and Sugiyama \(2013\)](#), the minimization of the empirical PU-SMI

risk $\widehat{R}_{\text{PU}}^{\text{SMI}}$ corresponds to the maximization of $-\widehat{R}_{\text{PU}}^{\text{SMI}}$ which is a lower-bound of PU-SMI as in the f -divergence estimation (Keziou, 2003; Nguyen et al., 2007). That is, the PU-SMI estimator can also be obtained by maximizing the lower-bound:

$$\widehat{\text{PU-SMI}} = \max_g -\frac{\theta_P}{\theta_N} \left(\widehat{R}_{\text{PU}}^{\text{SMI}}(g) + \frac{1}{2} \right) = -\frac{\theta_P}{\theta_N} \left(\widehat{R}_{\text{PU}}^{\text{SMI}}(\widehat{g}) + \frac{1}{2} \right). \quad (4.13)$$

Based on this fact, we derive convergence results of the parameters of the density-ratio model and the PU-SMI estimator in Section 4.5.

The PU-SMI risk minimization for estimating density-ratio does not involve class-prior θ_P (and θ_N), and PU-SMI is proportional to class-prior ratio θ_P/θ_N , meaning that the class-prior estimation is not necessary as long as our purpose is maximizing (minimizing) PU-SMI. For instance, in SMI-based dimensionality reduction discussed in Section 4.6.1, the ratio of the class-prior ratio θ_P/θ_N can be safely ignored from the PU-SMI maximization. Since class-prior estimation only from PU data is known to be challenging (du Plessis and Sugiyama, 2014; Ramaswamy et al., 2016; Jain et al., 2016; du Plessis et al., 2017), being able to avoid class-prior estimation is a significant advantage of the SMI-based method in practice.

4.4 Practical Implementation

In this section, we explain a practical implementation of the PU-SMI estimation method.

4.4.1 Neural Network and Linear-in-Parameter Models

When we use a neural network model, a solution can be typically obtained by backpropagation (Goodfellow et al., 2016). To this end, various deep learning implementations (Jia et al., 2014; Abadi et al., 2015; Tokui et al., 2015; Chen et al., 2015) may be used for obtaining a learned model.

Another candidate of the density-ratio model is a linear-in-parameter model:

$$g(\mathbf{x}) = \sum_{\ell=1}^b w_{\ell} \phi_{\ell}(\mathbf{x}) = \mathbf{w}^{\top} \boldsymbol{\phi}(\mathbf{x}), \quad (4.14)$$

where $\mathbf{w} := (w_1, \dots, w_b)^{\top}$ is a vector of parameters, and $\boldsymbol{\phi}(\mathbf{x}) := (\phi_1(\mathbf{x}), \dots, \phi_b(\mathbf{x}))^{\top}$ is a vector of basis functions. This model allows us to obtain an analytic-form solution with the ℓ_2 -regularizer and a PU-SMI estimator, similarly to the existing SMI estimation

method. Furthermore, the optimal convergence is theoretically guaranteed as shown in Section 4.5.

More specifically, with the ℓ_2 -regularizer, the optimization problem is given by

$$\hat{\mathbf{w}} := \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^\top \widehat{\mathbf{H}}_U \mathbf{w} - \mathbf{w}^\top \hat{\mathbf{h}}_P + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}, \quad (4.15)$$

where $\lambda \geq 0$ is the regularization parameter and

$$\widehat{\mathbf{H}}_U := \frac{1}{n_U} \Phi_U^\top \Phi_U, \quad (4.16)$$

$$\hat{\mathbf{h}}_P := \frac{1}{n_P} \Phi_P^\top \mathbf{1}_{n_P}, \quad (4.17)$$

$$\Phi_P := (\phi(\mathbf{x}_1^P), \dots, \phi(\mathbf{x}_{n_P}^P))^\top, \quad (4.18)$$

$$\Phi_U := (\phi(\mathbf{x}_1^U), \dots, \phi(\mathbf{x}_{n_U}^U))^\top, \quad (4.19)$$

where $\mathbf{1}_b$ is the b -dimensional vector whose elements are all ones. The solution can be obtained analytically by differentiating the objective function with respect to \mathbf{w} and set it to zero:

$$\hat{\mathbf{w}} = (\widehat{\mathbf{H}}_U + \lambda \mathbf{I}_b)^{-1} \hat{\mathbf{h}}_P, \quad (4.20)$$

where \mathbf{I}_b is the $b \times b$ identity matrix. Finally, with the obtained estimator, we can compute an SMI approximator only from positive and unlabeled data:

$$\widehat{\text{PU-SMI}} = \frac{\theta_P}{\theta_N} \left(\hat{\mathbf{w}}^\top \hat{\mathbf{h}}_P - \frac{1}{2} \hat{\mathbf{w}}^\top \widehat{\mathbf{H}}_U \hat{\mathbf{w}} - \frac{1}{2} \right). \quad (4.21)$$

4.4.2 Model Selection

In practice, the performance of the PU-SMI estimator depends on hyperparameters such as the regularization parameter and the number of epochs for backpropagation.

One way is to prepare a validation dataset and tune hyperparameters based on it. With positive and unlabeled samples for validation, we compute the empirical PU-SMI risk in Eq. (4.9) and choose the best model that minimizes the empirical PU-SMI risk from all model candidates.

Another way is to use cross-validation. The procedure is the same as cross-validation for the PU-AUC optimization in Section 3.4.3. The notable difference is that cross-validation for the PU-SMI estimation does not require an estimate of the class-prior (see

the discussion in Section 4.3.2).

4.5 Theoretical Analysis

In this section, we derive the convergence rate of the parameter of the density-ratio model and the PU-SMI approximator based on the *perturbation analysis of optimization problems* (Bonnans and Cominetti, 1996; Bonnans and Shapiro, 1998).

Assume the linear-in-parameter model in Eq.(4.14) for the density-ratio $\frac{p(\mathbf{x}|y=+1)}{p(\mathbf{x})}$. Without loss of generality, we assume that the basis functions $\{\phi_\ell(\mathbf{x})\}_{\ell=1}^b$ are linearly independent over the marginal density $p(\mathbf{x}) > 0$ and satisfy $0 \leq \phi_\ell(\mathbf{x}) \leq 1$ for all $\ell = 1, \dots, b$ and $\mathbf{x} \in \mathbb{R}^d$. Moreover, we assume that there exists a constant M such that $\|\mathbf{w}\|_2 \leq M$, i.e., \mathbf{w} is regularized.

Let

$$\mathbf{w}^* := \underset{\mathbf{w}}{\operatorname{argmin}} R_{\text{PU}}^{\text{SMI}}(\mathbf{w}) \quad (4.22)$$

be the ideal estimate of the parameter of the density-ratio model. Let \mathcal{O}_p denote the order in probability. Then we have the following convergence result (its proof is given in Section 4.8):

Theorem 9. *As $n_P, n_U \rightarrow \infty$, we have*

$$\|\widehat{\mathbf{w}} - \mathbf{w}^*\|_2 = \mathcal{O}_p\left(\frac{1}{\sqrt{n_P}} + \frac{1}{\sqrt{n_U}}\right), \quad (4.23)$$

$$|\text{PU-SMI} - \widehat{\text{PU-SMI}}| = \mathcal{O}_p\left(\frac{1}{\sqrt{n_P}} + \frac{1}{\sqrt{n_U}}\right), \quad (4.24)$$

provided that $\lambda = \mathcal{O}_p(1/\sqrt{\min(n_P, n_U)})$.¹

Theorem 9 guarantees that the convergence of the density-ratio estimator and the PU-SMI estimator. In our setting, since n_P and n_U can increase independently, this is the optimal convergence rate without any additional assumption (Kanamori et al., 2009, 2012).

Note that Theorem 9 shows that both positive and unlabeled data contribute to convergence, implying that unlabeled data is directly used in the estimation rather than extracting the information of a data structure, such as the cluster structure frequently assumed in semi-supervised learning (Chapelle et al., 2006).

¹ The regularization parameter λ vanishes with the order of $\min(n_P, n_U)$, i.e., the convergence rate is dominated by the smaller size of positive and unlabeled data.

4.6 Application based on PU-SMI

In this section, we show two applications of PU-SMI, along the line of existing SMI-based applications.

4.6.1 Dimension Reduction

Here, we extend the existing SMI-based dimension reduction (Suzuki and Sugiyama, 2013), called *least-squares dimension reduction* (LSDR), to the PU learning setting.

Algorithm: Let $v: \mathbb{R}^d \rightarrow \mathbb{R}^m$, where $m < d$, be a mapping function from an input vector to its low-dimensional representation. Following the information-maximization principle (Linsker, 1988), we maximize PU-SMI with respect to the mapping function to find low-dimensional representation that maximally preserves dependency between input and output.

To obtain such a mapping function, we propose an alternating optimization procedure to maximize SMI—First, we estimate PU-SMI, and then we maximize the estimated PU-SMI with respect to the mapping function. More specifically, we solve the objective function in Eq. (4.9) to approximate PU-SMI while fixing the mapping function \hat{v} :

$$\min_g \hat{R}_{\text{PU}}^{\text{SMI}}(g \circ \hat{v}), \quad (4.25)$$

where \hat{v} is a current estimate of the mapping function and “ \circ ” denotes the function composition, i.e., $(g \circ v)(x) = g(v(x))$. Note that the input dimension of g is redefined as m , i.e., $g: \mathbb{R}^m \rightarrow \mathbb{R}$, accordingly. Then, we maximize the estimated PU-SMI with respect to v while fixing \hat{g} :

$$\max_v \widehat{\text{PU-SMI}}(\hat{g} \circ v), \quad (4.26)$$

where \hat{g} is a current estimate of the density ratio.

As discussed in Section 4.3.2, an advantage of our dimension reduction method is that the whole procedure can be executed without estimating the class-prior, unlike existing PU learning methods (Elkan and Noto, 2008; du Plessis et al., 2015). This can be easily confirmed by Eq. (4.13):

$$\operatorname{argmax}_v \widehat{\text{PU-SMI}}(\hat{g} \circ v) = \operatorname{argmax}_v -\hat{R}_{\text{PU}}^{\text{SMI}}(\hat{g} \circ v). \quad (4.27)$$

Algorithm 1 PUDR

Require: $\{\mathbf{x}_i^P\}_{i=1}^{n_P}$, $\{\mathbf{x}_k^U\}_{k=1}^{n_U}$, g , and \mathbf{v}

- 1: Initialize g and \mathbf{v}
- 2: **repeat**
- 3: Update $\hat{g} \leftarrow \hat{g} - \varepsilon_g \nabla_g \hat{R}_{PU}^{SMI}(g \circ \hat{\mathbf{v}})$
- 4: Update $\hat{\mathbf{v}} \leftarrow \hat{\mathbf{v}} - \varepsilon_v \nabla_v \hat{R}_{PU}^{SMI}(\hat{g} \circ \mathbf{v})$
- 5: **until** stopping conditions meet

Moreover, minimization of the squared error does not require an estimate of the class-prior. Thus, our dimension reduction method does not need for a class-prior estimation in advance.

We refer to the dimension reduction method for PU data as *PUDR*. A pseudo-code of PUDR is summarized in Algorithm 1.

Linear Transformation: A choice of a mapping function is to use a linear transformation similarly to the existing SMI-based dimension reduction method (Suzuki and Sugiyama, 2013). As in the existing method, we search a projection matrix $\mathbf{W} \in \mathbb{R}^{d \times m}$ that satisfies $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_m$ to avoid degenerated solutions. Note that if a linear transformation is used, the maximization of SMI corresponds to *sufficient dimension reduction* (Li, 1991).

Non-Linear Transformation: Another choice is to use non-linear transformation for capturing complex structures of the data. To obtain low-dimensional representation from a non-linear mapping, we use a neural network because we can use several useful techniques developed recently. Specifically, we regard $g \circ \mathbf{v}$ as a neural network. We set the last layer as a fully-connected layer, which corresponds to g and the parameter \mathbf{w} of the linear-in-parameter model. Layers before the last layer can be regarded as the basis function ϕ and it can be anything; we could use, for instance, a convolution layer or a pretrained network.

4.6.2 Independence Test

Here, we extend the existing SMI-based independence test, called *least-squares independence test* (LSIT), to the PU learning setting.

Algorithm 2 PUIT

Require: $\{\mathbf{x}_i^P\}_{i=1}^{n_P}$, $\{\mathbf{x}_k^U\}_{k=1}^{n_U}$, and T

- 1: Prepare a memory h of size T
- 2: Compute $\widehat{\text{PU-SMI}}$ with $\{\mathbf{x}_i^P\}_{i=1}^{n_P}$ and $\{\mathbf{x}_k^U\}_{k=1}^{n_U}$
- 3: $i \leftarrow 0$
- 4: **repeat**
- 5: Generate random labels $\{\tilde{y}_k\}_{k=1}^{n_U}$
- 6: Compute $\widetilde{\text{SMI}}$ with $\{(\mathbf{x}_k^U, \tilde{y}_k)\}_{k=1}^{n_U}$
- 7: Store $\widetilde{\text{SMI}}$ in $h[i]$
- 8: $i \leftarrow i + 1$
- 9: **until** $i \leq T$
- 10: Construct a distribution of $\widetilde{\text{SMI}}$ with h
- 11: Approximate p -value with $\widehat{\text{PU-SMI}}$ and the distribution of $\widetilde{\text{SMI}}$

In the independence test, we test whether a pair of random variables are independent or not. Since SMI takes zero if and only if two random variables are independent, we utilize the property to detect independence of variables. Similarly to Sugiyama and Suzuki (2011), we employ a permutation test (Efron and Tibshirani, 1994) for numerically computing a p -value for the test.

In the original SMI-based test, we first estimate SMI from paired samples; then we repeatedly compute SMI with samples whose labels are randomly shuffled to construct a distribution of the SMI approximator when two random variables are independent. However, unlike the fully supervised case, we do not have labels to be shuffled in our setting. To conduct the permutation test, we propose the following procedure. We assign random labels $\{\tilde{y}_k\}_{k=1}^{n_U}$ to unlabeled data $\{\mathbf{x}_k^U\}_{k=1}^{n_U}$ and compute $\widetilde{\text{SMI}}$ from the randomly labeled data $\{(\mathbf{x}_k^U, \tilde{y}_k)\}_{k=1}^{n_U}$. We compute $\widetilde{\text{SMI}}$ many times and then construct the distribution of $\widetilde{\text{SMI}}$. Then, we approximate a p -value from the constructed distribution by comparing the estimate of $\widehat{\text{PU-SMI}}$ computed from the obtained data. Finally, we test the null hypothesis (the pair of random variables are independent) with the obtained p -value. We refer to the above independence test for PU data as *PUIT*. A pseudo-code of PUIT is summarized in Algorithm 2.

Note that the class proportion of random labels $\{\tilde{y}_k\}_{k=1}^{n_U}$ is determined by the estimated class-prior from PU data so that the class-prior of randomly paired samples $\{(\mathbf{x}_k^U, \tilde{y}_k)\}_{k=1}^{n_U}$ is the same as that of the marginal distribution $p(\mathbf{x})$, i.e., $p(\mathbf{x} \mid y)p(y) = p(\mathbf{x})p(y)$.

4.7 Experiments

In this section, we experimentally investigate the behavior of the PU-SMI estimator and evaluate the performance of the dimension reduction and independence testing methods on various benchmark datasets: the MNIST (LeCun et al., 1998), the CIFAR10 (Torralba et al., 2008), the SVHN (Netzer et al., 2011) datasets, and several datasets obtained from the *LIBSVM* webpage (Chang and Lin, 2011).

4.7.1 Experimental Analysis

First, we investigate the estimation accuracy of the PU-SMI estimator.

We used the linear-in-parameter model with the Gaussian basis function

$$\phi_\ell(\mathbf{x}) := \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_\ell\|^2}{2\sigma^2}\right), \quad (4.28)$$

where $\sigma > 0$ is the bandwidth and $\{\mathbf{x}_\ell\}_{\ell=1}^b$ are the Gaussian centers randomly sampled from $\{\mathbf{x}_k^U\}_{k=1}^{n_U}$. The number of basis functions was set at $b = \min(n_U, 200)$. The candidates of the Gaussian bandwidth were

$$\text{median}(\{\|\mathbf{x}_i - \mathbf{x}_j\|\}_{i,j=1}^n) \times \left\{\frac{1}{8}, \frac{1}{4}, \frac{1}{2}, 1, 2\right\} \quad (4.29)$$

and that of the regularization parameter were $\{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$. The bandwidth and regularization parameter were determined by five-fold cross-validation. We varied the number of positive (unlabeled) samples from 10 to 200, with the number of unlabeled (positive) samples fixed. The class-prior was assumed to be known in this illustrative experiment and set at $\theta_P = 0.5$.

Figure 4.1 summarizes the average and standard error of the squared error of PU-SMI over 50 trials². The results show that the mean squared error decreased both when the number of positive samples was increased and the number of unlabeled samples was increased. Thus, both positive and unlabeled data contributed to improving the estimation accuracy of SMI, which well agrees with the theoretical result in Section 4.5.

²True SMI was computed by a supervised SMI estimator (Suzuki et al., 2009) with a sufficiently large number of positive and negative samples.

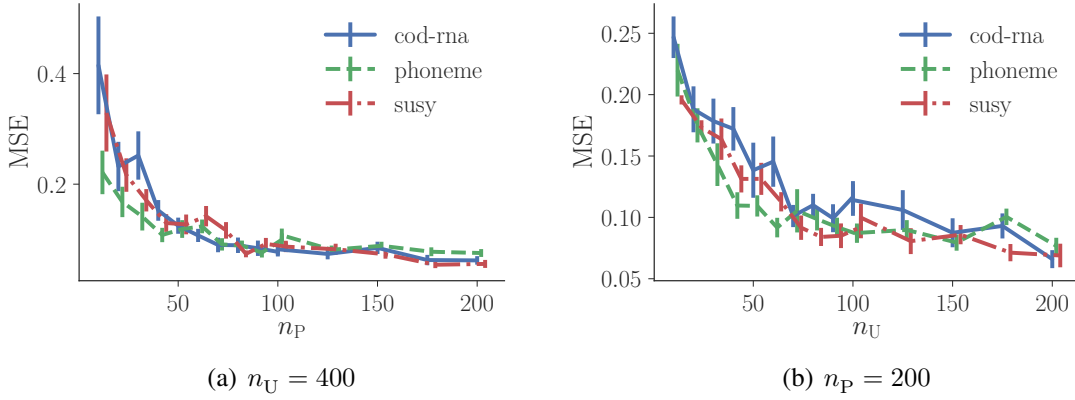


FIGURE 4.1: Average and standard error of the squared error of PU-SMI over 50 trials. (a) n_P is increased while $n_U = 400$ is fixed. (b) n_U is increased while $n_P = 200$ is fixed. The results show that both positive and unlabeled samples contribute to improving the estimation accuracy of SMI.

4.7.2 Dimension Reduction

Next, we evaluate the performance of the dimension reduction method based on PU-SMI (PUDR).

Illustration: We first illustrate how the PUDR method works on an artificial dataset. As a reference, we also show the result of the existing SMI-based dimension reduction method (LSDR) (Suzuki and Sugiyama, 2013) that uses both positive and negative data. For this illustrative experiment, we used a linear transformation as a projection matrix (i.e., linear dimension reduction).

For the artificial data, we generated samples from the following densities:

$$p(\mathbf{x} \mid y = +1) = N\left(\mathbf{x}; \begin{pmatrix} -3 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}\right), \quad (4.30)$$

$$p(\mathbf{x} \mid y = -1) = N\left(\mathbf{x}; \begin{pmatrix} 3 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}\right), \quad (4.31)$$

where $N(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the normal density with the mean vector $\boldsymbol{\mu}$ and the covariance matrix $\boldsymbol{\Sigma}$. From the densities, we drew $n_P = 200$ positive and $n_U = 400$ unlabeled samples for PUDR, and $n_P = 200$ positive and $n_N = 200$ negative samples for LSDR. Other experimental settings were the same as those in Section 4.7.1.

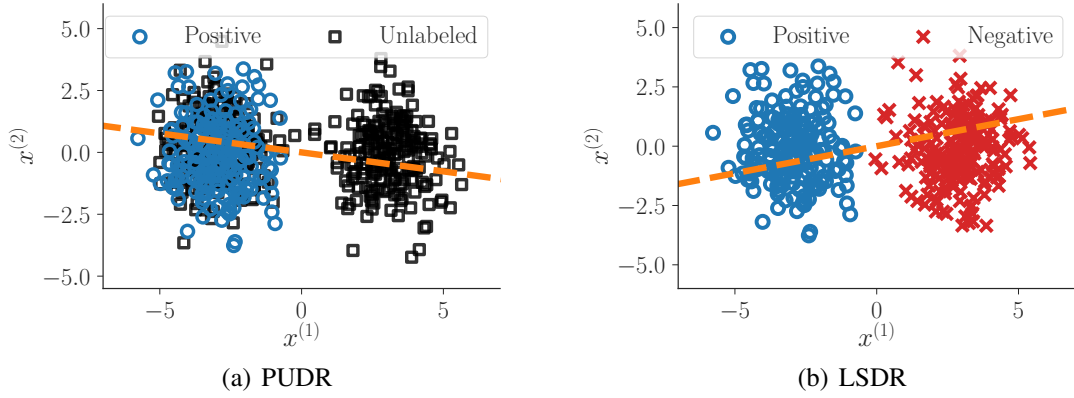


FIGURE 4.2: Illustration of linear dimension reduction. (a) Estimated subspace obtained by PUDR, where positive “ \circ ” and unlabeled samples “ \square ” are used. (b) Estimated subspace obtained by LSDR, where positive “ \circ ” and negative samples “ \times ” are used. The results show that from only PU data, the PUDR method gave a subspace similar to that obtained by its supervised counterpart.

We plot the estimated subspaces obtained by the PUDR and LSDR methods in Figure 4.2, showing that from only PU samples, the PUDR method obtained a subspace similar to that obtained by LSDR using positive and negative data.

Benchmark Data: Here, we apply the PUDR method to benchmark datasets. To obtain low-dimensional representation, we used a fully-connected neural network with four layers ($d=60-20-1$), the sigmoid function was used as the activation function, and *batch normalization* (Ioffe and Szegedy, 2015) was applied to all hidden layers. Stochastic gradient descent was used for optimization with learning rate 0.01. Also, weight decay with 0.001 and gradient noise with 0.01 were applied. We iteratively updated g (the last layer) with four minibatches and v (layers before the last layer) with one minibatch.

We compared the accuracy of class-prior estimation with and without dimension reduction. For comparison, we also considered principal component analysis (PCA) with different numbers of components: $\lfloor d/4 \rfloor$, $\lfloor d/2 \rfloor$, and $\lfloor 3d/4 \rfloor$, where $\lfloor \cdot \rfloor$ is the floor function. As a class-prior estimation method, we used the method based on the *kernel mean embedding* (KMM) method proposed by Ramaswamy et al. (2016). We used the mushrooms, the a9a, the MNIST, the CIFAR10, and the SVHN datasets. For multi-class classification datasets, we divided the whole classes into 2 groups to make binary classification

TABLE 4.1: Mean absolute error (with standard error) between the estimated class-prior and the true value over 20 trials. KM means the class-prior estimation method based on kernel mean embedding, PCA is the principal component analysis. The boldface denotes the best and comparable approach in terms of the average absolute error according to the t-test at the significance level 5%.

Dataset	θ_P	KM	PCA+KM			PUDR+KM
			$\lfloor d/4 \rfloor$	$\lfloor d/2 \rfloor$	$\lfloor 3d/4 \rfloor$	
mushrooms	0.3	0.05 (0.00)	0.03 (0.00)	0.05 (0.00)	0.05 (0.00)	0.03 (0.00)
	0.5	0.04 (0.01)	0.04 (0.00)	0.04 (0.01)	0.04 (0.01)	0.04 (0.01)
	0.7	0.03 (0.01)	0.03 (0.00)	0.03 (0.01)	0.03 (0.01)	0.05 (0.01)
a9a	0.3	0.11 (0.01)	0.11 (0.01)	0.11 (0.01)	0.11 (0.01)	0.04 (0.01)
	0.5	0.10 (0.01)	0.10 (0.01)	0.10 (0.01)	0.10 (0.01)	0.04 (0.01)
	0.7	0.08 (0.01)	0.08 (0.01)	0.08 (0.01)	0.08 (0.01)	0.03 (0.01)
MNIST	0.3	0.10 (0.00)	0.10 (0.00)	0.10 (0.00)	0.10 (0.00)	0.06 (0.01)
	0.5	0.12 (0.01)	0.12 (0.01)	0.12 (0.01)	0.12 (0.01)	0.07 (0.01)
	0.7	0.53 (0.07)	0.53 (0.07)	0.53 (0.07)	0.53 (0.07)	0.10 (0.01)
CIFAR10	0.3	0.30 (0.02)	0.57 (0.00)	0.56 (0.01)	0.50 (0.02)	0.17 (0.01)
	0.5	0.48 (0.01)	0.37 (0.00)	0.34 (0.01)	0.35 (0.01)	0.06 (0.01)
	0.7	0.65 (0.03)	0.17 (0.00)	0.15 (0.01)	0.16 (0.01)	0.12 (0.01)
SVHN	0.3	0.08 (0.03)	0.57 (0.00)	0.52 (0.03)	0.53 (0.01)	0.32 (0.02)
	0.5	0.04 (0.00)	0.37 (0.00)	0.29 (0.02)	0.34 (0.01)	0.19 (0.02)
	0.7	0.36 (0.08)	0.17 (0.00)	0.13 (0.02)	0.18 (0.02)	0.06 (0.01)

tasks. From the datasets, we drew $n_P = 1000$ positive and $n_U = 2000$ unlabeled samples. For model selection, we used validation samples of size $n_P = 50$ and $n_U = 200$.

Table 4.1 lists the mean absolute error (with standard error) between an estimated class-prior and the true value. Overall, the PUDR method tends to outperform other methods, meaning that the PUDR method provided useful low-dimensional representation except the SVHN dataset when $\theta_P = 0.3$ and 0.5. For the mushrooms, a9a, and MNIST datasets, applying the unsupervised dimension reduction method, PCA, did not improve the estimation accuracy. In fact, PCA performed poorly on the CIFAR10 and SVHN datasets, in particular, when $\theta_P = 0.3$.

4.7.3 Independence Test

Finally, we evaluate the performance of the SMI-based independence test (PUIT).

We computed the frequency of the *type-II* error, i.e., the number of times the null

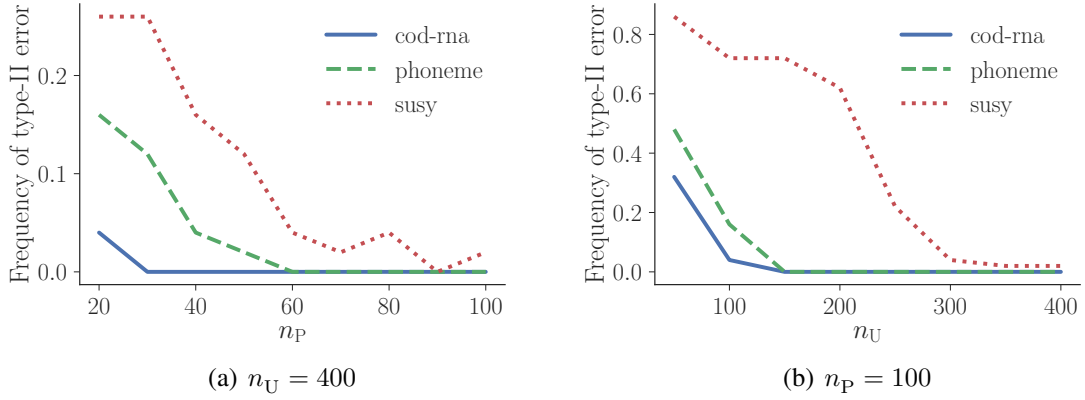


FIGURE 4.3: Frequency of the type-II error (the number of times the null hypothesis is accepted when the alternative hypothesis is true) of the PU-SMI based independence test with the significance level 5% over 50 trials. (a) n_P is increased while $n_U = 400$ is fixed. (b) n_U is increased while $n_P = 100$ is fixed. In both cases, the frequency of the type-II error decreased as the number of positive/unlabeled samples increases.

hypothesis (two random variables are independent) is accepted when the alternative hypothesis is true. Similarly to Section 4.7.1, we used classification datasets. Thus, the two random variables (pattern \mathbf{x} and class label y) are highly dependent, i.e., the frequency of the type-II error is expected to be low. Other experimental settings were the same as those in Section 4.7.1.

Figure 4.3 summarizes the frequency of the type-II error by the PUIT method at the significance level 5% over 50 trials. In both cases, the frequency decreased as the number of positive and unlabeled samples increased, showing that the PUIT method could detect statistical dependency well only from PU data.

4.8 Proof of Theorem 9

The idea of the proof is to regard the approximated squared error as perturbed optimization of expected one. Recall the linear-in-parameter model $g(\mathbf{x}) = \sum_{\ell=1}^b w_{\ell} \phi_{\ell}(\mathbf{x}) = \mathbf{w}^{\top} \boldsymbol{\phi}(\mathbf{x})$, where $0 \leq \phi_{\ell} \leq 1$ for all $\ell = 1, \dots, b$ and $\mathbf{x} \in \mathbb{R}^d$, and there exists a constant such that $\|\mathbf{w}\|_2 \leq M$. Moreover, we assume that the basis functions $\{\phi_{\ell}(\mathbf{x})\}_{\ell=1}^b$ are linearly independent over the marginal density $p(\mathbf{x}) > 0$.

Let

$$R_{\text{PU}}^{\text{SMI}}(\mathbf{w}) := \frac{1}{2} \mathbf{w}^\top \mathbf{H}_U \mathbf{w}^\top - \mathbf{w}^\top \mathbf{h}_P, \quad (4.32)$$

$$\widehat{R}_{\text{PU}}^{\text{SMI}, \lambda}(\mathbf{w}) := \frac{1}{2} \mathbf{w}^\top \widehat{\mathbf{H}}_U \mathbf{w}^\top - \mathbf{w}^\top \widehat{\mathbf{h}}_P + \frac{\lambda_{\text{PU}}}{2} \mathbf{w}^\top \mathbf{w} \quad (4.33)$$

be the expected squared error and the ℓ_2 -regularized empirical squared error, respectively.

To begin with, the following second order growth condition is proved:

Lemma 10. *Let ϵ be the smallest eigenvalue of \mathbf{H}_U . The second order growth condition holds*

$$R_{\text{PU}}^{\text{SMI}}(\mathbf{w}) \geq R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*) + \epsilon \|\mathbf{w} - \mathbf{w}^*\|_2^2. \quad (4.34)$$

Proof. Since \mathbf{H}_U is positive definite provided the linearly independent basis functions over $p(\mathbf{x})$, $R_{\text{PU}}^{\text{SMI}}(\mathbf{w})$ is strongly convex with parameter at least ϵ . Thus, we have

$$\begin{aligned} R_{\text{PU}}^{\text{SMI}}(\mathbf{w}) &\geq R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*) + \nabla R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*)^\top (\mathbf{w} - \mathbf{w}^*) + \epsilon \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\geq R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*) + \epsilon \|\mathbf{w} - \mathbf{w}^*\|_2^2, \end{aligned} \quad (4.35)$$

where the optimality condition $\nabla R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*) = \mathbf{0}$ is used. \square

Next, let us define a set of perturbation parameters as

$$\mathbf{u} := \{\mathbf{U}_U, \mathbf{u}_P \mid \mathbf{U}_U \in \mathbb{S}_+^b, \mathbf{u}_P \in \mathbb{R}^b\}, \quad (4.36)$$

where $\mathbb{S}_+^b \subset \mathbb{R}^{b \times b}$ is the cone of $b \times b$ positive semi-definite matrices. Our perturbed objective function and the solution are given by

$$R_{\text{PU}}^{\text{SMI}}(\mathbf{w}, \mathbf{u}) := \frac{1}{2} \mathbf{w}^\top (\mathbf{H}_U + \mathbf{U}_U) \mathbf{w} - \mathbf{w}^\top (\mathbf{h}_P + \mathbf{u}_P), \quad (4.37)$$

$$\mathbf{w}(\mathbf{u}) := \underset{\mathbf{w}}{\operatorname{argmin}} R_{\text{PU}}^{\text{SMI}}(\mathbf{w}, \mathbf{u}), \quad (4.38)$$

where

$$\mathbf{H}_U := \int \phi(\mathbf{x}) \phi(\mathbf{x})^\top p(\mathbf{x}) d\mathbf{x}, \quad (4.39)$$

$$\mathbf{h}_P := \int \phi(\mathbf{x}) p(\mathbf{x} \mid y = +1) d\mathbf{x}. \quad (4.40)$$

It can be easily confirmed $R_{\text{PU}}^{\text{SMI}}(\mathbf{w}) = R_{\text{PU}}^{\text{SMI}}(\mathbf{w}, \mathbf{0})$. Then, the following Lemma holds:

Lemma 11. $R_{\text{PU}}^{\text{SMI}}(\cdot, \mathbf{u}) - R_{\text{PU}}^{\text{SMI}}(\cdot)$ is Lipschitz continuous modulus $\omega(\mathbf{u}) = \mathcal{O}(\|\mathbf{U}_U\|_{\text{Fro}} + \|\mathbf{u}_P\|_2)$ on a sufficiently small neighborhood of \mathbf{w}^* , where $\|\cdot\|_{\text{Fro}}$ is the Frobenius norm.

Proof. Firstly, we have

$$R_{\text{PU}}^{\text{SMI}}(\mathbf{w}, \mathbf{u}) - R_{\text{PU}}^{\text{SMI}}(\mathbf{w}) = \frac{1}{2} \mathbf{w}^\top \mathbf{U}_U \mathbf{w} - \mathbf{w}^\top \mathbf{u}_P. \quad (4.41)$$

The partial gradient is given by

$$\frac{\partial}{\partial \mathbf{w}} (R_{\text{PU}}^{\text{SMI}}(\mathbf{w}, \mathbf{u}) - R_{\text{PU}}^{\text{SMI}}(\mathbf{w})) = \mathbf{U}_U \mathbf{w} - \mathbf{u}_P. \quad (4.42)$$

Let us define the δ -ball of \mathbf{w}^* as $B_\delta(\mathbf{w}^*) := \{\mathbf{w} \mid \|\mathbf{w} - \mathbf{w}^*\|_2 \leq \delta\}$. For any $\mathbf{w} \in B_\delta(\mathbf{w}^*)$, we can immediately show

$$\|\mathbf{w}\|_2 \leq \|\mathbf{w} - \mathbf{w}^*\|_2 + \|\mathbf{w}^*\|_2 \leq \delta + M. \quad (4.43)$$

Therefore, we have

$$\left\| \frac{\partial}{\partial \mathbf{w}} (R_{\text{PU}}^{\text{SMI}}(\mathbf{w}, \mathbf{u}) - R_{\text{PU}}^{\text{SMI}}(\mathbf{w})) \right\|_2 \leq (\delta + M) \|\mathbf{U}_U\|_{\text{Fro}} + \|\mathbf{u}_P\|_2. \quad (4.44)$$

This means that $R_{\text{PU}}^{\text{SMI}}(\cdot, \mathbf{u}) - R_{\text{PU}}^{\text{SMI}}(\cdot)$ is Lipschitz continuous on $B_\delta(\mathbf{w}^*)$ with a Lipschitz constant of order $\mathcal{O}(\|\mathbf{U}_U\|_{\text{Fro}} + \|\mathbf{u}_P\|_2)$. \square

Finally, we prove Theorem 9.

Proof. According to the *central limit theorem*, and our assumption about λ_{PU} , we have

$$\|\mathbf{U}_U\|_{\text{Fro}} = \mathcal{O}_p(1/\sqrt{n_U}), \quad \|\mathbf{u}_P\|_2 = \mathcal{O}_p(1/\sqrt{n_P})$$

as $n_P, n_U \rightarrow \infty$. Thus, by using Lemma 10, Lemma 11, and Proposition 6.1 in [Bonnans and Shapiro \(1998\)](#), we have the first half of Theorem 9:

$$\begin{aligned} \|\hat{\mathbf{w}} - \mathbf{w}^*\|_2 &\leq \epsilon^{-1} \omega(\mathbf{u}) \\ &= \mathcal{O}(\|\mathbf{U}_U\|_{\text{Fro}} + \|\mathbf{u}_P\|_2) \\ &= \mathcal{O}_p(1/\sqrt{n_P} + 1/\sqrt{n_U}). \end{aligned} \quad (4.45)$$

Next, we prove the latter half of Theorem 9. For the squared errors, we have

$$|\widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\widehat{\mathbf{w}}) - R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*)| \leq |\widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\widehat{\mathbf{w}}) - \widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\mathbf{w}^*)| + |\widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\mathbf{w}^*) - R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*)|. \quad (4.46)$$

Here, we have

$$\widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\widehat{\mathbf{w}}) - \widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\mathbf{w}^*) = \frac{1}{2}(\widehat{\mathbf{w}} + \mathbf{w}^*)^\top \widehat{\mathbf{H}}_{\text{U}}(\widehat{\mathbf{w}} - \mathbf{w}^*) - (\widehat{\mathbf{w}} - \mathbf{w}^*)^\top \widehat{\mathbf{h}}_{\text{P}}, \quad (4.47)$$

$$\widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\mathbf{w}^*) - R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*) = \frac{1}{2}\mathbf{w}^{*\top} \mathbf{U}_{\text{U}} \mathbf{w}^* - \mathbf{u}_{\text{P}} \mathbf{w}^* + \frac{\lambda_{\text{PU}}}{2} \mathbf{w}^{*\top} \mathbf{w}^*. \quad (4.48)$$

Since $0 \leq \phi_\ell(\mathbf{x}) \leq 1$, $\|\mathbf{w}^*\|_2 \leq M$, and $\lambda_{\text{PU}} = \mathcal{O}_p(1/\sqrt{\min(n_{\text{P}}, n_{\text{U}})})$, it leads to

$$\begin{aligned} |\widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\widehat{\mathbf{w}}) - R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*)| &\leq |\widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\widehat{\mathbf{w}}) - \widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\mathbf{w}^*)| + |\widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\mathbf{w}^*) - R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*)| \\ &\leq \mathcal{O}_p(\|\widehat{\mathbf{w}} - \mathbf{w}^*\|_2) + \mathcal{O}_p(\|\mathbf{U}_{\text{U}}\|_{\text{Fro}} + \|\mathbf{u}_{\text{P}}\|_2 + \lambda_{\text{PU}}) \\ &= \mathcal{O}_p(1/\sqrt{n_{\text{P}}} + 1/\sqrt{n_{\text{U}}}). \end{aligned} \quad (4.49)$$

As discussed in Section 4.3.2, we have

$$\text{PU-SMI} = -\frac{\theta_{\text{P}}}{\theta_{\text{N}}} \left(R_{\text{PU}}^{\text{SMI}}(\mathbf{w}^*) + \frac{1}{2} \right), \quad (4.50)$$

$$\widehat{\text{PU-SMI}} = -\frac{\theta_{\text{P}}}{\theta_{\text{N}}} \left(\widehat{R}_{\text{PU}}^{\text{SMI},\lambda}(\widehat{\mathbf{w}}) + \frac{1}{2} \right). \quad (4.51)$$

The final result indicates

$$|\text{PU-SMI} - \widehat{\text{PU-SMI}}| = \mathcal{O}_p(1/\sqrt{n_{\text{P}}} + 1/\sqrt{n_{\text{U}}}). \quad (4.52)$$

□

Chapter 5

Learning from Positive, Negative, and Unlabeled Data

In this chapter, we discuss a classification problem in a semi-supervised learning setup.

5.1 Introduction

In real-world machine learning applications, collecting a large amount of labeled data is a critical bottleneck owing to laborious manual annotation. On the other hand, unlabeled data can often be collected automatically and abundantly, e.g., by a web crawler. This fact has led to the development of various semi-supervised learning algorithms over the past decades.

To leverage unlabeled data in training, most of the existing semi-supervised learning methods rely on particular assumptions on the data distribution ([Chapelle et al., 2006](#)). For example, the *cluster assumption* requires that samples in the same cluster tend to have the same label. However, such strong distributional assumptions are rarely satisfied in practice. Once the data do not meet the required distributional assumption, the performance of the existing methods degrades, and sometimes it is even worse than that of supervised learning methods ([Cozman et al., 2003](#); [Sokolovska et al., 2008](#)).

In this chapter, we introduce a semi-supervised learning approach based on PU learning. Unlike most of the existing methods, our semi-supervised learning approach does not require the strong distributional assumptions. Furthermore, this approach can be applied any learning scenario as long as a risk function in PU learning is available. As the concrete implementation, we present two methods based on the misclassification rate and AUC. Additionally, without the distributional assumptions, we derive the generalization error bounds and show that the variance of the empirical risk estimator becomes smaller than its supervised counterpart. Experimental results report the usefulness of the methods.

5.2 Problem Setting

Let random variables $\mathbf{x} \in \mathbb{R}^d$ and $y \in \{+1, -1\}$ be equipped with probability density $p(\mathbf{x}, y)$, where d is a positive integer. Let us consider a binary classification problem from \mathbf{x} to y , given three sets of samples called the *positive* (P), *negative* (N), and *unlabeled* (U) data:

$$\{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = +1), \quad (5.1)$$

$$\{\mathbf{x}_i^N\}_{i=1}^{n_N} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = -1), \quad (5.2)$$

$$\{\mathbf{x}_i^U\}_{i=1}^{n_U} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}) = \theta_P p(\mathbf{x} \mid y = +1) + \theta_N p(\mathbf{x} \mid y = -1), \quad (5.3)$$

where

$$\theta_P := p(y = +1) \text{ and } \theta_N := p(y = -1)$$

are the class-prior probabilities for the positive and negative classes such that $\theta_P + \theta_N = 1$. Let E_P , E_N , and E_U be the expectations over $p(\mathbf{x} \mid y = +1)$, $p(\mathbf{x} \mid y = -1)$, and $p(\mathbf{x})$, respectively.

Furthermore, let $g: \mathbb{R}^d \rightarrow \mathbb{R}$ be an arbitrary real-valued decision function for binary classification, and classification is performed based on its sign. Let $\ell: \mathbb{R} \rightarrow \mathbb{R}$ be a loss function. The goal is to obtain an accurate classifier from P, N, and U samples.

Note that positive and negative (labeled) data are assumed to be drawn from corresponding conditional densities compared with standard supervised and semi-supervised learning settings in which labeled data are assumed to be drawn from the joint density $p(\mathbf{x}, y)$. In fact, our setting would be more practical than the existing one. Imagine that you will annotate samples collected from your domain of interest. Then, you may annotate some portion of the obtained data, e.g., 100 positive and 100 negative samples; then, the class-ratio of labeled data would be different from that of the underlying joint distribution. In this sense, our problem setting would be more practical.

5.3 Semi-Supervised Learning based on PU Learning

In this section, we introduce a semi-supervised learning approach based on PU learning.

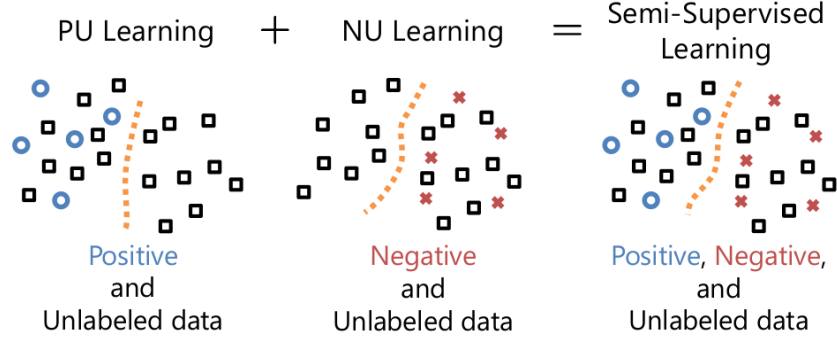


FIGURE 5.1: Illustration of PUNU Learning

5.3.1 PUNU Learning

The first idea to utilize positive, negative, and unlabeled data is to combine the PU and NU risks (see Figure 5.1). Let R_{PU} and R_{NU} be the risks in PU and NU learning, respectively.¹ For $\gamma \in [0, 1]$, let us consider a linear combination of the PU and NU risks (the PUNU risk):

$$R_{\text{PUNU}}^\gamma(g) := (1 - \gamma)R_{\text{PU}}(g) + \gamma R_{\text{NU}}(g). \quad (5.4)$$

We refer to this combined approach as *PUNU learning*.

Since this approach is independent of a classification measure, we can utilize any measures as long as the risk in PU (NU) learning is equivalent to that in PN learning. Below, we show two realizations of PUNU learning.

PUNU Classification: The PUNU risk in misclassification rate minimization (the PUNU-MR risk) is expressed as

$$R_{\text{PUNU}}^{\text{MR}, \gamma}(g) := (1 - \gamma)R_{\text{PU}}^{\text{MR}}(g) + \gamma R_{\text{NU}}^{\text{MR}}(g), \quad (5.5)$$

where $R_{\text{PU}}^{\text{MR}}$ and $R_{\text{NU}}^{\text{MR}}$ are the PU and NU risks in misclassification rate minimization. We refer to the method minimizing the PUNU-MR risk as *PUNU classification*.

¹ Depending on tasks, R_{PU} (R_{NU}) takes a specific form. For example, R_{PU} is replaced with $R_{\text{PU}}^{\text{MR}}$ in misclassification rate minimization, $R_{\text{PU}}^{\text{AUC}}$ in AUC optimization, $R_{\text{PU}}^{\text{SMI}}$ in an SMI estimation, respectively.

PUNU-AUC Optimization: The PUNU risk in AUC maximization (the PUNU-AUC risk) is expressed as

$$R_{\text{PUNU}}^{\text{AUC}, \gamma}(g) := (1 - \gamma)R_{\text{PU}}^{\text{AUC}}(g) + \gamma R_{\text{NU}}^{\text{AUC}}(g), \quad (5.6)$$

where $R_{\text{PU}}^{\text{AUC}}$ and $R_{\text{NU}}^{\text{AUC}}$ are the PU and NU risks in AUC maximization. We refer to the method minimizing the PUNU-AUC risk as *PUNU-AUC optimization*.

5.3.2 PNU Learning

Another possibility of using positive, negative, and unlabeled data is to combine the PN and PU/NU risks (see Figure 5.2). Let R_{PN} be the risk in PN learning.² For $\gamma \in [0, 1]$, let us consider linear combinations of the PN and PU/NU risks (the PNPU and PNNU risks):

$$R_{\text{PNPU}}^{\gamma}(g) := (1 - \gamma)R_{\text{PN}}(g) + \gamma R_{\text{PU}}(g), \quad (5.7)$$

$$R_{\text{PNNU}}^{\gamma}(g) := (1 - \gamma)R_{\text{PN}}(g) + \gamma R_{\text{NU}}(g). \quad (5.8)$$

We refer to the method minimizing the PNPU (resp., PNNU) risks as *PNPU* (resp., *PNNU*) *learning*. In practice, we combine the PNPU and PNNU risks and adaptively choose one of them with a new trade-off parameter $\eta \in [-1, 1]$ as

$$R_{\text{PNU}}^{\eta}(g) := \begin{cases} R_{\text{PNPU}}^{\eta}(g) & (\eta \geq 0), \\ R_{\text{PNNU}}^{-\eta}(g) & (\eta < 0). \end{cases} \quad (5.9)$$

We refer to this combined approach as *PNU learning*. Clearly, PNU learning with $\eta = -1, 0, +1$ corresponds to NU, PN, and PU learning. As η gets large (resp., small), the effect of the positive (resp., negative) class is more emphasized.

Similarly to PUNU learning, we show two realizations of PNU learning.

PNU Classification: The PNU risk in misclassification rate minimization (the PNU-MR risk) is expressed as

$$R_{\text{PNU}}^{\text{MR}, \eta}(g) := \begin{cases} R_{\text{PNPU}}^{\text{MR}, \eta}(g) & (\eta \geq 0), \\ R_{\text{PNNU}}^{\text{MR}, -\eta}(g) & (\eta < 0), \end{cases} \quad (5.10)$$

² Depending on tasks, R_{PN} takes a specific form. For example, R_{PN} is replaced with $R_{\text{PU}}^{\text{MR}}$ in misclassification rate minimization, $R_{\text{PU}}^{\text{AUC}}$ in AUC optimization, $R_{\text{PN}}^{\text{SMI}}$ in an SMI estimation, respectively.

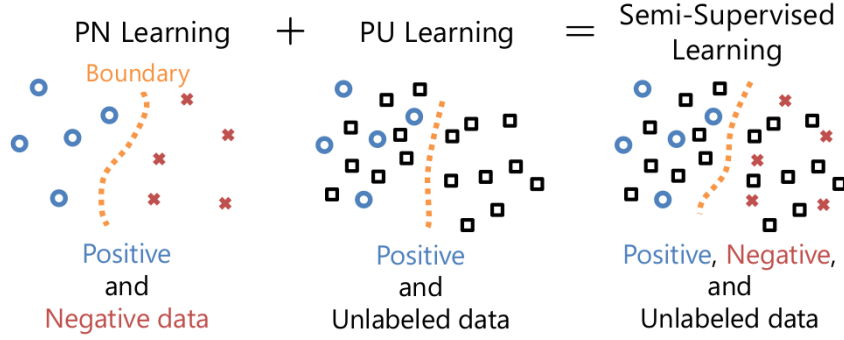


FIGURE 5.2: Illustration of PNU Learning

where

$$R_{\text{PNPU}}^{\text{MR},\gamma}(g) := (1 - \gamma)R_{\text{PN}}^{\text{MR}}(g) + \gamma R_{\text{PU}}^{\text{MR}}(g), \quad (5.11)$$

$$R_{\text{PNNU}}^{\text{MR},\gamma}(g) := (1 - \gamma)R_{\text{PN}}^{\text{MR}}(g) + \gamma R_{\text{NU}}^{\text{MR}}(g). \quad (5.12)$$

We refer to the method minimizing the PNU-MR risk as *PNU classification*.

PNU-AUC Optimization: The PNU risk in AUC maximization (the PNU-AUC risk) is expressed as

$$R_{\text{PNU}}^{\text{AUC},\eta}(g) := \begin{cases} R_{\text{PNPU}}^{\text{AUC},\eta}(g) & (\eta \geq 0), \\ R_{\text{PNNU}}^{\text{AUC},-\eta}(g) & (\eta < 0), \end{cases} \quad (5.13)$$

where

$$R_{\text{PNPU}}^{\text{AUC},\gamma}(g) := (1 - \gamma)R_{\text{PN}}^{\text{AUC}}(g) + \gamma R_{\text{PU}}^{\text{AUC}}(g), \quad (5.14)$$

$$R_{\text{PNNU}}^{\text{AUC},\gamma}(g) := (1 - \gamma)R_{\text{PN}}^{\text{AUC}}(g) + \gamma R_{\text{NU}}^{\text{AUC}}(g). \quad (5.15)$$

We refer to the method minimizing the PNU-AUC risk as *PNU-AUC optimization*.

In summary, we have the following lemma.

Lemma 12. *The following semi-supervised learning risks are equivalent:*

$$R^{\text{MR}}(g) = R_{\text{PUNU}}^{\text{MR},\gamma}(g) = R_{\text{PNPU}}^{\text{MR},\gamma}(g) = R_{\text{PNNU}}^{\text{MR},\gamma}(g) = R_{\text{PNU}}^{\text{MR},\eta}(g), \quad (5.16)$$

$$R^{\text{AUC}}(g) = R_{\text{PUNU}}^{\text{AUC},\gamma}(g) = R_{\text{PNPU}}^{\text{AUC},\gamma}(g) = R_{\text{PNNU}}^{\text{AUC},\gamma}(g) = R_{\text{PNU}}^{\text{AUC},\eta}(g). \quad (5.17)$$

It might seem that there is no reason for using different risk functions. However, the performance among risk estimators is different. In practice, the use of one risk function is better than others depending on, e.g., the size of samples in practice, similarly to PN, PU, and NU learning cases discussed in Section 2.5.2. In the next section, we discuss which approach, PUNU or PNU learning, is promising.

5.3.3 PUNU vs. PNU Learning

Here, we discuss which approach, PUNU or PNU learning, is more promising. Specifically, we focus on the misclassification rate as the classification measure and compare PUNU and PNU classification based on state-of-the-art theoretical comparisons (Niu et al., 2016), which are gave on estimation error bounds.

Let \hat{g}_{PN} , \hat{g}_{PU} , and \hat{g}_{NU} be the minimizers of $\hat{R}_{\text{PN}}^{\text{MR}}(g)$, $\hat{R}_{\text{PU}}^{\text{MR}}(g)$, and $\hat{R}_{\text{NU}}^{\text{MR}}(g)$, respectively. Let

$$\rho_{\text{PU},\text{PN}} := \frac{\theta_{\text{P}}/\sqrt{n_{\text{P}}} + 1/\sqrt{n_{\text{U}}}}{\theta_{\text{N}}/\sqrt{n_{\text{N}}}}, \quad (5.18)$$

$$\rho_{\text{NU},\text{PN}} := \frac{\theta_{\text{N}}/\sqrt{n_{\text{N}}} + 1/\sqrt{n_{\text{U}}}}{\theta_{\text{P}}/\sqrt{n_{\text{P}}}}. \quad (5.19)$$

The finite-sample comparisons (Niu et al., 2016) state that if $\rho_{\text{PU},\text{PN}} > 1$ (resp., $\rho_{\text{NU},\text{PN}} > 1$), PN classification is more promising than PU (resp., NU) classification in the sense that $R^{\text{MR}}(\hat{g}_{\text{PN}}) < R^{\text{MR}}(\hat{g}_{\text{PU}})$ (resp., $R^{\text{MR}}(\hat{g}_{\text{PN}}) < R^{\text{MR}}(\hat{g}_{\text{NU}})$); otherwise PU (resp., NU) classification is more promising than PN classification.

Suppose that n_{U} is not sufficiently large against n_{P} and n_{N} . According to the finite-sample comparisons, PN classification is most promising, and either PU or NU classification is the second best, i.e.,

$$R^{\text{MR}}(\hat{g}_{\text{PN}}) < R^{\text{MR}}(\hat{g}_{\text{PU}}) < R^{\text{MR}}(\hat{g}_{\text{NU}}) \quad (5.20)$$

or

$$R^{\text{MR}}(\hat{g}_{\text{PN}}) < R^{\text{MR}}(\hat{g}_{\text{NU}}) < R^{\text{MR}}(\hat{g}_{\text{PU}}). \quad (5.21)$$

On the other hand, if n_U is sufficiently large ($n_U \rightarrow \infty$, which is faster than $n_P, n_N \rightarrow \infty$), we have the asymptotic comparisons:

$$\rho_{PU,PN}^* = \lim_{n_P, n_N, n_U \rightarrow \infty} \rho_{PU,PN} = \frac{\theta_P / \sqrt{n_P}}{\theta_N / \sqrt{n_N}}, \quad (5.22)$$

$$\rho_{NU,PN}^* = \lim_{n_P, n_N, n_U \rightarrow \infty} \rho_{NU,PN} = \frac{\theta_N / \sqrt{n_N}}{\theta_P / \sqrt{n_P}}. \quad (5.23)$$

From the above results, we have

$$\rho_{PU,PN}^* \cdot \rho_{NU,PN}^* = 1. \quad (5.24)$$

If $\rho_{PU,PN}^* < 1$, then $\rho_{NU,PN}^* > 1$, implying that PU (resp., PN) classification is more promising than PN (resp., NU) classification, i.e.,

$$R^{\text{MR}}(\hat{g}_{PU}) < R^{\text{MR}}(\hat{g}_{PN}) < R^{\text{MR}}(\hat{g}_{NU}). \quad (5.25)$$

Similarly, when $\rho_{PU,PN}^* > 1$ and $\rho_{NU,PN}^* < 1$,

$$R^{\text{MR}}(\hat{g}_{NU}) < R^{\text{MR}}(\hat{g}_{PN}) < R^{\text{MR}}(\hat{g}_{PU}). \quad (5.26)$$

In real-world applications, since we do not know whether the number of unlabeled samples is sufficiently large or not, a practical approach would be to combine the best methods in both the finite-sample and asymptotic cases (See Eqs. (5.20), (5.21), (5.25), and (5.26)). From this viewpoint, PNU classification is the combination of the best methods in both cases, but PUNU classification is not. Additionally, PUNU classification includes the worst one in its combination in both cases.

In summary, PNU classification would be more promising than PUNU classification, as demonstrated in the experiments in Section 5.6.2.

5.4 Practical Implementation

In this section, practical implementation of the semi-supervised learning methods is explained.

5.4.1 General Case

We have so far only considered the true risk R (with respect to the expectations over true data distributions). When a classifier is trained with samples in practice, we use the empirical risk \widehat{R} , where the expectations are replaced with corresponding sample averages.

More specifically, we use a linear-in-parameter model given by

$$g(\mathbf{x}) = \sum_{j=1}^b w_j \phi_j(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}),$$

where $^\top$ denotes the transpose, b is the number of basis functions, $\mathbf{w} = (w_1, \dots, w_b)^\top$ is a parameter vector, and $\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_b(\mathbf{x}))^\top$ is a basis function vector. The parameter vector \mathbf{w} is learned in order to minimize the ℓ_2 -regularized empirical risk:

$$\min_{\mathbf{w}} \left[\widehat{R}(g) + \lambda \mathbf{w}^\top \mathbf{w} \right], \quad (5.27)$$

where $\lambda \geq 0$ is the regularization parameter.

5.4.2 PNU Classification

As given by Eq. (5.10), the PNU classification method consists of the PNPU-MR risk and the PNNU-MR risk. For the squared loss $\ell_S(m) := (1 - m)^2$, the empirical PNPU-MR risk can be expressed as

$$\begin{aligned} \widehat{R}_{\text{PNPU}}^{\text{MR}, \gamma}(g) &= \frac{(1 - \gamma)\theta_P}{n_P} \sum_{i=1}^{n_P} \ell_S(\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_i^P)) + \frac{(1 - \gamma)\theta_N}{n_N} \sum_{j=1}^{n_N} \ell_S(-\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_j^N)) \\ &\quad + \frac{\gamma\theta_P}{n_P} \sum_{i=1}^{n_P} \widetilde{\ell}_S(\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_i^P)) + \frac{\gamma}{n_U} \ell_S(-\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_k^U)) \\ &= (1 - \gamma)\mathbf{w}^\top \widehat{\mathbf{H}}_{\text{PN}}^{\text{MR}} \mathbf{w} - 2(1 - \gamma)\mathbf{w}^\top \widehat{\mathbf{h}}_{\text{PN}}^{\text{MR}} \\ &\quad + \gamma\mathbf{w}^\top \widehat{\mathbf{H}}_U^{\text{MR}} \mathbf{w} - 2\gamma\mathbf{w}^\top \widehat{\mathbf{h}}_{\text{PU}}^{\text{MR}} + 1, \end{aligned} \quad (5.28)$$

where

$$\hat{\mathbf{h}}_{\text{PN}}^{\text{MR}} := \frac{\theta_{\text{P}}}{n_{\text{P}}} \Phi_{\text{P}}^{\top} \mathbf{1}_{n_{\text{P}}} - \frac{\theta_{\text{N}}}{n_{\text{N}}} \Phi_{\text{N}}^{\top} \mathbf{1}_{n_{\text{N}}}, \quad (5.29)$$

$$\widehat{\mathbf{H}}_{\text{PN}}^{\text{MR}} := \frac{\theta_{\text{P}}}{n_{\text{P}}} \Phi_{\text{P}}^{\top} \Phi_{\text{P}} + \frac{\theta_{\text{N}}}{n_{\text{N}}} \Phi_{\text{N}}^{\top} \Phi_{\text{N}}, \quad (5.30)$$

$$\hat{\mathbf{h}}_{\text{PU}}^{\text{MR}} := \frac{2\theta_{\text{P}}}{n_{\text{P}}} \Phi_{\text{P}}^{\top} \mathbf{1}_{n_{\text{P}}} - \frac{1}{n_{\text{U}}} \Phi_{\text{U}}^{\top} \mathbf{1}_{n_{\text{U}}}, \quad (5.31)$$

$$\widehat{\mathbf{H}}_{\text{U}}^{\text{MR}} := \frac{1}{n_{\text{U}}} \Phi_{\text{U}}^{\top} \Phi_{\text{U}}, \quad (5.32)$$

$$\Phi_{\text{P}} := (\phi(\mathbf{x}_1^{\text{P}}), \dots, \phi(\mathbf{x}_{n_{\text{P}}}^{\text{P}}))^{\top}, \quad (5.33)$$

$$\Phi_{\text{N}} := (\phi(\mathbf{x}_1^{\text{N}}), \dots, \phi(\mathbf{x}_{n_{\text{N}}}^{\text{N}}))^{\top}, \quad (5.34)$$

$$\Phi_{\text{U}} := (\phi(\mathbf{x}_1^{\text{U}}), \dots, \phi(\mathbf{x}_{n_{\text{U}}}^{\text{U}}))^{\top}, \quad (5.35)$$

where $\mathbf{1}_b$ is a b -dimension vector whose elements are all one. The solution for the ℓ_2 -regularized PNPU classification can be analytically obtained by

$$\hat{\mathbf{w}}_{\text{PNPU}}^{\text{MR}, \gamma} := \left((1 - \gamma) \widehat{\mathbf{H}}_{\text{PN}}^{\text{MR}} + \gamma \widehat{\mathbf{H}}_{\text{U}}^{\text{MR}} + \lambda \mathbf{I}_b \right)^{-1} \left((1 - \gamma) \hat{\mathbf{h}}_{\text{PN}}^{\text{MR}} + \gamma \hat{\mathbf{h}}_{\text{PU}}^{\text{MR}} \right), \quad (5.36)$$

where $\lambda \geq 0$ is the regularization parameter. Similarly, the solution for the ℓ_2 -regularized PNNU classification can be obtained by

$$\hat{\mathbf{w}}_{\text{PNNU}}^{\text{MR}, \gamma} := \left((1 - \gamma) \widehat{\mathbf{H}}_{\text{PN}}^{\text{MR}} + \gamma \widehat{\mathbf{H}}_{\text{U}}^{\text{MR}} + \lambda \mathbf{I}_b \right)^{-1} \left((1 - \gamma) \hat{\mathbf{h}}_{\text{PN}}^{\text{MR}} + \gamma \hat{\mathbf{h}}_{\text{NU}}^{\text{MR}} \right), \quad (5.37)$$

where

$$\hat{\mathbf{h}}_{\text{NU}}^{\text{MR}} := \frac{1}{n_{\text{U}}} \Phi_{\text{U}}^{\top} \mathbf{1}_{n_{\text{U}}} - \frac{2\theta_{\text{N}}}{n_{\text{N}}} \Phi_{\text{N}}^{\top} \mathbf{1}_{n_{\text{N}}}. \quad (5.38)$$

5.4.3 PNU-AUC Optimization

As given by Eq. (5.13), the PNU-AUC optimization method consists of the PNPU-AUC and PNNU-AUC risks. For the squared loss $\ell_S(m) := (1 - m)^2$, the empirical PNPU-AUC risk can be expressed as

$$\begin{aligned}
\widehat{R}_{\text{PNPU}}^{\text{AUC}, \gamma}(g) &= \frac{1 - \gamma}{n_P n_N} \sum_{i=1}^{n_P} \sum_{j=1}^{n_N} \ell_S(\mathbf{w}^\top \bar{\phi}(\mathbf{x}_i^P, \mathbf{x}_j^N)) + \frac{\gamma}{\theta_N n_P n_U} \sum_{i=1}^{n_P} \sum_{k=1}^{n_U} \ell_S(\mathbf{w}^\top \bar{\phi}(\mathbf{x}_i^P, \mathbf{x}_k^U)) \\
&\quad - \frac{\gamma \theta_P}{\theta_N n_P (n_P - 1)} \sum_{i=1}^{n_P} \sum_{i'=1}^{n_P} \ell_S(\mathbf{w}^\top \bar{\phi}(\mathbf{x}_i^P, \mathbf{x}_{i'}^P)) + \frac{\gamma \theta_P}{\theta_N (n_P - 1)} \\
&= (1 - \gamma) - 2(1 - \gamma) \mathbf{w}^\top \widehat{\mathbf{h}}_{\text{PN}}^{\text{AUC}} + (1 - \gamma) \mathbf{w}^\top \widehat{\mathbf{H}}_{\text{PN}}^{\text{AUC}} \mathbf{w} \\
&\quad + \gamma - 2\gamma \mathbf{w}^\top \widehat{\mathbf{h}}_{\text{PU}}^{\text{AUC}} + \gamma \mathbf{w}^\top \widehat{\mathbf{H}}_{\text{PU}}^{\text{AUC}} \mathbf{w} - \gamma \mathbf{w}^\top \widehat{\mathbf{H}}_{\text{PP}}^{\text{AUC}} \mathbf{w}, \tag{5.39}
\end{aligned}$$

where

$$\widehat{\mathbf{h}}_{\text{PN}}^{\text{AUC}} := \frac{1}{n_P} \Phi_P^\top \mathbf{1}_{n_P} - \frac{1}{n_N} \Phi_N^\top \mathbf{1}_{n_N}, \tag{5.40}$$

$$\begin{aligned}
\widehat{\mathbf{H}}_{\text{PN}}^{\text{AUC}} &:= \frac{1}{n_P} \Phi_P^\top \Phi_P - \frac{1}{n_P n_N} \Phi_P^\top \mathbf{1}_{n_P} \mathbf{1}_{n_N}^\top \Phi_N \\
&\quad - \frac{1}{n_P n_N} \Phi_N^\top \mathbf{1}_{n_N} \mathbf{1}_{n_P}^\top \Phi_P + \frac{1}{n_N} \Phi_N^\top \Phi_N, \tag{5.41}
\end{aligned}$$

$$\widehat{\mathbf{h}}_{\text{PU}}^{\text{AUC}} := \frac{1}{\theta_N n_P} \Phi_P^\top \mathbf{1}_{n_P} - \frac{1}{\theta_N n_U} \Phi_U^\top \mathbf{1}_{n_U}, \tag{5.42}$$

$$\begin{aligned}
\widehat{\mathbf{H}}_{\text{PU}}^{\text{AUC}} &:= \frac{1}{\theta_N n_P} \Phi_P^\top \Phi_P - \frac{1}{\theta_N n_P n_U} \Phi_U^\top \mathbf{1}_{n_U} \mathbf{1}_{n_P}^\top \Phi_P \\
&\quad - \frac{1}{\theta_N n_P n_U} \Phi_P^\top \mathbf{1}_{n_P} \mathbf{1}_{n_U}^\top \Phi_U + \frac{1}{\theta_N n_U} \Phi_U^\top \Phi_U, \tag{5.43}
\end{aligned}$$

$$\widehat{\mathbf{H}}_{\text{PP}}^{\text{AUC}} := \frac{2\theta_P}{\theta_N (n_P - 1)} \Phi_P^\top \Phi_P - \frac{2\theta_P}{\theta_N n_P (n_P - 1)} \Phi_P^\top \mathbf{1}_{n_P} \mathbf{1}_{n_P}^\top \Phi_P, \tag{5.44}$$

$$\Phi_P := (\phi(\mathbf{x}_1^P), \dots, \phi(\mathbf{x}_{n_P}^P))^\top, \tag{5.45}$$

$$\Phi_U := (\phi(\mathbf{x}_1^U), \dots, \phi(\mathbf{x}_{n_U}^U))^\top, \tag{5.46}$$

$$\Phi_N := (\phi(\mathbf{x}_1^N), \dots, \phi(\mathbf{x}_{n_N}^N))^\top. \tag{5.47}$$

The solution for the ℓ_2 -regularized PNPU-AUC optimization can be analytically obtained by

$$\hat{\mathbf{w}}_{\text{PNPU}}^{\text{AUC},\gamma} := \left((1-\gamma)\widehat{\mathbf{H}}_{\text{PN}}^{\text{AUC}} + \gamma\widehat{\mathbf{H}}_{\text{PU}}^{\text{AUC}} - \gamma\widehat{\mathbf{H}}_{\text{PP}}^{\text{AUC}} + \lambda\mathbf{I}_b \right)^{-1} \left((1-\gamma)\hat{\mathbf{h}}_{\text{PN}}^{\text{AUC}} + \gamma\hat{\mathbf{h}}_{\text{PU}}^{\text{AUC}} \right). \quad (5.48)$$

Similarly, the solution for the ℓ_2 -regularized PNNU-AUC optimization can be obtained by

$$\hat{\mathbf{w}}_{\text{PNNU}}^{\text{AUC},\gamma} := \left((1-\gamma)\widehat{\mathbf{H}}_{\text{PN}}^{\text{AUC}} + \gamma\widehat{\mathbf{H}}_{\text{NU}}^{\text{AUC}} - \gamma\widehat{\mathbf{H}}_{\text{NN}}^{\text{AUC}} + \lambda\mathbf{I}_b \right)^{-1} \left((1-\gamma)\hat{\mathbf{h}}_{\text{PN}}^{\text{AUC}} + \gamma\hat{\mathbf{h}}_{\text{NU}}^{\text{AUC}} \right), \quad (5.49)$$

where

$$\hat{\mathbf{h}}_{\text{NU}}^{\text{AUC}} := \frac{1}{\theta_{\text{P}}n_{\text{U}}} \boldsymbol{\Phi}_{\text{U}}^{\top} \mathbf{1}_{n_{\text{U}}} - \frac{1}{\theta_{\text{P}}n_{\text{N}}} \boldsymbol{\Phi}_{\text{N}}^{\top} \mathbf{1}_{n_{\text{N}}}, \quad (5.50)$$

$$\begin{aligned} \widehat{\mathbf{H}}_{\text{NU}}^{\text{AUC}} &:= \frac{\theta_{\text{N}}}{\theta_{\text{P}}n_{\text{N}}} \boldsymbol{\Phi}_{\text{N}}^{\top} \boldsymbol{\Phi}_{\text{N}} - \frac{\theta_{\text{N}}}{\theta_{\text{P}}n_{\text{N}}n_{\text{U}}} \boldsymbol{\Phi}_{\text{U}}^{\top} \mathbf{1}_{n_{\text{U}}} \mathbf{1}_{n_{\text{N}}}^{\top} \boldsymbol{\Phi}_{\text{N}} \\ &\quad - \frac{\theta_{\text{N}}}{\theta_{\text{P}}n_{\text{N}}n_{\text{U}}} \boldsymbol{\Phi}_{\text{N}}^{\top} \mathbf{1}_{n_{\text{N}}} \mathbf{1}_{n_{\text{U}}}^{\top} \boldsymbol{\Phi}_{\text{U}} + \frac{\theta_{\text{N}}}{\theta_{\text{P}}n_{\text{U}}} \boldsymbol{\Phi}_{\text{U}}^{\top} \boldsymbol{\Phi}_{\text{U}}, \end{aligned} \quad (5.51)$$

$$\widehat{\mathbf{H}}_{\text{NN}}^{\text{AUC}} := \frac{2\theta_{\text{N}}}{\theta_{\text{P}}(n_{\text{N}}-1)} \boldsymbol{\Phi}_{\text{N}}^{\top} \boldsymbol{\Phi}_{\text{N}} - \frac{2\theta_{\text{N}}}{\theta_{\text{P}}n_{\text{N}}(n_{\text{N}}-1)} \boldsymbol{\Phi}_{\text{N}}^{\top} \mathbf{1}_{n_{\text{N}}} \mathbf{1}_{n_{\text{N}}}^{\top} \boldsymbol{\Phi}_{\text{N}}. \quad (5.52)$$

The computational complexity of computing $\hat{\mathbf{h}}_{\text{PN}}^{\text{AUC}}$ and $\widehat{\mathbf{H}}_{\text{PN}}^{\text{AUC}}$ are $\mathcal{O}((n_{\text{P}}+n_{\text{N}})b)$ and $\mathcal{O}((n_{\text{P}}+n_{\text{N}})b^2)$, respectively. Then, obtaining the solution $\hat{\mathbf{w}}_{\text{PNPU}}^{\text{AUC}}$ ($\hat{\mathbf{w}}_{\text{PNNU}}^{\text{AUC}}$) requires the computational complexity of $\mathcal{O}(b^3)$. Including the computational complexity of computing $\hat{\mathbf{h}}_{\text{PU}}^{\text{AUC}}$, $\widehat{\mathbf{H}}_{\text{PU}}^{\text{AUC}}$, and $\widehat{\mathbf{H}}_{\text{PP}}^{\text{AUC}}$, the total computational complexity of the PNPU-AUC optimization method is $\mathcal{O}((n_{\text{P}}+n_{\text{N}}+n_{\text{U}})b^2+b^3)$. Similarly, the total computational complexity of the PNNU-AUC optimization method is $\mathcal{O}((n_{\text{P}}+n_{\text{N}}+n_{\text{U}})b^2+b^3)$. In total, the computational complexity of the PNU-AUC optimization method is $\mathcal{O}((n_{\text{P}}+n_{\text{N}}+n_{\text{U}})b^2+b^3)$.

5.4.4 Model Selection

To tune the hyperparameters such as the regularization parameter λ , we use cross-validation. Unlike most of the semi-supervised learning methods, the semi-supervised learning approach based on PU learning can utilize unlabeled data for the risk estimation, meaning that unlabeled data can be utilized for computing a score of cross-validation.

PNU Classification: For the PNU classification method, we use the PNU risk in Eq. (5.10) with the zero-one loss as the score. To this end, however, we need to fix the combination parameter η in advance and then, we tune the hyperparameters including the combination parameter η . More specifically, let $\bar{\eta} \in [-1, 1]$ be the predefined combination parameter. We conduct cross-validation with respect to $R_{\text{PNU}}^{\text{MR}, \bar{\eta}}(g)$ for tuning the hyperparameters. Since the PNU risk is equivalent to the PN risk for any $\bar{\eta}$, we can choose any $\bar{\eta}$ in principle. However, when the empirical PNU risk is used in practice, choice of $\bar{\eta}$ may affect the performance of cross-validation.

Here, based on theoretical results of variance reduction given in Section 5.5.2, we give a practical method to determine $\bar{\eta}$. The variance reduction implies that the empirical PNU risk is more reliable than the PN risk because the variance of the empirical PNU risk becomes smaller than the PN risk, given a combination parameter in a certain interval. This means that the empirical PNU risk approximates the expected risk more accurately than the risk computed by only labeled data; therefore, hyperparameters will be chosen more reliably by the empirical PNU risk.

Assuming the standard deviations (in Section 5.5.2) satisfying $\sigma_P(g) = \sigma_N(g)$, we can obtain simpler expressions of Eqs. (5.74) and (5.75) as

$$\bar{\gamma}_{\text{PNPU}}^{\text{MR}} := \frac{1 - \theta_P^2 n_N / (\theta_N^2 n_P)}{1 + \theta_P^2 n_N / (\theta_N^2 n_P)}, \quad (5.53)$$

$$\bar{\gamma}_{\text{PNNU}}^{\text{MR}} := \frac{1 - \theta_N^2 n_P / (\theta_P^2 n_N)}{1 + \theta_N^2 n_P / (\theta_P^2 n_N)}. \quad (5.54)$$

They can be computed simply from the number of samples and the (estimated) class-prior. Finally, to select the combination parameter η for training a classifier, we use $\hat{R}_{\text{PNPU}}^{\bar{\gamma}_{\text{PNPU}}^{\text{MR}}}$ for $\eta \geq 0$, and $\hat{R}_{\text{PNNU}}^{\bar{\gamma}_{\text{PNNU}}^{\text{MR}}}$ for $\eta < 0$.

PNU-AUC Optimization: Similarly to the misclassification rate case, to conduct cross-validation, we need to fix the combination parameter in advance. Assuming the covariances (in Section 5.5.2), e.g., $\tau_{\text{PN}, \text{PP}}(g)$, are small enough to be neglected and $\tau_{\text{PN}}(g) = \tau_{\text{PP}}(g) = \tau_{\text{NN}}(g)$, we can obtain simpler expressions of Eqs. (5.85) and (5.86) as

$$\bar{\gamma}_{\text{PNPU}}^{\text{AUC}} := \frac{1}{1 + \theta_P^2 n_N / (\theta_N^2 n_P)}, \quad (5.55)$$

$$\bar{\gamma}_{\text{PNNU}}^{\text{AUC}} := \frac{1}{1 + \theta_N^2 n_P / (\theta_P^2 n_N)}. \quad (5.56)$$

Similarly to the misclassification rate case, they can be computed simply from the number of samples and the (estimated) class-prior. Finally, to select the combination parameter η , we use $\hat{R}_{\text{PNPU}}^{\text{AUC}}$ for $\eta \geq 0$, and $\hat{R}_{\text{PNNU}}^{\text{AUC}}$ for $\eta < 0$.

5.5 Theoretical Analyses

In this section, we theoretically analyze the behavior of the empirical risks of our semi-supervised learning methods. We first derive generalization error bounds and then discuss variance reduction.

5.5.1 Generalization Error Bounds

[Rigollet \(2007\)](#) proved generalization error bounds for semi-supervised learning by mathematically modeling the cluster assumption. As a consequence, the generalization error bounds hold when the cluster assumption is satisfied; otherwise, the generalization error bounds cannot be guaranteed. In contrast, we give generalization error bounds to our risk estimators without such a strong distributional assumption.

Similarly to Section 3.5, we assume that a classifier g is the linear-in-parameter model and it belongs to a function class of bounded hyperplanes \mathcal{G} :

$$\mathcal{G} := \{g(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) \mid \|\mathbf{w}\| \leq C_w; \forall \mathbf{x}: \|\phi(\mathbf{x})\| \leq C_\phi\}, \quad (5.57)$$

where $C_w > 0$ and $C_\phi > 0$ are certain positive constants.

Misclassification Rate Case: Denote by

$$I^{\text{MR}}(g) = \mathbb{E}_{p(\mathbf{x}, y)}[\ell_{0-1}(yg(\mathbf{x}))] \quad (5.58)$$

the risk of g for binary classification, i.e., the generalization error of g . In the following, we study upper bounds of $I^{\text{MR}}(g)$ holding uniformly for all $g \in \mathcal{G}$. We respectively focus on the (*scaled*) *ramp and squared losses* for the non-convex and convex methods. Similar results can be obtained with a little more effort if other eligible losses are used.

First, we consider the case that the scaled ramp loss function is used. Let $R_{\text{N-PUNU}}^{\text{MR}, \gamma}$, $R_{\text{N-PNPU}}^{\text{MR}, \gamma}$, and $R_{\text{N-PNNU}}^{\text{MR}, \gamma}$ be the PUNU-MR, PNPU-MR, and PNNU-MR risks equipped with the scaled ramp loss. A key observation is that $\ell_{0-1}(m) \leq 2\ell_{\text{R}}(m)$, and consequently

$I^{\text{MR}}(g) \leq 2R^{\text{MR}}(g)$. Note that by definition we have

$$R_{\text{N-PUNU}}^{\text{MR},\gamma}(g) = R_{\text{N-PNPU}}^{\text{MR},\gamma}(g) = R_{\text{N-PNNU}}^{\text{MR},\gamma}(g) = R^{\text{MR}}(g).$$

The theorem below can be proven using the Rademacher analysis (see, for example, [Mohri et al., 2012](#); [Ledoux and Talagrand, 1991](#)). The proof is available in Section 5.8.1.

Theorem 13. *Let $\ell_{\text{R}}(m)$ be the loss for defining the empirical risks. For any $\delta > 0$, the following inequalities hold separately with probability at least $1 - \delta$ for all $g \in \mathcal{G}$:*

$$I^{\text{MR}}(g) \leq 2\hat{R}_{\text{N-PUNU}}^{\text{MR},\gamma}(g) + c_1^{\text{MR}}(\delta) \cdot \left(\frac{2(1-\gamma)\theta_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{2\gamma\theta_{\text{N}}}{\sqrt{n_{\text{N}}}} + \frac{|2\gamma-1|}{\sqrt{n_{\text{U}}}} \right), \quad (5.59)$$

$$I^{\text{MR}}(g) \leq 2\hat{R}_{\text{N-PNPU}}^{\text{MR},\gamma}(g) + c_1^{\text{MR}}(\delta) \cdot \left(\frac{(1+\gamma)\theta_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{(1-\gamma)\theta_{\text{N}}}{\sqrt{n_{\text{N}}}} + \frac{\gamma}{\sqrt{n_{\text{U}}}} \right), \quad (5.60)$$

$$I^{\text{MR}}(g) \leq 2\hat{R}_{\text{N-PNNU}}^{\text{MR},\gamma}(g) + c_1^{\text{MR}}(\delta) \cdot \left(\frac{(1-\gamma)\theta_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{(1+\gamma)\theta_{\text{N}}}{\sqrt{n_{\text{N}}}} + \frac{\gamma}{\sqrt{n_{\text{U}}}} \right), \quad (5.61)$$

where $c_1^{\text{MR}}(\delta) = 2C_w C_\phi + \sqrt{2 \ln(3/\delta)}$.

Theorem 13 guarantees that when $\ell_{\text{R}}(m)$ is used, $I^{\text{MR}}(g)$ can be bounded from above by two times the empirical risks, i.e., $2\hat{R}_{\text{N-PUNU}}^{\text{MR},\gamma}(g)$, $2\hat{R}_{\text{N-PNPU}}^{\text{MR},\gamma}(g)$, and $2\hat{R}_{\text{N-PNNU}}^{\text{MR},\gamma}(g)$, plus the corresponding confidence terms of order

$$\mathcal{O}_p\left(\frac{1}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{N}}}} + \frac{1}{\sqrt{n_{\text{U}}}}\right).$$

Since n_{P} , n_{N} , and n_{U} can increase independently, this is already the optimal convergence rate without any additional assumption ([Vapnik, 1998](#); [Mendelson, 2008](#)).

Next, we prove the generalization error bounds when the squared loss function is used. Let $R_{\text{C-PUNU}}^{\text{MR},\gamma}$, $R_{\text{C-PNPU}}^{\text{MR},\gamma}$, and $R_{\text{C-PNNU}}^{\text{MR},\gamma}$ be the PUNU-MR, PNPU-MR, and PNNU-MR risks equipped with the squared loss. Analogously, we have $\ell_{0-1}(m) \leq 4\ell_{\text{S}}(m)$ for the squared loss. However, it is too loose when $|m| \gg 0$. Fortunately, we do not have to use $\ell_{\text{S}}(m)$ if we work on the generalization error rather than the estimation error. To this end, we define the *truncated (scaled) squared loss* $\ell_{\text{TS}}(m)$ as

$$\ell_{\text{TS}}(m) = \begin{cases} \ell_{\text{S}}(m) & 0 < m \leq 1, \\ \ell_{0-1}(m)/4 & \text{otherwise,} \end{cases} \quad (5.62)$$

so that $\ell_{0-1}(m) \leq 4\ell_{\text{TS}}(m)$ is much tighter. For $\ell_{\text{TS}}(m)$, $R_{\text{C-PU}}^{\text{MR}}(g)$ and $R_{\text{C-NU}}^{\text{MR}}(g)$ need to be redefined as follows (see Section 2.5 for details):

$$R_{\text{C-PU}}^{\text{MR}}(g) := \theta_{\text{P}} \mathbb{E}_{\text{P}}[\tilde{\ell}_{\text{TS}}(g(\mathbf{x}))] + \mathbb{E}_{\text{U}}[\ell_{\text{TS}}(-g(\mathbf{x}))], \quad (5.63)$$

$$R_{\text{C-NU}}^{\text{MR}}(g) := \theta_{\text{N}} \mathbb{E}_{\text{N}}[\tilde{\ell}_{\text{TS}}(-g(\mathbf{x}))] + \mathbb{E}_{\text{U}}[\ell_{\text{TS}}(g(\mathbf{x}))], \quad (5.64)$$

where $\tilde{\ell}_{\text{TS}}(m) = \ell_{\text{TS}}(m) - \ell_{\text{TS}}(-m)$ is the composite loss of truncated scaled squared loss. The condition $\tilde{\ell}_{\text{TS}}(m) \neq -m$ means the loss of convexity, but the equivalence is not lost; indeed, we still have

$$R_{\text{C-PUNU}}^{\text{MR},\gamma}(g) = R_{\text{C-PNPU}}^{\text{MR},\gamma}(g) = R_{\text{C-PNNU}}^{\text{MR},\gamma}(g) = R^{\text{MR}}(g).$$

Then, we have the following theorem (its proof is available in Section 5.8.2):

Theorem 14. *Let $\ell_{\text{TS}}(m)$ be the loss for defining the empirical risks (where $R_{\text{C-PU}}^{\text{MR}}(g)$ and $R_{\text{C-NU}}^{\text{MR}}(g)$ are redefined). For any $\delta > 0$, the following inequalities hold separately with probability at least $1 - \delta$ for all $g \in \mathcal{G}$:*

$$I^{\text{MR}}(g) \leq 4\hat{R}_{\text{C-PUNU}}^{\text{MR},\gamma}(g) + c_2^{\text{MR}}(\delta) \cdot \left(\frac{(1-\gamma)\theta_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{\gamma\theta_{\text{N}}}{\sqrt{n_{\text{N}}}} + \frac{1}{\sqrt{n_{\text{U}}}} \right), \quad (5.65)$$

$$I^{\text{MR}}(g) \leq 4\hat{R}_{\text{C-PNPU}}^{\text{MR},\gamma}(g) + c_2^{\text{MR}}(\delta) \cdot \left(\frac{\theta_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{(1-\gamma)\theta_{\text{N}}}{\sqrt{n_{\text{N}}}} + \frac{\gamma}{\sqrt{n_{\text{U}}}} \right), \quad (5.66)$$

$$I^{\text{MR}}(g) \leq 4\hat{R}_{\text{C-PNNU}}^{\text{MR},\gamma}(g) + c_2^{\text{MR}}(\delta) \cdot \left(\frac{(1-\gamma)\theta_{\text{P}}}{\sqrt{n_{\text{P}}}} + \frac{\theta_{\text{N}}}{\sqrt{n_{\text{N}}}} + \frac{\gamma}{\sqrt{n_{\text{U}}}} \right), \quad (5.67)$$

where $c_2^{\text{MR}}(\delta) = 4C_w C_\phi + \sqrt{2 \ln(4/\delta)}$.

Theorem 14 ensures that when $\ell_{\text{TS}}(m)$ is used (for evaluating the empirical risks rather than learning the empirical risk minimizers), $I^{\text{MR}}(g)$ can be bounded from above by four times the empirical risks plus confidence terms in the optimal parametric rate. As $\ell_{\text{TS}}(m) \leq \ell_{\text{S}}(m)$, Theorem 14 is valid (but weaker) if all empirical risks are w.r.t. $\ell_{\text{S}}(m)$.

AUC Case: We then prove the generalization error bounds of our risks for imbalanced binary classification. Similarly to Section 3.5, we assume that a surrogate loss is bounded from above by C_ℓ and satisfying $\ell_{0-1}(m) \leq \ell(m)$, and denote the Lipschitz constant by L .

Let

$$I^{\text{AUC}}(g) = \mathbb{E}_{\text{P}}[\mathbb{E}_{\text{N}}[\ell_{0-1}(g(\mathbf{x}^{\text{P}}) - g(\mathbf{x}^{\text{N}}))]] \quad (5.68)$$

be the generalization error of g in AUC optimization. For convenience, we define

$$c^{\text{AUC}}(\delta) := 2\sqrt{2}LC_\ell C_w C_\phi + \frac{3}{2}\sqrt{2\log(2/\delta)}. \quad (5.69)$$

In the following, we prove the generalization error bounds of the semi-supervised AUC optimization method.

For the PNPU-AUC and PNNU-AUC risks, we prove the following generalization error bounds (its proof is also available in Section 5.8.3):

Theorem 15. *For any $\delta > 0$, the following inequalities hold separately with probability at least $1 - \delta$ for all $g \in \mathcal{G}$:*

$$\begin{aligned} I^{\text{AUC}}(g) &\leq \widehat{R}_{\text{PNPU}}^{\text{AUC},\gamma}(g) + c^{\text{AUC}}(\delta/3) \\ &\quad \times \left(\frac{1-\gamma}{\sqrt{\min(n_P, n_N)}} + \frac{\gamma}{\theta_N \sqrt{\min(n_P, n_U)}} + \frac{\gamma\theta_P}{\theta_N \sqrt{n_P}} \right), \end{aligned} \quad (5.70)$$

$$\begin{aligned} I^{\text{AUC}}(g) &\leq \widehat{R}_{\text{PNNU}}^{\text{AUC},\gamma}(g) + c^{\text{AUC}}(\delta/3) \\ &\quad \times \left(\frac{1-\gamma}{\sqrt{\min(n_P, n_N)}} + \frac{\gamma}{\theta_P \sqrt{\min(n_N, n_U)}} + \frac{\gamma\theta_N}{\theta_P \sqrt{n_N}} \right), \end{aligned} \quad (5.71)$$

where $\widehat{R}_{\text{PNPU}}^{\text{AUC},\gamma}$ and $\widehat{R}_{\text{PNNU}}^{\text{AUC},\gamma}$ are unbiased empirical risk estimators corresponding to $R_{\text{PNPU}}^{\text{AUC},\gamma}$ and $R_{\text{PNNU}}^{\text{AUC},\gamma}$, respectively.

Theorem 15 guarantees that $I^{\text{AUC}}(g)$ can be bounded from above by the empirical risk, $\widehat{R}^{\text{AUC}}(g)$, plus the confidence terms of order

$$\mathcal{O}_p\left(\frac{1}{\sqrt{n_P}} + \frac{1}{\sqrt{n_N}} + \frac{1}{\sqrt{n_U}}\right).$$

Again, since n_P , n_N , and n_U can increase independently in our setting, this is the optimal convergence rate without any additional assumptions.

5.5.2 Variance Reduction

The risk estimators presented in Section 5.3 are all unbiased provided that the unbiased risk in PU learning is available with a specific classification measure. In this section, we explore the variance of the risk estimators. We show that unlabeled samples help reduce the variance of our risk estimators compared with its supervised counterpart $\widehat{R}_{\text{PN}}(g)$. For the sake of simplicity, we assume that $n_U \rightarrow \infty$, to illustrate the maximum variance reduction that could be achieved.

Misclassification Rate Case: For the sake of simplicity, we define

$$\sigma_P^2(g) := \text{Var}_P[\ell(g(\mathbf{x}))], \quad \sigma_N^2(g) := \text{Var}_N[\ell(-g(\mathbf{x}))],$$

where Var_P and Var_N denote the variance over $p(\mathbf{x} \mid y = +1)$ and $p(\mathbf{x} \mid y = -1)$. Moreover, denote by

$$\psi_P := \theta_P^2 \sigma_P^2(g)/n_P, \quad \psi_N := \theta_N^2 \sigma_N^2(g)/n_N$$

for short, and let Var be the variance over

$$\begin{aligned} & p(\mathbf{x}_1 \mid y = +1) \cdots p(\mathbf{x}_{n_P} \mid y = +1) \\ & \times p(\mathbf{x}_1 \mid y = -1) \cdots p(\mathbf{x}_{n_N} \mid y = -1) \\ & \times p(\mathbf{x}_1) \cdots p(\mathbf{x}_{n_U}). \end{aligned}$$

We have the following theorems (its proofs are available in Section 5.8.4):

Theorem 16. Assume $n_U \rightarrow \infty$. For any fixed g , let

$$\gamma_{N\text{-}P\text{UNU}}^{\text{MR}} = \underset{\gamma}{\text{argmin}} \text{Var}[\widehat{R}_{N\text{-}P\text{UNU}}^{\text{MR},\gamma}(g)] = \frac{\psi_P}{\psi_P + \psi_N}. \quad (5.72)$$

Then, we have $\gamma_{N\text{-}P\text{UNU}}^{\text{MR}} \in [0, 1]$. Further,

$$\text{Var}[\widehat{R}_{N\text{-}P\text{UNU}}^{\text{MR},\gamma}(g)] < \text{Var}[\widehat{R}_{PN}^{\text{MR}}(g)] \quad (5.73)$$

for all $\gamma \in (2\gamma_{N\text{-}P\text{UNU}}^{\text{MR}} - 1/2, 1/2)$ if $\psi_P < \psi_N$, or for all $\gamma \in (1/2, 2\gamma_{N\text{-}P\text{UNU}}^{\text{MR}} - 1/2)$ if $\psi_P > \psi_N$.³

Theorem 16 guarantees that the variance is always reduced by $\widehat{R}_{N\text{-}P\text{UNU}}^{\text{MR},\gamma}(g)$ if γ is close to $\gamma_{N\text{-}P\text{UNU}}^{\text{MR}}$, which is optimal for variance reduction. The interval of such good γ values has the length

$$\min \left(\frac{|\psi_P - \psi_N|}{\psi_P + \psi_N}, \frac{1}{2} \right).$$

In particular, if $3\psi_P \leq \psi_N$ or $\psi_P \geq 3\psi_N$, the length is $1/2$.

³Being fixed means g is determined before seeing the data for evaluating the empirical risk. For example, if g is trained by some learning method, and the empirical risk is subsequently evaluated on the independent validation/test data, g is regarded as fixed in the evaluation.

Theorem 17. Assume $n_U \rightarrow \infty$. For any fixed g , let

$$\gamma_{N-PNPU}^{MR} = \underset{\gamma}{\operatorname{argmin}} \operatorname{Var}[\hat{R}_{N-PNPU}^{MR,\gamma}(g)] = \frac{\psi_N - \psi_P}{\psi_P + \psi_N}, \quad (5.74)$$

$$\gamma_{N-PNNU}^{MR} = \underset{\gamma}{\operatorname{argmin}} \operatorname{Var}[\hat{R}_{N-PNNU}^{MR,\gamma}(g)] = \frac{\psi_P - \psi_N}{\psi_P + \psi_N}. \quad (5.75)$$

Then, we have $\gamma_{N-PNPU}^{MR} \in [0, 1]$ if $\psi_P \leq \psi_N$ or $\gamma_{N-PNNU}^{MR} \in [0, 1]$ if $\psi_P \geq \psi_N$. Additionally,

$$\operatorname{Var}[\hat{R}_{N-PNPU}^{MR,\gamma}(g)] < \operatorname{Var}[\hat{R}_{PN}^{MR}(g)] \quad (5.76)$$

for all $\gamma \in (0, 2\gamma_{N-PNPU}^{MR})$ if $\psi_P < \psi_N$, or

$$\operatorname{Var}[\hat{R}_{N-PNNU}^{MR,\gamma}(g)] < \operatorname{Var}[\hat{R}_{PN}^{MR}(g)] \quad (5.77)$$

for all $\gamma \in (0, 2\gamma_{N-PNNU}^{MR})$ if $\psi_P > \psi_N$.

Theorem 17 implies that the variance of $\hat{R}_{PN}^{MR}(g)$ is reduced by either $\hat{R}_{N-PNPU}^{MR,\gamma}(g)$ if $\psi_P \leq \psi_N$ or $\hat{R}_{N-PNNU}^{MR,\gamma}(g)$ if $\psi_P \geq \psi_N$, where γ should be close to γ_{N-PNPU}^{MR} or γ_{N-PNNU}^{MR} . The range of such good γ values is of length

$$\min\left(\frac{2|\psi_P - \psi_N|}{\psi_P + \psi_N}, 1\right).$$

In particular, if $3\psi_P \leq \psi_N$, $\hat{R}_{N-PNPU}^{MR,\gamma}(g)$ given any $\gamma \in (0, 1)$ can reduce the variance, and if $\psi_P \geq 3\psi_N$, $\hat{R}_{N-PNNU}^{MR,\gamma}(g)$ given any $\gamma \in (0, 1)$ can reduce the variance.

As a corollary of Theorems 16 and 17, the minimum variance achievable by $\hat{R}_{N-PUNU}^{MR,\gamma}(g)$, $\hat{R}_{N-PNPU}^{MR,\gamma}(g)$, and $\hat{R}_{N-PNNU}^{MR,\gamma}(g)$ at their optimal γ_{N-PUNU}^{MR} , γ_{N-PNPU}^{MR} , and γ_{N-PNNU}^{MR} is exactly the same, namely, $4\psi_P\psi_N/(\psi_P + \psi_N)$. Nevertheless, $\hat{R}_{N-PNPU}^{MR,\gamma}(g)$ and $\hat{R}_{N-PNNU}^{MR,\gamma}(g)$ have a much wider range of nice γ values than $\hat{R}_{N-PUNU}^{MR,\gamma}(g)$.

If we further assume that $\sigma_P(g) = \sigma_N(g)$, the condition in Theorems 16 and 17 as to whether $\psi_P \leq \psi_N$ or $\psi_P \geq \psi_N$ will be independent of g . Also, it will coincide with the condition in Theorem 7 in Niu et al. (2016) where the minimizers of $\hat{R}_{PN}^{MR}(g)$, $\hat{R}_{PU}^{MR}(g)$ and $\hat{R}_{NU}^{MR}(g)$ are compared.

A remark is that learning is uninvolved in Theorems 16 and 17, such that $\ell(m)$ can be any loss that satisfies $\ell(m) + \ell(-m) = 1$, and g can be any fixed decision function. For instance, we may adopt $\ell_{0-1}(m)$ and pick some g resulted from some other learning methods. As a consequence, the variance of $\hat{I}_{PN}^{MR}(g)$ over the validation data can be reduced,

and then the cross-validation should be more stable, given that n_U is sufficiently large. Therefore, even without being minimized, our proposed risk estimators are themselves of practical importance.

AUC Case: Similarly to the misclassification rate case, we here investigate if the PNU-AUC risk estimator has smaller variance than its supervised counterpart.

Let us introduce the following variances and covariances:

$$\tau_{PN}^2(g) = \text{Var}_{PN}[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^N))], \quad (5.78)$$

$$\tau_{PP}^2(g) = \text{Var}_{PP}[\ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P))], \quad (5.79)$$

$$\tau_{NN}^2(g) = \text{Var}_{NN}[\ell(g(\mathbf{x}^N) - g(\bar{\mathbf{x}}^N))], \quad (5.80)$$

$$\tau_{PN,PP}(g) = \text{Cov}_{PN,PP}[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^N)), \ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P))], \quad (5.81)$$

$$\tau_{PN,NN}(g) = \text{Cov}_{PN,NN}[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^N)), \ell(g(\mathbf{x}^N) - g(\bar{\mathbf{x}}^N))], \quad (5.82)$$

$$\tau_{PU,PP}(g) = \text{Cov}_{PU,PP}[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^U)), \ell(g(\mathbf{x}^P) - g(\bar{\mathbf{x}}^P))], \quad (5.83)$$

$$\tau_{NU,NN}(g) = \text{Cov}_{NU,NN}[\ell(g(\mathbf{x}^P) - g(\mathbf{x}^U)), \ell(g(\mathbf{x}^N) - g(\bar{\mathbf{x}}^N))], \quad (5.84)$$

where Var_{PN} , Var_{PP} , and Var_{NN} are the variances over

$$\begin{aligned} & p(\mathbf{x}^P | y = +1)p(\mathbf{x}^N | y = -1), \\ & p(\mathbf{x}^P | y = +1)p(\bar{\mathbf{x}}^P | y = +1), \\ & p(\mathbf{x}^N | y = -1)p(\bar{\mathbf{x}}^N | y = -1), \end{aligned}$$

respectively. $\text{Cov}_{PN,PP}$, $\text{Cov}_{PN,NN}$, $\text{Cov}_{PU,PP}$, and $\text{Cov}_{NU,NN}$ are the covariances over

$$\begin{aligned} & p(\mathbf{x}^P | y = +1)p(\mathbf{x}^N | y = -1)p(\bar{\mathbf{x}}^P | y = +1), \\ & p(\mathbf{x}^P | y = +1)p(\mathbf{x}^N | y = -1)p(\bar{\mathbf{x}}^N | y = -1), \\ & p(\mathbf{x}^P | y = +1)p(\mathbf{x}^U)p(\bar{\mathbf{x}}^P | y = +1), \\ & p(\mathbf{x}^N | y = -1)p(\mathbf{x}^U)p(\bar{\mathbf{x}}^N | y = -1), \end{aligned}$$

respectively.

Then, we have the following theorem (its proof is available in Section 5.8.5):

Theorem 18. Assume $n_U \rightarrow \infty$. For any fixed g , the minimizers of the variance of the empirical PNP-U-AUC and PNN-U-AUC risks are respectively obtained by

$$\gamma_{\text{PNPU}}^{\text{AUC}} = \underset{\gamma}{\operatorname{argmin}} \operatorname{Var}[\hat{R}_{\text{PNPU}}^{\text{AUC},\gamma}(g)] = \frac{\bar{\psi}_{\text{PN}} - \bar{\psi}_{\text{PP}}/2}{\bar{\psi}_{\text{PN}} + \bar{\psi}_{\text{PU}} - \bar{\psi}_{\text{PP}}}, \quad (5.85)$$

$$\gamma_{\text{PNNU}}^{\text{AUC}} = \underset{\gamma}{\operatorname{argmin}} \operatorname{Var}[\hat{R}_{\text{PNNU}}^{\text{AUC},\gamma}(g)] = \frac{\bar{\psi}_{\text{PN}} - \bar{\psi}_{\text{NN}}/2}{\bar{\psi}_{\text{PN}} + \bar{\psi}_{\text{NU}} - \bar{\psi}_{\text{NN}}}, \quad (5.86)$$

where

$$\bar{\psi}_{\text{PN}} = \frac{1}{n_{\text{P}}n_{\text{N}}} \tau_{\text{PN}}^2(g), \quad (5.87)$$

$$\bar{\psi}_{\text{PU}} = \frac{\theta_{\text{P}}^2}{\theta_{\text{N}}^2 n_{\text{P}}^2} \tau_{\text{PP}}^2(g) - \frac{\theta_{\text{P}}}{\theta_{\text{N}}^2 n_{\text{P}}} \tau_{\text{PU,PP}}(g), \quad (5.88)$$

$$\bar{\psi}_{\text{PP}} = \frac{1}{\theta_{\text{N}} n_{\text{P}}} \tau_{\text{PN,PU}}(g) - \frac{\theta_{\text{P}}}{\theta_{\text{N}} n_{\text{P}}} \tau_{\text{PN,PP}}(g), \quad (5.89)$$

$$\bar{\psi}_{\text{NU}} = \frac{\theta_{\text{N}}^2}{\theta_{\text{P}}^2 n_{\text{N}}^2} \tau_{\text{NN}}^2(g) - \frac{\theta_{\text{N}}}{\theta_{\text{P}}^2 n_{\text{N}}} \tau_{\text{NU,NN}}(g), \quad (5.90)$$

$$\bar{\psi}_{\text{NN}} = \frac{1}{\theta_{\text{P}} n_{\text{N}}} \tau_{\text{PN,NU}}(g) - \frac{\theta_{\text{N}}}{\theta_{\text{P}} n_{\text{N}}} \tau_{\text{PN,NN}}(g). \quad (5.91)$$

Additionally, we have

$$\operatorname{Var}[\hat{R}_{\text{PNPU}}^{\text{AUC},\gamma}(g)] < \operatorname{Var}[\hat{R}_{\text{PN}}^{\text{AUC}}(g)] \quad (5.92)$$

for any $\gamma \in (0, 2\gamma_{\text{PNPU}}^{\text{AUC}})$ if $\bar{\psi}_{\text{PN}} + \bar{\psi}_{\text{PU}} > \bar{\psi}_{\text{PP}}$ and $2\bar{\psi}_{\text{PN}} > \bar{\psi}_{\text{PP}}$. Similarly, we have

$$\operatorname{Var}[\hat{R}_{\text{PNNU}}^{\text{AUC},\gamma}(g)] < \operatorname{Var}[\hat{R}_{\text{PN}}^{\text{AUC}}(g)] \quad (5.93)$$

for any $\gamma \in (0, 2\gamma_{\text{PNNU}}^{\text{AUC}})$ if $\bar{\psi}_{\text{PN}} + \bar{\psi}_{\text{NU}} > \bar{\psi}_{\text{NN}}$ and $2\bar{\psi}_{\text{PN}} > \bar{\psi}_{\text{NN}}$.

This theorem means that, if γ is chosen appropriately, our proposed risk estimators, $\hat{R}_{\text{PNPU}}^{\text{AUC},\gamma}$ and $\hat{R}_{\text{PNNU}}^{\text{AUC},\gamma}$, have smaller variance than the standard supervised risk estimator $\hat{R}_{\text{PN}}^{\text{AUC}}$. A practical consequence of Theorem 18 is that when we conduct cross-validation for hyperparameter selection, we may use the PNU-AUC risk estimators $\hat{R}_{\text{PNPU}}^{\text{AUC},\gamma}$ and $\hat{R}_{\text{PNNU}}^{\text{AUC},\gamma}$ instead of the standard supervised risk estimator $\hat{R}_{\text{PN}}^{\text{AUC}}$ since they are more stable (see Section 5.4.4 for details).

5.6 Experiments on Balanced Classification

In this section, we first numerically analyze our semi-supervised learning approach with the misclassification rate and then compare the methods against existing methods. All experiments were carried out using a PC equipped with two 2.60GHz Intel® Xeon® E5-2640 v3 CPUs.

5.6.1 Experimental Analyses

Here, we numerically analyze the behavior of the PNU classification method.

Common Setup: As a classifier, we use the Gaussian kernel model:

$$g(\mathbf{x}) = \sum_{i=1}^n w_i \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2\sigma^2}\right),$$

where $n = n_P + n_N$, $\{w_i\}_{i=1}^n$ are the parameters, $\{\mathbf{x}_i\}_{i=1}^n = \mathcal{X}_P \cup \mathcal{X}_N$, and $\sigma > 0$ is the Gaussian bandwidth. The bandwidth candidates are

$$\text{median}(\|\mathbf{x}_i - \mathbf{x}_j\|_{i,j=1}^n) \times \left\{ \frac{1}{8}, \frac{1}{4}, \frac{1}{2}, 1, \frac{3}{2}, 2 \right\}.$$

We denote the classifier trained by minimizing the empirical PN risk by \hat{g}_{PN} . The number of labeled samples for training is 20, where the class-prior was 0.5. Note that the class-prior of test data was the same as that of unlabeled data.

Effect of Different Loss Functions: We first compare the performance of our semi-supervised classification methods with different loss functions. Specifically, we compare the PNU classification method with the ramp loss (RL), squared loss (SL), logistic loss (LL), and double hinge loss (DH) and the PUNU classification method with the squared loss and logistic loss.

Table 5.1 summarizes the average and standard error of the misclassification rates over 30 trials, showing that the clear difference was not observed between loss functions for both the PNU and PUNU classification methods. As discussed in Section 5.3.3, the PNU classification method tends to outperform the PUNU classification method. This tendency will be confirmed by the experiments on more datasets in Section 5.6.2.

Figure 5.3 plots the computation time, showing that the PNU and PUNU classification methods with the squared loss were the fastest. The methods with the logistic loss were

TABLE 5.1: Average and standard error of the misclassification rates of each method over 30 trials for benchmark datasets. The boldface denotes the best and comparable methods in terms of the average misclassification rate according to the t-test at the significance level 5%.

Dataset	θ_p	$\hat{\theta}_p$	$n_L = 50, n_U = 300$					
			PNU(RL)	PNU(DH)	PNU(SL)	PNU(LL)	PUNU(SL)	PUNU(LL)
Phoneme	0.5	0.58 (0.02)	28.4 (1.0)	27.7 (0.7)	27.9 (0.9)	28.2 (0.8)	28.7 (0.9)	28.1 (0.8)
cod-rna	0.5	0.63 (0.02)	16.4 (2.1)	16.9 (1.9)	16.5 (2.0)	17.4 (1.7)	20.6 (1.7)	23.3 (1.8)
Waveform	0.5	0.52 (0.01)	14.9 (0.6)	13.7 (0.5)	14.3 (0.5)	13.3 (0.4)	15.8 (0.6)	13.6 (0.5)
Spambase	0.5	0.62 (0.02)	24.6 (1.3)	23.1 (1.1)	24.7 (1.4)	24.5 (1.4)	25.3 (1.1)	24.4 (1.0)
phishing	0.5	0.58 (0.01)	13.3 (0.6)	13.7 (0.7)	13.3 (0.7)	14.9 (0.8)	15.8 (0.7)	17.4 (0.7)
Coil2	0.5	0.62 (0.02)	27.3 (1.4)	26.7 (1.4)	27.4 (1.3)	27.6 (1.4)	27.4 (1.5)	28.4 (1.4)

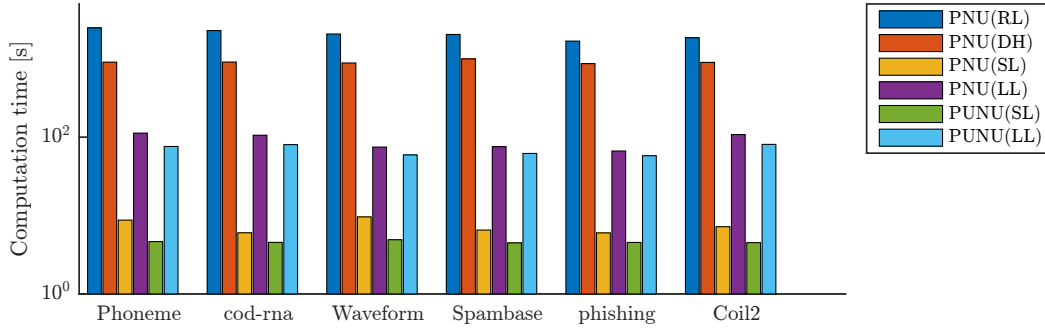


FIGURE 5.3: The average computation time over 30 trials for benchmark datasets.

slower than those with the squared loss. The ramp loss and double-hinge loss methods were much slower than other methods.

In summary, from the viewpoint of performance and computational efficiency, the squared loss was shown to be advantageous and thus we use it in the subsequent experiments.

Variance Reduction in Practice: Next, we numerically investigate how many unlabeled samples are sufficient in practice so that the variance of the empirical PNU risk is smaller than that of the PN risk: $\text{Var}[\hat{R}_{\text{PNU}}^{\text{MR},\eta}(g)] < \text{Var}[\hat{R}_{\text{PN}}^{\text{MR}}(g)]$ given a fixed classifier g .

As the fixed classifier, we used the classifier \hat{g}_{PN} , where the hyperparameters were determined by five-fold cross-validation. To compute the variance of the empirical PN and PNU risks, $\text{Var}[\hat{R}_{\text{PN}}^{\text{MR}}(\hat{g}_{\text{PN}})]$ and $\text{Var}[\hat{R}_{\text{PNU}}^{\text{MR},\eta}(\hat{g}_{\text{PN}})]$, we repeatedly drew additional $n_p^V = 10$ positive, $n_N^V = 10$ negative, and n_U^V unlabeled samples from the rest of the dataset. The

additional samples were also used for estimating $\hat{\sigma}_P(\hat{g}_{PN})$ and $\hat{\sigma}_N(\hat{g}_{PN})$ to compute η , i.e., γ as in Section 5.4.4.

Figure 5.4 shows the ratio between the variance of the empirical PNU risk and that of the PN risk, $\text{Var}[\hat{R}_{\text{PNU}}^{\text{MR},\eta}(\hat{g}_{PN})] / \text{Var}[\hat{R}_{\text{PN}}^{\text{MR}}(\hat{g}_{PN})]$, where the number of unlabeled samples for validation n_U^V increases from 10 to 300. With a rather small number of unlabeled samples, we see that the ratio becomes less than 1, meaning that the variance of the empirical PNU risk becomes smaller than that of the PN risk. Although the variance reduction is proved given an infinite number of unlabeled samples, in practice, a finite number of samples is sufficient for reducing the variance.

Compared to cases when $\theta_P = 0.3$ and 0.7 , the effect of variance reduction is small when $\theta_P = 0.5$ because $\gamma_{\text{N-PNPU}}^{\text{MR}} \approx \gamma_{\text{N-PNNU}}^{\text{MR}} \approx 0$ if we assume $\sigma_P(g) \approx \sigma_N(g)$ when $n_P \approx n_N$ and $\theta_P = 0.5$ (i.e., $\psi_P \approx \psi_N$. See Theorem 17). That is, the PNU risk is dominated by the PN risk, implying that $\text{Var}[\hat{R}_{\text{PNU}}^{\text{MR},\eta}(g)] \approx \text{Var}[\hat{R}_{\text{PN}}^{\text{MR}}(g)]$. It is worth noting that the class-prior is not the only factor for variance reduction; for example, the variance reduction will be large if $\theta_P = 0.5$, $n_P \gg n_N$, and $\sigma_P(g) \approx \sigma_N(g)$ due to $\gamma_{\text{N-PNPU}}^{\text{MR}} \not\approx 0$ (because $\psi_P \ll \psi_N$).

PNU Risk in Validation: As discussed in Section 5.5.2, the empirical PNU risk will be a reliable validation score since its variance is smaller than the empirical PN risk. Here, we show that the empirical PNU risk is a promising alternative to a validation score.

To focus on the effect of the risk function for validation only, we fix a training method. More specifically, we prepared two classifiers trained by the same risk, e.g., the empirical PN risk. We then tune the classifiers with the empirical PN and PNU risks denoted by $\hat{g}_{\text{PN}}^{\text{PN}}$ and $\hat{g}_{\text{PN}}^{\text{PNU}}$, respectively. For validation, we draw additional $n_P^V = 10$ positive, $n_N^V = 10$ negative, and n_U^V unlabeled samples from the rest of the dataset.

Figure 5.5 shows the ratio between the misclassification rate of $\hat{g}_{\text{PN}}^{\text{PNU}}$ and that of $\hat{g}_{\text{PN}}^{\text{PN}}$: $R^{\text{MR}}(\hat{g}_{\text{PN}}^{\text{PNU}}) / R^{\text{MR}}(\hat{g}_{\text{PN}}^{\text{PN}})$. With a rather small number of unlabeled samples, the ratio becomes less than 1, i.e., $\hat{g}_{\text{PN}}^{\text{PNU}}$ achieves better performance than $\hat{g}_{\text{PN}}^{\text{PN}}$. In particular, when $\theta_P = 0.3$ and 0.7 , $\hat{g}_{\text{PN}}^{\text{PNU}}$ improved substantially; the large variance reduction tends to give the large improvement (cf. Figure 5.4). This result shows that the use of the empirical PNU risk for validation improved the classification performance given a rather small number of unlabeled samples.

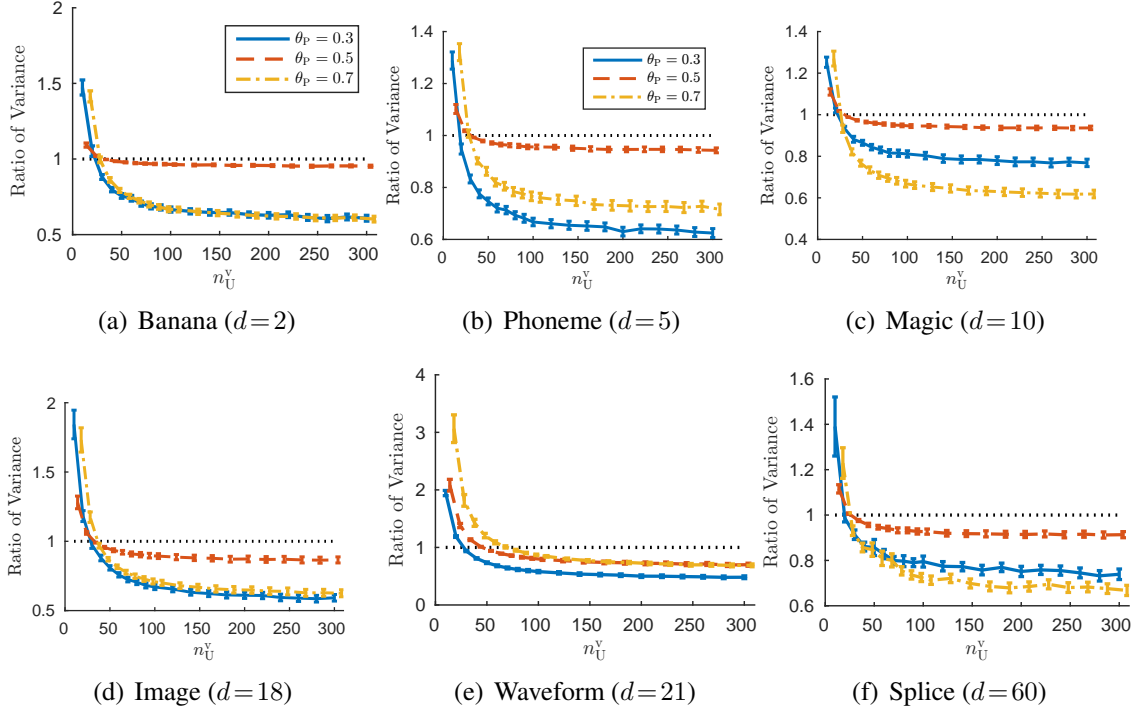


FIGURE 5.4: Average and standard error of the ratio between the variance of empirical PNU risk and that of PN risk, $\text{Var}[\hat{R}_{\text{PNU}}^{\text{MR},\eta}(\hat{g}_{\text{PN}})] / \text{Var}[\hat{R}_{\text{PN}}^{\text{MR}}(\hat{g}_{\text{PN}})]$, as a function of the number of unlabeled samples over 100 trials. Although the variance reduction is proved for an infinite number of samples, it can be observed with a finite number of samples.

5.6.2 Comparison with Existing Methods

In this subsection, we numerically compare our methods against existing semi-supervised classification methods.

Common Setup: We compare our methods against five conventional semi-supervised classification methods: the *entropy regularization* (ER) (Grandvalet and Bengio, 2004), the *Laplacian support vector machine* (LapSVM) (Belkin et al., 2006; Melacci and Belkin, 2011), the *squared-loss mutual information regularization* (SMIR) (Niu et al., 2013), the *weakly labeled support vector machine* (WellSVM) (Li et al., 2013), and the *safe semi-supervised support vector machine* (S4VM) (Li and Zhou, 2015). See Section 2.6 for the details of these methods.

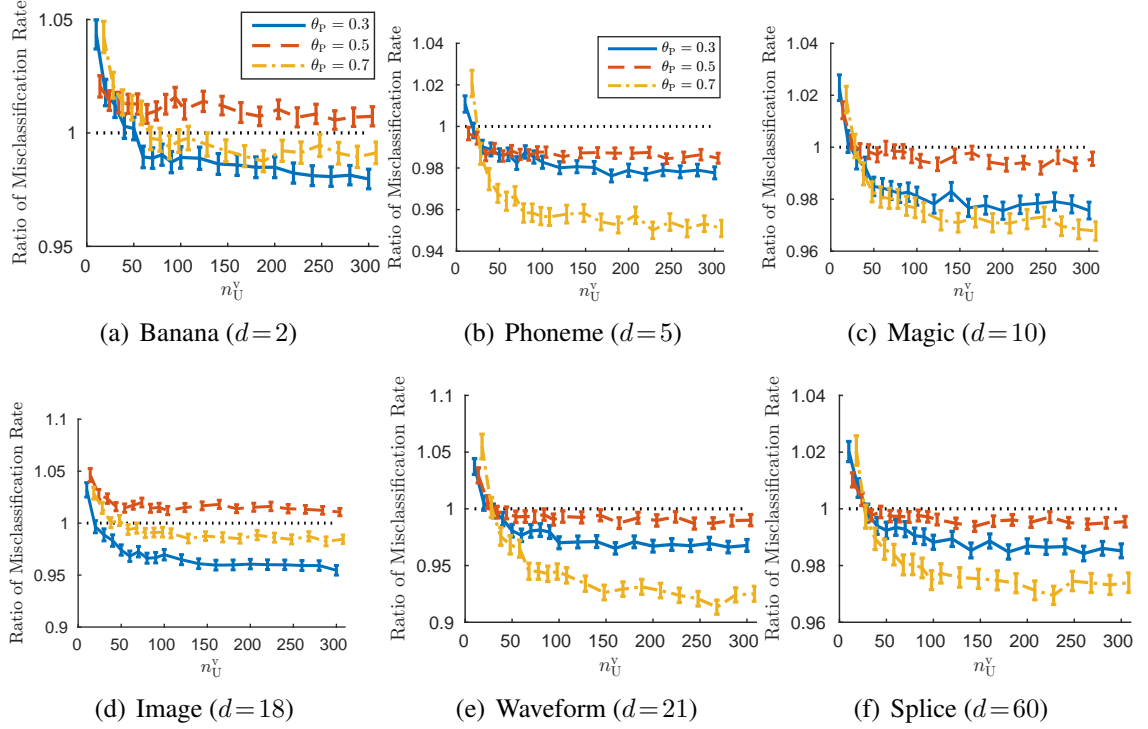


FIGURE 5.5: Average and standard error of the ratio between the misclassification rates of \hat{g}_{PN}^{PNU} and \hat{g}_{PN}^{PN} as a function of unlabeled samples over 1000 trials. In many cases, the ratio becomes less than 1, implying that the PNU risk is a promising alternative to the standard PN risk in validation if unlabeled data are available.

DataSets: We used sixteen benchmark datasets taken from the *UCI Machine Learning Repository* (Lichman, 2013), the *Semi-Supervised Learning book* (Chapelle et al., 2006), the *LIBSVM* (Chang and Lin, 2011), the *ELENA Project*,⁴ and the paper by Chapelle and Zien (2005).⁵ Each feature was scaled to $[0, 1]$. Similarly to the setting in Section 5.6.1, we used the Gaussian kernel model for all methods. The training data is $\{\mathbf{x}_i\}_{i=1}^n = \mathcal{X}_P \cup \mathcal{X}_N \cup \mathcal{X}_U$, where $n = n_P + n_N + n_U$. We selected all hyper-parameters with validation samples of size 20 ($n_P^V = n_N^V = 10$). For training, we drew n_L labeled and $n_U = 300$ unlabeled samples. The class-prior of labeled data was set at 0.7 and that of unlabeled samples was set at $\theta_P = 0.5$ that were assumed to be known. In practice, the class-prior, θ_P , can be estimated by methods proposed, e.g., by Blanchard et al. (2010), Ramaswamy et al. (2016), or Kawakubo et al. (2016).

⁴ <https://www.elen.ucl.ac.be/neural-nets/Research/Projects/ELENA/elena.htm>

⁵ <http://olivier.chapelle.cc/lds/>

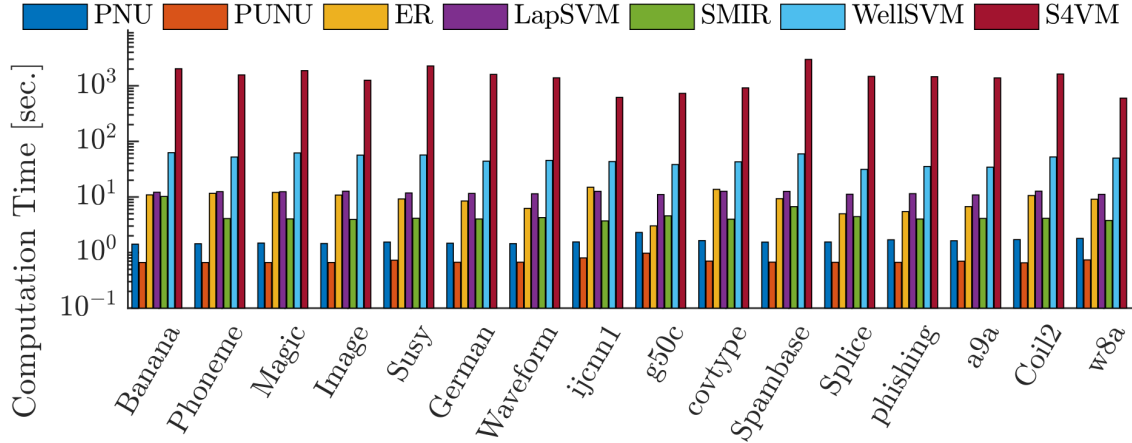


FIGURE 5.6: Average computation time over 50 trials for benchmark datasets when $n_L = 50$.

Table 5.2 lists the average with standard error of the misclassification rates over 50 trials and the number of best/comparable performances of each method in the bottom row. The superior performance of the PNU classification method over the PUNU classification method agrees well with the discussion in Section 5.3.3. For the g50c dataset, which well satisfies the low-density separation principle, WellSVM achieved the best performance. However, in the Banana dataset, where the two classes are highly overlapped, the performance of WellSVM was worse than the other methods. On the other hand, the PNU classification method achieved consistently better performance than or comparable performance to the other methods and its performance did not degenerate considerably across datasets. These results show that the idea of using PU classification in semi-supervised classification is promising.

Figure 5.6 plots the computation time, showing that our method was computationally efficient compared with the other methods.

5.6.3 Image Classification:

Next, we apply the PNU classification method to image classification. We used the *Places* 205 dataset (Zhou et al., 2014), which contains 2.5 million images in 205 scene classes. As a feature vector, we used a 4096-dimensional feature vector extracted from each image by *AlexNet* (Krizhevsky et al., 2012) under the framework of *Caffe*,⁶ which is available on the project website.⁷ We chose two scenes with similar names to construct binary

⁶ <http://caffe.berkeleyvision.org/>

⁷ <http://places.csail.mit.edu/>

TABLE 5.2: Average and standard error of the misclassification rates of each method over 50 trials for benchmark datasets. Boldface numbers denote the best and comparable methods in terms of average misclassification rate according to a t-test at a significance level of 5%. The bottom row gives the number of best/comparable cases of each method.

Dataset	n_L	PNU	PUNU	ER	LapSVM	SMIR	WellSVM	S4VM
Banana $d = 2$	10	30.1 (1.0)	32.1 (1.1)	35.8 (1.0)	36.9 (1.0)	37.7 (1.1)	41.8 (0.6)	45.3 (1.0)
	50	19.0 (0.6)	26.4 (1.2)	20.6 (0.7)	21.3 (0.7)	21.1 (1.0)	42.6 (0.5)	38.7 (0.9)
Phoneme $d = 5$	10	32.5 (0.8)	33.5 (1.0)	33.4 (1.2)	36.5 (1.5)	36.4 (1.2)	28.4 (0.6)	33.7 (1.4)
	50	28.1 (0.5)	32.8 (0.9)	27.8 (0.6)	27.0 (0.8)	28.6 (1.0)	26.8 (0.4)	25.1 (0.2)
Magic $d = 10$	10	31.7 (0.8)	34.1 (0.9)	34.2 (1.1)	37.9 (1.3)	36.0 (1.2)	30.1 (0.8)	33.3 (0.9)
	50	29.9 (0.8)	33.4 (0.9)	30.9 (0.5)	31.0 (0.9)	30.8 (0.9)	28.8 (0.8)	29.2 (0.4)
Image $d = 18$	10	29.8 (0.9)	31.7 (0.8)	33.7 (1.1)	36.6 (1.2)	36.7 (1.2)	34.7 (1.1)	35.9 (1.0)
	50	20.7 (0.8)	26.6 (1.1)	20.8 (0.8)	20.3 (1.0)	20.9 (0.9)	27.2 (1.0)	23.2 (0.7)
Susy $d = 18$	10	44.6 (0.6)	45.0 (0.6)	47.7 (0.4)	48.2 (0.4)	45.1 (0.7)	48.0 (0.3)	46.8 (0.3)
	50	38.9 (0.6)	41.5 (0.6)	37.9 (0.7)	43.1 (0.6)	43.9 (0.8)	43.8 (0.7)	42.1 (0.4)
German $d = 20$	10	40.8 (0.9)	42.4 (0.7)	43.6 (0.9)	45.9 (0.7)	46.2 (0.8)	42.4 (0.8)	42.0 (0.7)
	50	36.2 (0.8)	39.0 (0.8)	38.9 (0.6)	40.6 (0.6)	38.4 (1.1)	38.5 (1.0)	34.9 (0.5)
Waveform $d = 21$	10	17.4 (0.6)	18.0 (0.9)	18.5 (0.6)	24.9 (1.4)	18.0 (1.0)	16.7 (0.6)	20.8 (0.8)
	50	16.3 (0.6)	23.7 (1.2)	14.2 (0.4)	18.1 (0.8)	15.4 (0.6)	15.5 (0.5)	15.3 (0.3)
ijcnn1 $d = 22$	10	43.6 (0.6)	40.3 (1.0)	49.7 (0.1)	49.2 (0.3)	44.0 (1.0)	45.9 (0.7)	49.3 (0.8)
	50	34.5 (0.8)	37.1 (0.9)	35.5 (0.8)	33.4 (1.1)	49.4 (0.3)	46.2 (0.8)	48.6 (0.4)
g50c $d = 50$	10	11.4 (0.6)	12.5 (0.6)	23.3 (2.3)	39.8 (1.6)	21.9 (1.3)	6.6 (0.4)	27.0 (1.4)
	50	12.5 (1.1)	10.1 (0.6)	8.7 (0.4)	22.5 (1.5)	10.6 (0.6)	7.4 (0.4)	12.1 (0.5)
covtype $d = 54$	10	46.2 (0.4)	46.0 (0.4)	46.0 (0.5)	47.1 (0.5)	47.9 (0.5)	46.9 (0.6)	46.4 (0.4)
	50	41.3 (0.5)	42.3 (0.5)	41.0 (0.4)	41.5 (0.5)	46.2 (0.8)	43.6 (0.6)	40.8 (0.4)
Spambase $d = 57$	10	27.2 (0.9)	28.1 (1.1)	31.8 (1.4)	39.7 (1.4)	30.9 (1.3)	23.8 (0.8)	36.1 (1.5)
	50	23.4 (1.0)	26.6 (1.0)	22.1 (0.7)	28.5 (1.3)	20.9 (0.5)	19.1 (0.4)	24.5 (0.9)
Splice $d = 60$	10	38.3 (0.8)	39.3 (0.8)	43.9 (0.8)	47.9 (0.5)	41.6 (0.7)	42.0 (1.0)	42.4 (0.6)
	50	30.6 (0.8)	34.7 (0.9)	30.9 (0.8)	38.8 (1.0)	30.6 (0.9)	40.9 (0.8)	35.9 (0.7)
phishing $d = 68$	10	24.2 (1.2)	25.8 (1.0)	27.3 (1.6)	37.2 (1.6)	27.6 (1.6)	27.5 (1.4)	31.7 (1.3)
	50	15.8 (0.6)	18.3 (0.8)	15.4 (0.5)	21.1 (1.3)	14.7 (0.8)	17.2 (0.7)	16.7 (0.8)
a9a $d = 83$	10	31.4 (0.9)	31.3 (1.0)	34.3 (1.2)	41.0 (1.1)	37.3 (1.3)	33.1 (1.2)	34.3 (1.2)
	50	27.9 (0.6)	29.9 (0.8)	28.6 (0.7)	33.3 (1.0)	26.9 (0.7)	28.9 (0.8)	26.2 (0.4)
Coil2 $d = 241$	10	38.7 (0.8)	40.1 (0.8)	42.8 (0.7)	43.9 (0.8)	43.2 (0.8)	39.1 (0.9)	44.0 (0.8)
	50	23.2 (0.6)	30.5 (0.9)	23.6 (0.9)	22.8 (0.9)	25.1 (0.9)	22.6 (0.8)	25.4 (0.8)
w8a $d = 300$	10	35.9 (0.9)	33.6 (1.0)	41.6 (1.0)	46.6 (0.8)	39.4 (0.9)	42.1 (0.8)	43.0 (0.8)
	50	28.1 (0.7)	27.6 (0.6)	27.0 (0.9)	38.7 (0.8)	28.0 (0.9)	33.7 (0.8)	35.2 (1.0)
#Best/Comp.		23	13	11	4	9	13	7

classification tasks (see the description of datasets in Table 5.3). We drew 100 labeled and n_U unlabeled samples from each task; the class-prior of labeled and unlabeled data were respectively set at 0.5 and $\theta_P = m_P / (m_P + m_N)$, where m_P and m_N respectively denote the number of total samples in positive and negative scenes. We used a linear classifier $g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0$, where \mathbf{w} is the weight vector and w_0 is the offset (in the SMIR, the linear kernel model is used; see Section 2.6 for details).

TABLE 5.3: The description of the dataset used in the image classification experiment.

Dataset	Data sources	#Samples
Arts	Art Gallery	$(m_P = 15000)$
	vs. Art Studio	$(m_N = 15000)$
Deserts	Desert Sand	$(m_P = 15000)$
	vs. Desert Vegetation	$(m_N = 5556)$
Fields	Field Wild	$(m_P = 15000)$
	vs. Field Cultivated	$(m_N = 8117)$
Stadiums	Stadium Baseball	$(m_P = 15000)$
	vs. Stadium Football	$(m_N = 15000)$
Platforms	Subway Station	$(m_P = 5597)$
	vs. Train Station	$(m_N = 15000)$
Temples	Temple East Asia	$(m_P = 8691)$
	vs. Temple South Asia	$(m_N = 7178)$

The hyper-parameters in the PNU classification method were selected by five-fold cross-validation. The class-prior $p(y = +1) = \theta_P$ was estimated using the method based on energy distance minimization ([Kawakubo et al., 2016](#)).

Table 5.4 lists the average with standard error of the misclassification rates over 30 trials, where methods taking more than 2 hours were omitted and indicated as *N/A*. The results show that the PNU classification method was most effective. The average computation times are shown in Figure 5.7, revealing again that the PNU classification method was the fastest method.

TABLE 5.4: Average and standard error of misclassification rates over 30 trials for the Places 205 dataset. Boldface numbers denote the best and comparable methods in terms of the average misclassification rate according to a t-test at a significance level of 5%.

Dataset	n_U	θ_P	$\hat{\theta}_P$	PNU	ER	LapSVM	SMIR	WellSVM
Arts	1000	0.50	0.49 (0.01)	27.4 (1.3)	26.6 (0.5)	26.1 (0.7)	40.1 (3.9)	27.5 (0.5)
	5000	0.50	0.50 (0.01)	24.8 (0.6)	26.1 (0.5)	26.1 (0.4)	30.1 (1.6)	N/A
	10000	0.50	0.52 (0.01)	25.6 (0.7)	25.4 (0.5)	25.5 (0.6)	N/A	N/A
Deserts	1000	0.73	0.67 (0.01)	13.0 (0.5)	15.3 (0.6)	16.7 (0.8)	17.2 (0.8)	18.2 (0.7)
	5000	0.73	0.67 (0.01)	13.4 (0.4)	13.3 (0.5)	16.6 (0.6)	24.4 (0.6)	N/A
	10000	0.73	0.68 (0.01)	13.3 (0.5)	13.7 (0.6)	16.8 (0.8)	N/A	N/A
Fields	1000	0.65	0.57 (0.01)	22.4 (1.0)	26.2 (1.0)	26.6 (1.3)	28.2 (1.1)	26.6 (0.8)
	5000	0.65	0.57 (0.01)	20.6 (0.5)	22.6 (0.6)	24.7 (0.8)	29.6 (1.2)	N/A
	10000	0.65	0.57 (0.01)	21.6 (0.6)	22.5 (0.6)	25.0 (0.9)	N/A	N/A
Stadiums	1000	0.50	0.50 (0.01)	11.4 (0.4)	11.5 (0.5)	12.5 (0.5)	17.4 (3.6)	11.7 (0.4)
	5000	0.50	0.50 (0.01)	11.0 (0.5)	10.9 (0.3)	11.1 (0.3)	13.4 (0.7)	N/A
	10000	0.50	0.51 (0.00)	10.7 (0.3)	10.9 (0.3)	11.2 (0.2)	N/A	N/A
Platforms	1000	0.27	0.33 (0.01)	21.8 (0.5)	23.9 (0.6)	24.1 (0.5)	30.1 (2.3)	26.2 (0.8)
	5000	0.27	0.34 (0.01)	23.3 (0.8)	24.4 (0.7)	24.9 (0.7)	26.6 (0.3)	N/A
	10000	0.27	0.34 (0.01)	21.4 (0.5)	24.3 (0.6)	24.8 (0.5)	N/A	N/A
Temples	1000	0.55	0.51 (0.01)	43.9 (0.7)	43.9 (0.6)	43.4 (0.6)	50.7 (1.6)	44.3 (0.5)
	5000	0.55	0.54 (0.01)	43.4 (0.9)	43.0 (0.6)	43.1 (1.0)	43.6 (0.7)	N/A
	10000	0.55	0.50 (0.01)	45.2 (0.8)	44.4 (0.8)	44.2 (0.7)	N/A	N/A

5.7 Experiments on Imbalanced Classification

In this section, we numerically investigate the behavior of the PNU-AUC optimization method and evaluate their performance on various datasets.

As the classifier, we used the linear-in-parameter model. In all experiments except text classification tasks, we used the Gaussian kernel basis function expressed as $\phi_\ell(\mathbf{x}) = \exp(-\|\mathbf{x} - \mathbf{x}_\ell\|^2 / (2\sigma^2))$, where $\sigma > 0$ is the Gaussian bandwidth, $\{\mathbf{x}_\ell\}_{\ell=1}^b$ are the samples randomly selected from training samples $\{\mathbf{x}_i\}_{i=1}^n$ and n is the number of training samples. In text classification tasks, we used the linear kernel basis function: $\phi_\ell(\mathbf{x}) = \mathbf{x}^\top \mathbf{x}_\ell$. The number of basis functions was set at $b = \min(n, 200)$. The candidates of the Gaussian bandwidth were $\text{median}(\{\|\mathbf{x}_i - \mathbf{x}_j\|\}_{i,j=1}^n) \times \{1/8, 1/4, 1/2, 1, 2\}$ and those of the regularization parameter were $\{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$. All hyper-parameters were determined by five-fold cross-validation. As the loss function, we used the squared loss function $\ell_S(m) = (1 - m)^2$.

We used 15 benchmark datasets from the *IDA Benchmark Repository* (Rätsch et al.,

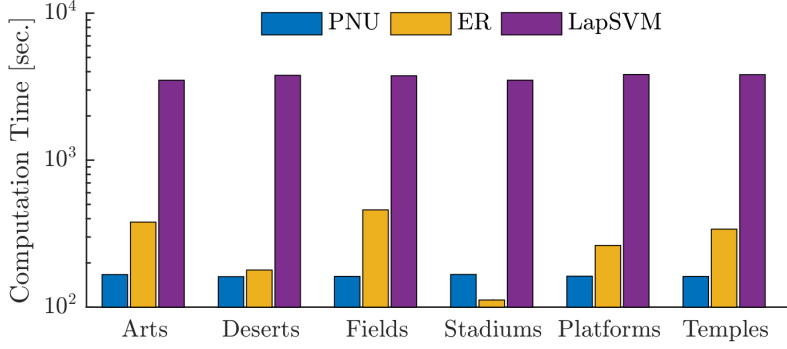


FIGURE 5.7: Average computation time over 30 trials for the Places 205 dataset when $n_U = 10000$.

2001), the *Semi-Supervised Learning* book (Chapelle et al., 2006), the *LIBSVM* (Chang and Lin, 2011), and the *UCI Machine Learning Repository* (Lichman, 2013). Furthermore, in text classification tasks, we used the *Reuters Corpus Volume I* dataset (Lewis et al., 2004), the *Amazon Review* dataset (Dredze et al., 2008), and the *20 Newsgroups* dataset (Lang, 1995). Since these text classification datasets are originally for multi-class classification, by dividing the whole classes into two group or choosing the two classes randomly, we constructed binary classification tasks: the rcv1, amazon2, and news20 datasets. The rcv1 and news20 datasets are available at the website of LIBSVM (Chang and Lin, 2011), and the amazon2 is designed by ourselves, which consists of the product reviews of books and music from the *Amazon* dataset (Blondel et al., 2013).

5.7.1 Experimental Analysis

Range of Variance Reduction: Firstly, we numerically confirm the effect of variance reduction. We compare the variance of the empirical PNU-AUC risk against the variance of the empirical PN-AUC risk, $\text{Var}[\hat{R}_{\text{PNU}}^{\text{AUC}, \eta}(g)]$ vs. $\text{Var}[\hat{R}_{\text{PN}}^{\text{AUC}}(g)]$, under a fixed classifier g .

As the fixed classifier, we used the minimizer of the empirical PN-AUC risk, denoted by \hat{g}_{PN} . The number of positive and negative samples for training varied as $(n_P, n_N) = (2, 8)$, $(10, 10)$, and $(18, 2)$. We then computed the variance of the empirical PN-AUC and PNU-AUC risks with additional 10 positive, 10 negative, and 300 unlabeled samples. As the dataset, we used the Banana dataset (Rätsch et al., 2001). In this experiment, the class-prior was set at $\theta_P = 0.1$ and assumed to be known.

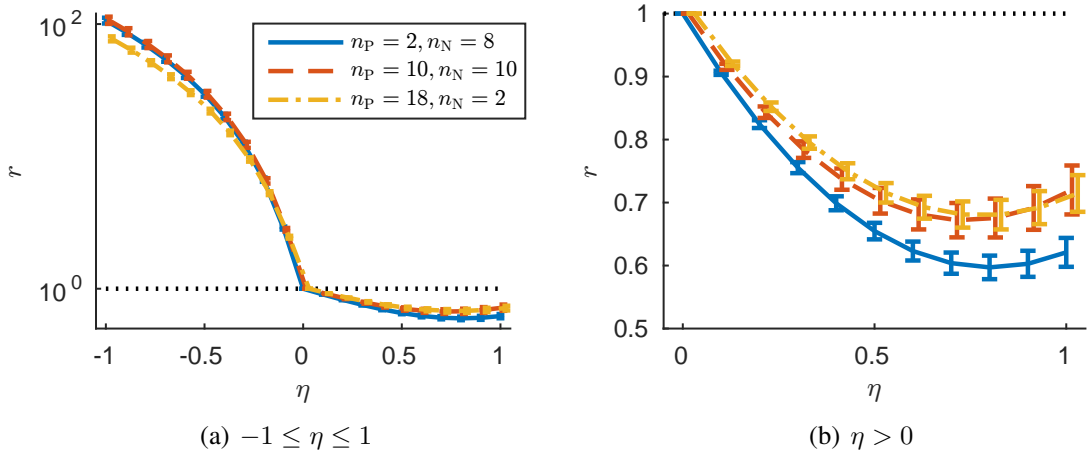


FIGURE 5.8: Average with standard error of the ratio between the variance of the empirical PNU risk and that of the PN risk, $r = \text{Var}[\hat{R}_{\text{PNU}}^{\text{AUC}, \eta}(\hat{g}_{\text{PN}})] / \text{Var}[\hat{R}_{\text{PN}}^{\text{AUC}}(\hat{g}_{\text{PN}})]$, as a function of the combination parameter η over 100 trials on the Banana dataset. The class-prior is $\theta_P = 0.1$ and the number of positive and negative samples varies as $(n_P, n_N) = (2, 8)$, $(10, 10)$, and $(18, 2)$. Left: values of r as a function of η . Right: values for $\eta > 0$ are magnified.

Figure 5.8 plots the value of the variance of the empirical PNU-AUC risk divided by that of the PN-AUC risk,

$$r := \frac{\text{Var}[\hat{R}_{\text{PNU}}^{\text{AUC}, \eta}(\hat{g}_{\text{PN}})]}{\text{Var}[\hat{R}_{\text{PN}}^{\text{AUC}}(\hat{g}_{\text{PN}})]}, \quad (5.94)$$

as a function of the combination parameter η under different numbers of positive and negative samples. The results show that $r < 1$ can be achieved by an appropriate choice of η , meaning that the variance of the empirical PNU-AUC risk can be smaller than that of the PN-AUC risk.

We then investigate how the class-prior affects the variance reduction. In this experiment, the number of positive and negative samples for \hat{g}_{PN} were set at $n_P = 10$ and $n_N = 10$, respectively. Figure 5.9 showed the values of r as a function of the combination parameter η under different class-priors. In particular, the variance can be reduced for $\eta > 0$ when the class-prior, θ_P , is 0.1 and 0.2. On the other hand, the range of the value of η that yields variance reduction becomes smaller when the class-prior is 0.3. However, this may not be that problematic in practice, because AUC optimization is effective when

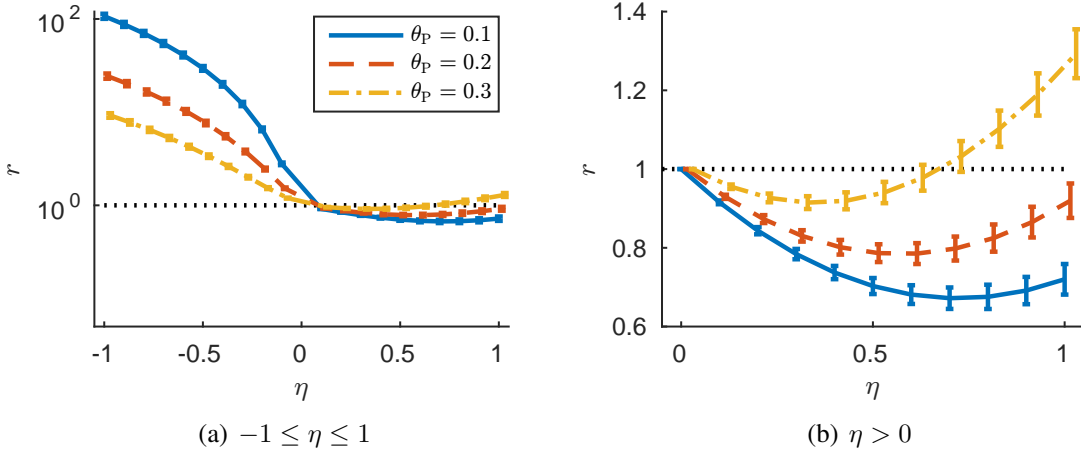


FIGURE 5.9: Average with standard error of the ratio between the variance of the PNU-AUC risk and that of the PN-AUC risk, $r = \text{Var}[\hat{R}_{\text{PNU}}^{\text{AUC}, \eta}(\hat{g}_{\text{PN}})] / \text{Var}[\hat{R}_{\text{PN}}^{\text{AUC}}(\hat{g}_{\text{PN}})]$, as a function of the combination parameter η over 100 trials on the Banana dataset. A class-prior varies as $\theta_P = 0.1, 0.2$, and 0.3 . Left: values of r as a function of η . Right: values for $\eta > 0$ are magnified. When $\theta_P = 0.1, 0.2$, the variance of the empirical PNU-AUC risk is smaller than that of the PN risk for $\eta > 0$.

two classes are highly imbalanced, i.e., the class-prior is far from 0.5; when the class-prior is close to 0.5, we may simply use the standard misclassification rate minimization approach.

Sensitivity Analysis: Secondly, we investigate the effect of the estimation accuracy of the class-prior for the PNU-AUC optimization method. Specifically, we used the class-prior $\hat{\theta}_P = \theta_P + \rho$, which is the addition of the true class-prior θ_P and noise $\rho \in \{-0.09, -0.08, \dots, 0.09\}$, as the estimated class-prior for the PNU-AUC optimization method. Under the different values of the class-prior $\theta_P = 0.1, 0.2$, and 0.3 , we trained a classifier with samples of size $n_P = \theta_P \cdot 50$, $n_N = \theta_N \cdot 50$, and $n_U = 1000$.

Figure 5.10 summarizes the average with standard error of the AUC as a function of the noise, showing that when $\theta_P = 0.2$ and 0.3 , the performance of the PNU-AUC optimization method is stable even when the estimated class-prior has some noise. In contrast, when $\theta_P = 0.1$, the performance largely decreases as the noise is close to $\rho = -0.09$. Since the true class-prior is small, it is sensitive to the negative bias. In particular, when $\rho = -0.09$, the ratio between the estimated and true class-priors is larger than other values. For instance, when $\rho = -0.09$ and $\theta_P = 0.2$, $\theta_P / \hat{\theta}_P \approx 1.8$, but when $\rho = -0.09$

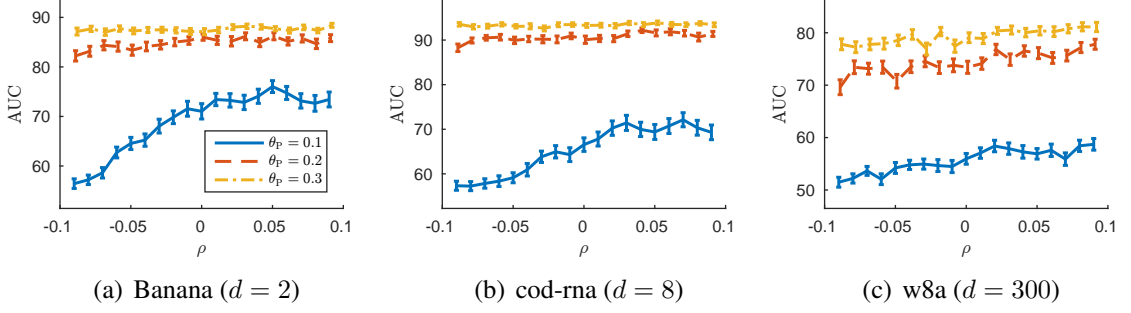


FIGURE 5.10: Average with standard error of the AUC as a function of the noise ρ over 100 trials. The PNU-AUC optimization method used the noisy class-prior $\hat{\theta}_P = \theta_P + \rho$ in training. The plots show that when $\theta_P = 0.2$ and 0.3 , the performance of the PNU-AUC optimization method is stable even when the estimated class-prior has some noise. However, when $\theta_P = 0.1$, as the noise is close to $\rho = -0.09$, the performance largely decreases.

Since the true class-prior is small, it is sensitive to the negative bias.

and $\theta_P = 0.1$, $\theta_P / \hat{\theta}_P \approx 10$. In contrast, the positive bias does not heavily affect the performance even when $\theta_P = 0.1$.

Scalability: Here, we report the scalability of the PNU-AUC optimization method. Specifically, we evaluated the AUC and computation time while increasing the number of unlabeled samples. We picked two large datasets: the SUSY and amazon2 datasets. The number of positive and negative samples were set at $n_P = 40$ and $n_N = 160$, respectively.

Figure 5.11 summarizes the average with standard error of the AUC and computation time as a function of the number of unlabeled samples. The AUC on the SUSY dataset slightly increased at $n_U = 1,000,000$, but the improvement on the amazon2 dataset was not noticeable or the performance decreased slightly. In this experiment, the increase of the size of unlabeled data did not improve the performance of the classifier significantly, but it did not affect adversely, i.e., it did not cause significant performance degeneration.

The result of computation time shows that the PNU-AUC optimization method can handle approximately 1,000,000 samples within reasonable computation time in this experiment. The longer computation time on the SUSY dataset before $n_U = 10,000$ is because we need to choose one additional hyperparameter, i.e., the bandwidth of the Gaussian kernel basis function, compared with the linear kernel basis function. However, the effect gradually decreases; after $n_U = 10,000$, the matrix multiplication of the high dimensional matrix on the amazon2 dataset ($d = 262,144$) requires more computation time

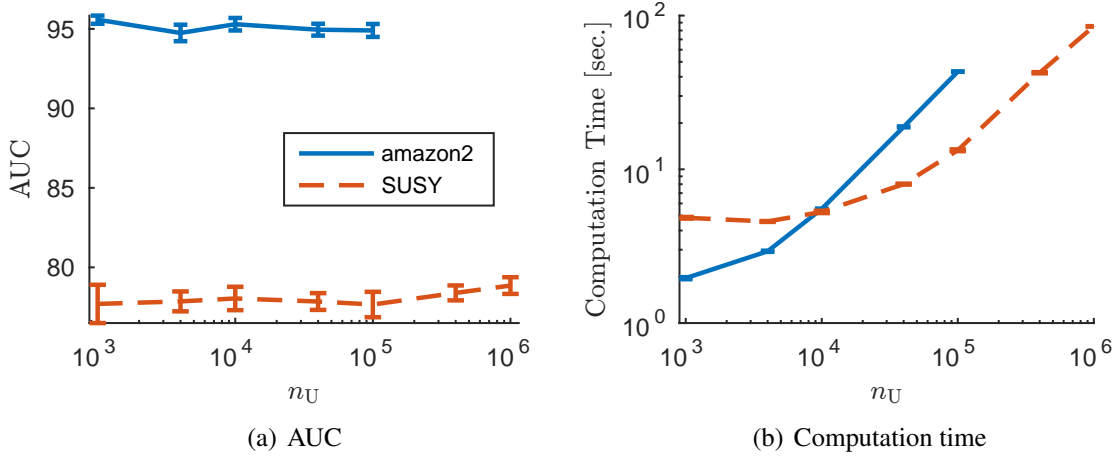


FIGURE 5.11: Average with standard error of the AUC as a function of the number of unlabeled data n_U over 20 trials.

than the SUSY dataset ($d = 18$).

5.7.2 Comparison with Existing Methods

In this subsection, we report the classification performance of the PNU-AUC optimization method.

Here, we compare the PNU-AUC optimization method against existing AUC optimization methods: the semi-supervised rankboost (SSRankboost) (Amini et al., 2008),⁸ the semi-supervised AUC-optimized logistic sigmoid (sAUC-LS) (Fujino and Ueda, 2016),⁹ and the optimum AUC with a generative model (OptAG) (Fujino and Ueda, 2016).

The classifier was trained with samples of size $n_P = \theta_P \cdot n_L$, $n_N = n_L - n_P$, and $n_U = 1000$, where n_L is the number of labeled samples. For the PNU-AUC optimization method, the squared loss function was used and the candidates of the combination parameter η were $\{-0.9, -0.8, \dots, 0.9\}$. For the class-prior estimation, we used the energy distance minimization method (Kawakubo et al., 2016). The results of the estimated class-prior are summarized in Table 5.5.

For SSRankboost, the discount factor and the number of neighbors were chosen from $\{10^{-3}, 10^{-2}, 10^{-1}\}$ and $\{2, 3, \dots, 7\}$, respectively. For sAUC-LS and OptAG, the regularization parameter for the entropy regularizer was chosen from $\{1, 10\}$. Furthermore,

⁸ We used the code available at <http://ama.liglab.fr/~amini/SSRankBoost/>

⁹ This method is equivalent to OptAG without a generative model, which only employs a discriminative model with the entropy minimization principle. To eliminate the adverse effect of the wrongly chosen generative model, we added this method for comparison.

TABLE 5.5: Average and standard error of the estimated class-prior over 50 trials on benchmark datasets in semi-supervised learning setting.

Dataset	n_L	$\theta_P = 0.1$	$\theta_P = 0.2$
Banana	50	0.12 (0.01)	0.21 (0.02)
($d = 2$)	100	0.10 (0.01)	0.20 (0.01)
skin_nonskin	50	0.11 (0.01)	0.21 (0.01)
($d = 3$)	100	0.10 (0.01)	0.20 (0.01)
cod-rna	50	0.12 (0.01)	0.22 (0.01)
($d = 8$)	100	0.12 (0.01)	0.21 (0.01)
Magic	50	0.09 (0.01)	0.17 (0.01)
($d = 10$)	100	0.07 (0.01)	0.20 (0.01)
Image	50	0.12 (0.01)	0.22 (0.01)
($d = 18$)	100	0.11 (0.01)	0.20 (0.01)
SUSY	50	0.10 (0.01)	0.20 (0.01)
($d = 18$)	100	0.10 (0.01)	0.19 (0.01)
Ringnorm	50	0.06 (0.00)	0.15 (0.00)
($d = 20$)	100	0.07 (0.00)	0.17 (0.00)
Twonorm	50	0.10 (0.00)	0.20 (0.00)
($d = 20$)	100	0.10 (0.00)	0.20 (0.00)
Waveform	50	0.11 (0.01)	0.20 (0.01)
($d = 21$)	100	0.09 (0.01)	0.19 (0.01)
covtype	50	0.09 (0.01)	0.20 (0.01)
($d = 54$)	100	0.09 (0.01)	0.18 (0.01)
phishing	50	0.10 (0.00)	0.20 (0.00)
($d = 68$)	100	0.10 (0.00)	0.20 (0.00)
a9a	50	0.10 (0.01)	0.20 (0.01)
($d = 83$)	100	0.10 (0.00)	0.21 (0.01)
mushrooms	50	0.10 (0.00)	0.20 (0.00)
($d = 112$)	100	0.10 (0.00)	0.20 (0.00)
USPS	50	0.10 (0.00)	0.20 (0.01)
($d = 241$)	100	0.09 (0.01)	0.19 (0.01)
w8a	50	0.10 (0.00)	0.19 (0.00)
($d = 300$)	100	0.09 (0.00)	0.20 (0.01)

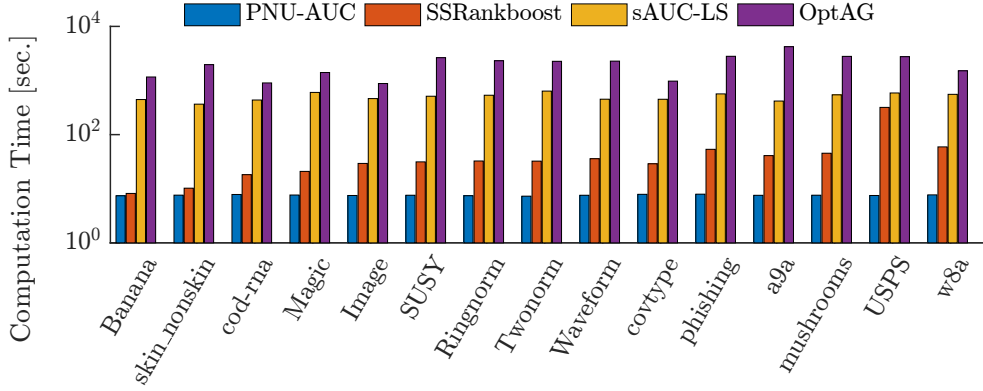


FIGURE 5.12: Average computation time of each method on benchmark datasets when $n_L = 100$ and $\theta_P = 0.1$ over 50 trials.

as the generative model of OptAG, we adapted the Gaussian distribution for the data distribution and the Gaussian and Gamma distributions for the prior of the data distribution.

Table 5.6 lists the average with standard error of the AUC over 50 trials, showing that the PNU-AUC optimization method achieves better performance than or comparable performance to the existing methods on many datasets. Figure 5.12 summarizes the average computation time over 50 trials. The computation time of the PNU-AUC optimization method includes both the class-prior estimation and the empirical risk minimization. The results show that even though the method involves the class-prior estimation, the computation time is relatively faster than SSRankboost and much faster than sAUC-LS and OptAG. The reason for longer computation time of sAUC-LS and OptAG is that their implementation is based on the logistic loss in which the number of operations for loss evaluation is $\mathcal{O}(n_P n_N + n_P n_U + n_N n_U)$, unlike the PNU-AUC optimization method with the squared loss in which the number of operations for loss evaluation is $\mathcal{O}(n_P + n_N + n_U)$ (cf. the discussion about the computational complexity in Section 3.4).

5.7.3 Text Classification

Finally, we report the performance of the PNU-AUC optimization method on text classification tasks.

A classifier was trained with samples of size $n_P = 20$, $n_N = 80$, and $n_U = 10,000$. The true class-prior was set at $\theta_P = 0.2$ and estimated by the method based on energy distance minimization (Kawakubo et al., 2016). For the generative model of OptAG, we employed naive Bayes (NB) multinomial models and a Dirichlet prior for the prior distribution of the NB model as described in Fujino and Ueda (2016).

TABLE 5.6: Average and standard error of the AUC over 50 trials on benchmark datasets. The boldface denotes the best and comparable methods in terms of the average AUC according to the t-test at the significance level 5%. The bottom row shows the number of best/comparable cases of each method. SSRboost is an abbreviation for SSRankboost.

Dataset	n_L	$\theta_P = 0.1$				$\theta_P = 0.2$			
		PNU-AUC	SSRboost	sAUC-LS	OptAG	PNU-AUC	SSRboost	sAUC-LS	OptAG
Banana ($d = 2$)	50	84.6 (1.3)	61.6 (1.0)	82.0 (1.3)	84.0 (1.3)	86.7 (0.9)	66.1 (0.9)	83.3 (1.5)	86.9 (0.7)
	100	88.4 (0.7)	66.2 (0.7)	84.0 (1.2)	88.7 (0.5)	91.2 (0.4)	69.1 (0.5)	88.5 (0.6)	89.9 (0.6)
skin_nonskin ($d = 3$)	50	96.5 (0.9)	92.2 (0.8)	96.4 (0.6)	97.8 (0.6)	98.1 (0.6)	93.7 (0.6)	98.8 (0.2)	99.1 (0.2)
	100	98.8 (0.3)	94.5 (0.4)	99.0 (0.1)	99.5 (0.0)	99.1 (0.3)	95.6 (0.2)	99.2 (0.3)	99.3 (0.2)
cod-rna ($d = 8$)	50	86.7 (1.4)	79.4 (1.0)	61.6 (1.5)	63.3 (1.4)	91.2 (1.0)	86.8 (0.6)	66.0 (1.5)	68.7 (1.3)
	100	93.4 (0.6)	88.7 (0.5)	67.3 (1.5)	67.7 (1.3)	95.6 (0.5)	91.7 (0.3)	74.1 (1.5)	75.9 (0.9)
Magic ($d = 10$)	50	77.0 (1.1)	74.8 (0.8)	73.7 (1.6)	77.4 (0.7)	77.2 (1.3)	78.2 (0.5)	76.5 (0.7)	75.4 (1.1)
	100	79.5 (0.5)	79.3 (0.6)	74.5 (1.6)	77.1 (1.0)	81.5 (0.4)	81.2 (0.4)	75.8 (1.0)	78.3 (0.5)
Image ($d = 18$)	50	81.6 (1.7)	70.0 (1.1)	76.5 (1.6)	80.7 (1.6)	86.1 (0.8)	78.8 (0.6)	81.5 (1.0)	82.9 (1.1)
	100	88.3 (0.6)	80.9 (0.7)	83.4 (1.2)	85.2 (1.0)	92.1 (0.4)	87.1 (0.5)	86.5 (0.6)	87.8 (0.4)
SUSY ($d = 18$)	50	63.0 (1.1)	67.7 (1.3)	56.2 (0.9)	56.2 (1.0)	64.7 (0.8)	70.9 (0.8)	55.6 (0.9)	57.7 (0.7)
	100	65.4 (1.1)	73.4 (0.5)	59.0 (0.9)	59.0 (0.8)	69.4 (1.0)	76.2 (0.4)	58.9 (0.7)	58.3 (0.6)
Ringnorm ($d = 20$)	50	98.3 (0.4)	79.9 (0.7)	98.7 (0.5)	99.0 (0.6)	98.4 (0.4)	86.6 (0.4)	99.5 (0.3)	99.1 (0.4)
	100	98.8 (0.3)	88.0 (0.4)	99.8 (0.0)	99.8 (0.0)	99.4 (0.2)	90.8 (0.3)	99.5 (0.4)	99.6 (0.2)
Twonorm ($d = 20$)	50	96.9 (0.6)	90.3 (0.4)	94.4 (1.1)	96.1 (0.3)	97.5 (0.5)	93.2 (0.3)	97.7 (0.2)	97.4 (0.2)
	100	98.6 (0.1)	94.7 (0.2)	96.6 (0.2)	96.5 (0.2)	99.0 (0.1)	96.8 (0.1)	98.3 (0.1)	98.0 (0.2)
Waveform ($d = 21$)	50	86.7 (1.3)	88.9 (0.4)	85.5 (1.5)	85.8 (0.8)	92.0 (0.6)	91.3 (0.2)	88.6 (0.8)	89.4 (0.6)
	100	92.8 (0.4)	91.8 (0.2)	87.8 (1.1)	87.7 (1.2)	94.7 (0.2)	93.1 (0.1)	90.1 (0.4)	89.7 (0.5)
covtype ($d = 54$)	50	57.8 (1.3)	63.1 (1.0)	55.7 (0.9)	58.9 (1.1)	60.3 (1.0)	65.6 (0.8)	56.4 (0.9)	57.8 (0.9)
	100	60.7 (1.1)	66.7 (0.8)	59.1 (0.9)	60.3 (0.9)	64.2 (0.6)	70.6 (0.6)	57.7 (0.9)	60.6 (0.7)
phishing ($d = 68$)	50	89.8 (1.1)	91.9 (0.6)	71.1 (1.6)	74.0 (0.9)	91.8 (0.7)	94.2 (0.3)	74.7 (1.5)	77.1 (1.2)
	100	94.6 (0.2)	94.8 (0.2)	76.3 (1.3)	76.9 (1.2)	94.9 (0.4)	96.3 (0.1)	80.5 (1.3)	82.0 (1.0)
a9a ($d = 83$)	50	69.4 (1.7)	75.1 (0.9)	75.3 (1.4)	73.8 (1.4)	78.2 (1.0)	79.3 (0.5)	78.3 (1.1)	79.2 (0.7)
	100	77.9 (1.1)	80.1 (0.5)	80.3 (0.5)	79.5 (0.8)	82.0 (0.5)	83.2 (0.3)	81.5 (0.5)	81.5 (0.7)
mushrooms ($d = 112$)	50	96.0 (0.7)	96.0 (0.6)	96.6 (0.5)	97.8 (0.3)	96.4 (0.7)	98.0 (0.2)	97.2 (0.5)	97.6 (0.8)
	100	98.1 (0.5)	98.2 (0.1)	97.9 (0.3)	98.9 (0.2)	98.7 (0.2)	98.7 (0.1)	98.2 (0.4)	99.1 (0.4)
USPS ($d = 241$)	50	85.6 (1.0)	73.9 (0.9)	79.5 (1.4)	82.7 (0.8)	88.2 (0.7)	80.5 (0.7)	81.2 (0.9)	85.1 (1.2)
	100	90.3 (0.5)	79.4 (0.7)	84.2 (1.0)	86.0 (1.0)	93.7 (0.5)	85.2 (0.5)	83.3 (0.9)	89.7 (0.3)
w8a ($d = 300$)	50	69.6 (1.1)	70.9 (0.9)	52.5 (0.9)	52.9 (1.1)	79.2 (0.9)	75.4 (0.9)	52.6 (0.9)	55.2 (0.9)
	100	79.9 (1.1)	77.1 (0.7)	53.1 (1.0)	53.9 (1.1)	85.6 (0.6)	81.0 (0.5)	55.7 (0.8)	56.2 (0.9)
#Best/Comp.		20	11	5	13	21	13	7	9

TABLE 5.7: Average with standard error of the AUC over 20 trials on the text classification datasets. The boldface denotes the best and comparable methods in terms of the average AUC according to the t-test at the significance level 5%.

Dataset	d	$\hat{\theta}_p$	PNU-AUC	SSRankboost	sAUC-LS	OptAG
rcv1	47236	0.21 (0.00)	91.1 (0.7)	76.6 (0.7)	79.6 (2.3)	79.4 (2.0)
amazon2	262144	0.21 (0.01)	87.2 (1.3)	68.7 (1.8)	55.9 (0.4)	56.8 (0.5)
news20	1355191	0.20 (0.00)	79.8 (1.1)	74.2 (2.0)	66.4 (1.6)	71.5 (1.2)

Table 5.7 lists the average with standard error of the AUC over 20 trials, showing that the PNU-AUC optimization method outperforms the existing methods. Figure 5.13 summarizes the average computation time of each method. These results show that the PNU-AUC optimization method achieves better performance with short computation time.

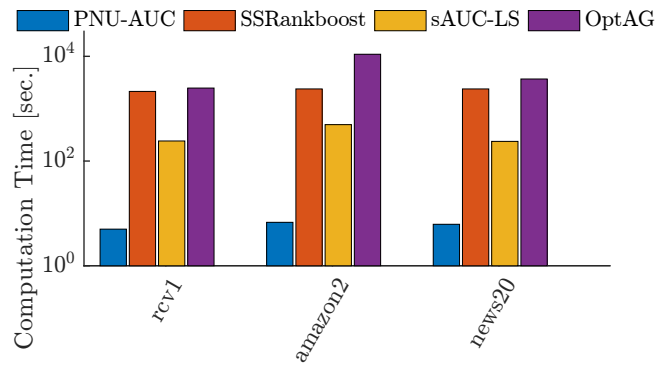


FIGURE 5.13: Average computation time of each method on the text classification datasets.

5.8 Proofs of Theorems

In this section, we give the proofs of Theorems in Section 5.5.

5.8.1 Proof of Theorem 13

To simplify the notation, let $R_P(g)$, $R_N(g)$, $R_{U,P}(g)$, and $R_{U,N}(g)$ be the risks of classifier g under loss ℓ :

$$R_P(g) := \mathbb{E}_P[\ell(g(\mathbf{x}))], \quad (5.95)$$

$$R_N(g) := \mathbb{E}_N[\ell(-g(\mathbf{x}))], \quad (5.96)$$

$$R_{U,P}(g) := \mathbb{E}_U[\ell(g(\mathbf{x}))], \quad (5.97)$$

$$R_{U,N}(g) := \mathbb{E}_U[\ell(-g(\mathbf{x}))], \quad (5.98)$$

Recall that

$$\begin{aligned} R_{N-PUNU}^{\text{MR},\gamma}(g) &= (1 - \gamma)R_{N-PU}^{\text{MR}}(g) + \gamma R_{N-NU}^{\text{MR}}(g) \\ &= (2 - 2\gamma)\theta_P R_P(g) + 2\gamma\theta_N R_N(g) \\ &\quad + (1 - \gamma)R_{U,N}(g) + \gamma R_{U,P}(g) + \text{Const}, \end{aligned} \quad (5.99)$$

$$\begin{aligned} R_{N-PNPU}^{\text{MR},\gamma}(g) &= (1 - \gamma)R_{PN}^{\text{MR}}(g) + \gamma R_{N-PU}^{\text{MR}}(g) \\ &= (1 + \gamma)\theta_P R_P(g) + (1 - \gamma)\theta_N R_N(g) + \gamma R_{U,N}(g) + \text{Const}, \end{aligned} \quad (5.100)$$

$$\begin{aligned} R_{N-PNNU}^{\text{MR},\gamma}(g) &= (1 - \gamma)R_{PN}^{\text{MR}}(g) + \gamma R_{N-NU}^{\text{MR}}(g) \\ &= (1 - \gamma)\theta_P R_P(g) + (1 + \gamma)\theta_N R_N(g) + \gamma R_{U,P}(g) + \text{Const}. \end{aligned} \quad (5.101)$$

Let $\hat{R}_P(g)$, $\hat{R}_N(g)$, $\hat{R}_{U,P}(g)$ and $\hat{R}_{U,N}(g)$ be the empirical risks. In order to prove Theorem 13, the following concentration lemma is needed:

Lemma 19. *For any $\delta > 0$, we have these uniform deviation bounds with probability at least $1 - \delta/3$:*

$$\sup_{g \in \mathcal{G}} (R_P(g) - \hat{R}_P(g)) \leq \frac{C_w C_\phi}{\sqrt{n_P}} + \sqrt{\frac{\ln(3/\delta)}{2n_P}}, \quad (5.102)$$

$$\sup_{g \in \mathcal{G}} (R_N(g) - \hat{R}_N(g)) \leq \frac{C_w C_\phi}{\sqrt{n_N}} + \sqrt{\frac{\ln(3/\delta)}{2n_N}}, \quad (5.103)$$

$$\sup_{g \in \mathcal{G}} (R_{U,P}(g) - \hat{R}_{U,P}(g)) \leq \frac{C_w C_\phi}{\sqrt{n_U}} + \sqrt{\frac{\ln(3/\delta)}{2n_U}}, \quad (5.104)$$

$$\sup_{g \in \mathcal{G}} (R_{U,N}(g) - \hat{R}_{U,N}(g)) \leq \frac{C_w C_\phi}{\sqrt{n_U}} + \sqrt{\frac{\ln(3/\delta)}{2n_U}}. \quad (5.105)$$

All inequalities in Lemma 19 are from the basic *uniform deviation bound* using the Rademacher complexity (Mohri et al., 2012), *Talagrand's contraction lemma* (Ledoux and Talagrand, 1991), as well as the fact that the Lipschitz constant of ℓ_R is $1/2$. For these reasons, the detailed proof of Lemma 19 is omitted.

Consider $R_{N\text{-PNPU}}^{\text{MR},\gamma}(g)$. It is clear that

$$\begin{aligned} \sup_{g \in \mathcal{G}} (R_{N\text{-PNPU}}^{\text{MR},\gamma}(g) - \widehat{R}_{N\text{-PNPU}}^{\text{MR},\gamma}(g)) &\leq (1 + \gamma)\theta_P \sup_{g \in \mathcal{G}} (R_P(g) - \widehat{R}_P(g)) \\ &\quad + (1 - \gamma)\theta_N \sup_{g \in \mathcal{G}} (R_N(g) - \widehat{R}_N(g)) \\ &\quad + \gamma \sup_{g \in \mathcal{G}} (R_{U,N}(g) - \widehat{R}_{U,N}(g)). \end{aligned} \quad (5.106)$$

Therefore, by applying Lemma 19, for any $\delta > 0$, it holds with probability at least $1 - \delta$ that

$$\sup_{g \in \mathcal{G}} (R_{N\text{-PNPU}}^{\text{MR},\gamma}(g) - \widehat{R}_{N\text{-PNPU}}^{\text{MR},\gamma}(g)) \leq \frac{1}{2} c_1^{\text{MR}}(\delta) \left(\frac{(1 + \gamma)\theta_P}{\sqrt{n_P}} + \frac{(1 - \gamma)\theta_N}{\sqrt{n_N}} + \frac{\gamma}{\sqrt{n_U}} \right). \quad (5.107)$$

Since $I^{\text{MR}}(g) \leq 2R_{N\text{-PNPU}}^{\text{MR},\gamma}(g)$, with the same probability,

$$\sup_{g \in \mathcal{G}} (I(g) - 2\widehat{R}_{N\text{-PNPU}}^{\text{MR},\gamma}(g)) \leq c_1^{\text{MR}}(\delta) \left(\frac{(1 + \gamma)\theta_P}{\sqrt{n_P}} + \frac{(1 - \gamma)\theta_N}{\sqrt{n_N}} + \frac{\gamma}{\sqrt{n_U}} \right). \quad (5.108)$$

Similarly,

$$\sup_{g \in \mathcal{G}} (I^{\text{MR}}(g) - 2\widehat{R}_{N\text{-PUNU}}^{\text{MR},\gamma}(g)) \leq c_1^{\text{MR}}(\delta) \left(\frac{(1 - \gamma)\theta_P}{\sqrt{n_P}} + \frac{(1 + \gamma)\theta_N}{\sqrt{n_N}} + \frac{\gamma}{\sqrt{n_U}} \right) \quad (5.109)$$

with probability at least $1 - \delta$.

Finally, $R_{N\text{-PUNU}}^{\text{MR},\gamma}(g)$ is slightly more involved, for that there are both $R_{U,P}(g)$ and $R_{U,N}(g)$. From $\ell_R(m) + \ell_R(-m) = 1$, we can know $R_{U,P}(g) + R_{U,N}(g) = 1$ and then

$$(1 - \gamma)R_{U,N}(g) + \gamma R_{U,P}(g) = \begin{cases} (2\gamma - 1)R_{U,P}(g) + \text{Const} & \gamma \geq 1/2, \\ (1 - 2\gamma)R_{U,N}(g) + \text{Const} & \gamma < 1/2. \end{cases} \quad (5.110)$$

As a result,

$$\sup_{g \in \mathcal{G}} (I^{\text{MR}}(g) - 2\widehat{R}_{N\text{-PUNU}}^{\text{MR},\gamma}(g)) \leq c_1^{\text{MR}}(\delta) \left(\frac{(2 - 2\gamma)\theta_P}{\sqrt{n_P}} + \frac{2\gamma}{\sqrt{n_N}} + \frac{|2\gamma - 1|}{\sqrt{n_U}} \right) \quad (5.111)$$

with probability at least $1 - \delta$.

5.8.2 Proof of Theorem 14

In fact,

$$\ell_{\text{TS}}(m) = \begin{cases} 1/4 & m \leq 0, \\ (m-1)^2/4 & 0 < m \leq 1, \\ 0 & m > 1, \end{cases} \quad (5.112)$$

and after plugging this $\ell_{\text{TS}}(m)$ into $\tilde{\ell}_{\text{TS}}(m)$,

$$\begin{aligned} \tilde{\ell}_{\text{TS}}(m) &= \ell_{\text{TS}}(m) - \ell_{\text{TS}}(-m) \\ &= \begin{cases} 1/4 & m \leq -1, \\ 1/4 - (m+1)^2/4 & -1 < m \leq 0, \\ (m-1)^2/4 - 1/4 & 0 < m \leq 1, \\ -1/4 & m > 1. \end{cases} \end{aligned} \quad (5.113)$$

It is easy to see that $\ell_{\text{TS}}(m)$ and $\tilde{\ell}_{\text{TS}}(m)$ are Lipschitz continuous with the same Lipschitz constant $1/2$.

Next, recall that

$$\begin{aligned} R_{\text{C-PUNU}}^{\text{MR},\gamma}(g) &= (1-\gamma)R_{\text{C-PU}}^{\text{MR}}(g) + \gamma R_{\text{C-NU}}^{\text{MR}}(g) \\ &= (1-\gamma)\theta_{\text{P}}R'_{\text{P}}(g) + \gamma\theta_{\text{N}}R'_{\text{N}}(g) + (1-\gamma)R_{\text{U,N}}(g) + \gamma R_{\text{U,P}}(g), \end{aligned} \quad (5.114)$$

$$\begin{aligned} R_{\text{C-PNPU}}^{\text{MR},\gamma}(g) &= (1-\gamma)R_{\text{PN}}^{\text{MR}}(g) + \gamma R_{\text{C-PU}}^{\text{MR}}(g) \\ &= (1-\gamma)\theta_{\text{P}}R_{\text{P}}(g) + (1-\gamma)\theta_{\text{N}}R_{\text{N}}(g) + \gamma\theta_{\text{P}}R'_{\text{P}}(g) + \gamma R_{\text{U,N}}(g), \end{aligned} \quad (5.115)$$

$$\begin{aligned} R_{\text{C-PNNU}}^{\text{MR},\gamma}(g) &= (1-\gamma)R_{\text{PN}}^{\text{MR}}(g) + \gamma R_{\text{C-NU}}^{\text{MR}}(g) \\ &= (1-\gamma)\theta_{\text{P}}R_{\text{P}}(g) + (1-\gamma)\theta_{\text{N}}R_{\text{N}}(g) + \gamma\theta_{\text{N}}R'_{\text{N}}(g) + \gamma R_{\text{U,P}}(g). \end{aligned} \quad (5.116)$$

Let $\hat{R}_{\text{P}}(g)$, $\hat{R}_{\text{N}}(g)$, $\hat{R}_{\text{U,P}}(g)$, $\hat{R}_{\text{U,N}}(g)$, $\hat{R}'_{\text{P}}(g)$ and $\hat{R}'_{\text{N}}(g)$ be the empirical risks. Again, the following concentration lemma is needed:

Lemma 20. *For any $\delta > 0$, we have these uniform deviation bounds with probability at least $1 - \delta/4$:*

$$\sup_{g \in \mathcal{G}} (R_P(g) - \hat{R}_P(g)) \leq \frac{C_w C_\phi}{\sqrt{n_P}} + \sqrt{\frac{\ln(4/\delta)}{32n_P}}, \quad (5.117)$$

$$\sup_{g \in \mathcal{G}} (R_N(g) - \hat{R}_N(g)) \leq \frac{C_w C_\phi}{\sqrt{n_N}} + \sqrt{\frac{\ln(4/\delta)}{32n_N}}, \quad (5.118)$$

$$\sup_{g \in \mathcal{G}} (R_{U,P}(g) - \hat{R}_{U,P}(g)) \leq \frac{C_w C_\phi}{\sqrt{n_U}} + \sqrt{\frac{\ln(4/\delta)}{32n_U}}, \quad (5.119)$$

$$\sup_{g \in \mathcal{G}} (R_{U,N}(g) - \hat{R}_{U,N}(g)) \leq \frac{C_w C_\phi}{\sqrt{n_U}} + \sqrt{\frac{\ln(4/\delta)}{32n_U}}, \quad (5.120)$$

$$\sup_{g \in \mathcal{G}} (R'_P(g) - \hat{R}'_P(g)) \leq \frac{C_w C_\phi}{\sqrt{n_P}} + \sqrt{\frac{\ln(4/\delta)}{8n_P}}, \quad (5.121)$$

$$\sup_{g \in \mathcal{G}} (R'_N(g) - \hat{R}'_N(g)) \leq \frac{C_w C_\phi}{\sqrt{n_N}} + \sqrt{\frac{\ln(4/\delta)}{8n_N}}. \quad (5.122)$$

The detailed proof of Lemma 20 is omitted for the same reason as Lemma 19. The difference is due to that $0 \leq \ell_{TS}(m) \leq 1/4$ and $-1/4 \leq \tilde{\ell}_{TS}(m) \leq 1/4$ whereas $0 \leq \ell_R(m) \leq 1$ just like $0 \leq \ell_{0-1}(m) \leq 1$. For convenience, we will relax $1/32$ to $1/8$ in the square root for $R_P(g)$, $R_N(g)$, $R_{U,P}(g)$, $R_{U,N}(g)$.

Consider $R_{C-PUNU}^{MR,\gamma}(g)$. By applying Lemma 20, for any $\delta > 0$, it holds with probability at least $1 - \delta$ that

$$\sup_{g \in \mathcal{G}} (R_{C-PUNU}^{MR,\gamma}(g) - \hat{R}_{C-PUNU}^{MR,\gamma}(g)) \leq \frac{1}{4} c_2^{MR}(\delta) \left(\frac{(1-\gamma)\theta_P}{\sqrt{n_P}} + \frac{\gamma\theta_N}{\sqrt{n_N}} + \frac{1}{\sqrt{n_U}} \right). \quad (5.123)$$

Since $I(g) \leq 4R_{C-PUNU}^{MR,\gamma}$, with the same probability,

$$\sup_{g \in \mathcal{G}} (I^{MR}(g) - 4\hat{R}_{C-PUNU}^{MR,\gamma}(g)) \leq c_2^{MR}(\delta) \left(\frac{(1-\gamma)\theta_P}{\sqrt{n_P}} + \frac{\gamma\theta_N}{\sqrt{n_N}} + \frac{1}{\sqrt{n_U}} \right). \quad (5.124)$$

The other two generalization error bounds can be proven similarly.

5.8.3 Proof of Theorem 15

Next, we prove the generalization error bounds of the PNPU-AUC and PNNU-AUC risks in Theorem 15. The proof techniques are similar to Theorem 5 for PU-AUC optimization in Chapter 3. We first prove the following risk bounds:

Lemma 21. *For any $\delta > 0$, the following inequalities hold separately with probability at least $1 - \delta$ for all $g \in \mathcal{G}$:*

$$\begin{aligned} R_{\text{PNPU}}^{\text{AUC},\gamma}(g) - \widehat{R}_{\text{PNPU}}^{\text{AUC},\gamma}(g) \\ \leq c^{\text{AUC}}(\delta/3) \left(\frac{1-\gamma}{\sqrt{\min(n_P, n_N)}} + \frac{\gamma}{\theta_N \sqrt{\min(n_P, n_U)}} + \frac{\theta_P \gamma}{\theta_N \sqrt{n_P}} \right), \end{aligned} \quad (5.125)$$

$$\begin{aligned} R_{\text{PNNU}}^{\text{AUC},\gamma}(g) - \widehat{R}_{\text{PNNU}}^{\text{AUC},\gamma}(g) \\ \leq c^{\text{AUC}}(\delta/3) \left(\frac{1-\gamma}{\sqrt{\min(n_P, n_N)}} + \frac{\gamma}{\theta_P \sqrt{\min(n_N, n_U)}} + \frac{\theta_N \gamma}{\theta_P \sqrt{n_N}} \right). \end{aligned} \quad (5.126)$$

Proof. Recall the PNPU-AUC and PNNU-AUC risks:

$$R_{\text{PNPU}}^{\text{AUC},\gamma}(g) := (1-\gamma)R_{\text{PN}}^{\text{AUC}}(g) + \gamma R_{\text{PU}}^{\text{AUC}}(g), \quad (5.127)$$

$$R_{\text{PNNU}}^{\text{AUC},\gamma}(g) := (1-\gamma)R_{\text{PN}}^{\text{AUC}}(g) + \gamma R_{\text{NU}}^{\text{AUC}}(g). \quad (5.128)$$

Based on Theorem 6 in Chapter 3, for any $\delta > 0$, we have these uniform deviation bounds with probability at least $1 - \delta/3$:

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left(\mathbb{E}_{\mathbf{P}}[\mathbb{E}_{\mathbf{N}}[\ell(g(\mathbf{x}^{\mathbf{P}}) - g(\mathbf{x}^{\mathbf{N}}))]] - \frac{1}{n_{\mathbf{P}}n_{\mathbf{N}}} \sum_{i=1}^{n_{\mathbf{P}}} \sum_{j=1}^{n_{\mathbf{N}}} \ell(g(\mathbf{x}_i^{\mathbf{P}}) - g(\mathbf{x}_j^{\mathbf{N}})) \right) \\ \leq \frac{c^{\text{AUC}}(\delta/3)}{\sqrt{\min(n_{\mathbf{P}}, n_{\mathbf{N}})}}, \end{aligned} \quad (5.129)$$

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left(\mathbb{E}_{\mathbf{P}}[\mathbb{E}_{\mathbf{U}}[\ell(g(\mathbf{x}^{\mathbf{P}}) - g(\mathbf{x}^{\mathbf{U}}))]] - \frac{1}{n_{\mathbf{P}}n_{\mathbf{U}}} \sum_{i=1}^{n_{\mathbf{P}}} \sum_{k=1}^{n_{\mathbf{U}}} \ell(g(\mathbf{x}_i^{\mathbf{P}}) - g(\mathbf{x}_k^{\mathbf{U}})) \right) \\ \leq \frac{c^{\text{AUC}}(\delta/3)}{\sqrt{\min(n_{\mathbf{P}}, n_{\mathbf{U}})}}, \end{aligned} \quad (5.130)$$

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left(\mathbb{E}_{\mathbf{U}}[\mathbb{E}_{\mathbf{N}}[\ell(g(\mathbf{x}^{\mathbf{U}}) - g(\mathbf{x}^{\mathbf{N}}))]] - \frac{1}{n_{\mathbf{N}}n_{\mathbf{U}}} \sum_{k=1}^{n_{\mathbf{U}}} \sum_{j=1}^{n_{\mathbf{N}}} \ell(g(\mathbf{x}_k^{\mathbf{U}}) - g(\mathbf{x}_j^{\mathbf{N}})) \right) \\ \leq \frac{c^{\text{AUC}}(\delta/3)}{\sqrt{\min(n_{\mathbf{N}}, n_{\mathbf{U}})}}, \end{aligned} \quad (5.131)$$

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left(\mathbb{E}_{\mathbf{P}}[\mathbb{E}_{\bar{\mathbf{P}}}[\ell(g(\mathbf{x}^{\mathbf{P}}) - g(\bar{\mathbf{x}}^{\mathbf{P}}))]] - \frac{1}{n_{\mathbf{P}}^2} \sum_{i=1}^{n_{\mathbf{P}}} \sum_{i'=1}^{n_{\mathbf{P}}} \ell(g(\mathbf{x}_i^{\mathbf{P}}) - g(\bar{\mathbf{x}}_{i'}^{\mathbf{P}})) \right) \\ \leq \frac{c^{\text{AUC}}(\delta/3)}{\sqrt{n_{\mathbf{P}}}}, \end{aligned} \quad (5.132)$$

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left(\mathbb{E}_{\mathbf{N}}[\mathbb{E}_{\bar{\mathbf{N}}}[\ell(g(\mathbf{x}^{\mathbf{N}}) - g(\bar{\mathbf{x}}^{\mathbf{N}}))]] - \frac{1}{n_{\mathbf{N}}^2} \sum_{j=1}^{n_{\mathbf{N}}} \sum_{j'=1}^{n_{\mathbf{N}}} \ell(g(\mathbf{x}_j^{\mathbf{N}}) - g(\bar{\mathbf{x}}_{j'}^{\mathbf{N}})) \right) \\ \leq \frac{c^{\text{AUC}}(\delta/3)}{\sqrt{n_{\mathbf{N}}}}. \end{aligned} \quad (5.133)$$

Combining three bounds from the above, for any $\delta > 0$, with probability $1 - \delta$, we have

$$\begin{aligned}
& \sup_{g \in \mathcal{G}} \left(R_{\text{PNPU}}^{\text{AUC}, \gamma}(g) - \widehat{R}_{\text{PNPU}}^{\text{AUC}, \gamma}(g) \right) \\
& \leq (1 - \gamma) \sup_{g \in \mathcal{G}} \left(\mathbb{E}_{\mathbf{P}}[\mathbb{E}_{\mathbf{N}}[\ell(g(\mathbf{x}^{\mathbf{P}}) - g(\mathbf{x}^{\mathbf{N}}))]] - \frac{1}{n_{\mathbf{P}} n_{\mathbf{N}}} \sum_{i=1}^{n_{\mathbf{P}}} \sum_{j=1}^{n_{\mathbf{N}}} \ell(g(\mathbf{x}_i^{\mathbf{P}}) - g(\mathbf{x}_j^{\mathbf{N}})) \right) \\
& \quad + \frac{\gamma}{\theta_{\mathbf{N}}} \sup_{g \in \mathcal{G}} \left(\mathbb{E}_{\mathbf{P}}[\mathbb{E}_{\mathbf{U}}[\ell(g(\mathbf{x}^{\mathbf{P}}) - g(\mathbf{x}^{\mathbf{U}}))]] - \frac{1}{n_{\mathbf{P}} n_{\mathbf{U}}} \sum_{i=1}^{n_{\mathbf{P}}} \sum_{k=1}^{n_{\mathbf{U}}} \ell(g(\mathbf{x}_i^{\mathbf{P}}) - g(\mathbf{x}_k^{\mathbf{U}})) \right) \\
& \quad + \frac{\gamma \theta_{\mathbf{P}}}{\theta_{\mathbf{N}}} \sup_{g \in \mathcal{G}} \left(\mathbb{E}_{\mathbf{P}}[\mathbb{E}_{\overline{\mathbf{P}}}[\ell(g(\mathbf{x}^{\mathbf{P}}) - g(\overline{\mathbf{x}}^{\mathbf{P}}))]] - \frac{1}{n_{\mathbf{P}}^2} \sum_{i=1}^{n_{\mathbf{P}}} \sum_{i'=1}^{n_{\mathbf{P}}} \ell(g(\mathbf{x}_i^{\mathbf{P}}) - g(\overline{\mathbf{x}}_{i'}^{\mathbf{P}})) \right) \\
& \leq c^{\text{AUC}}(\delta/3) \left(\frac{1 - \gamma}{\sqrt{\min(n_{\mathbf{P}}, n_{\mathbf{N}})}} + \frac{\gamma}{\theta_{\mathbf{N}} \sqrt{\min(n_{\mathbf{P}}, n_{\mathbf{U}})}} + \frac{\gamma \theta_{\mathbf{P}}}{\theta_{\mathbf{N}} \sqrt{n_{\mathbf{P}}}} \right). \tag{5.134}
\end{aligned}$$

This concludes the risk bounds of the PNPU-AUC risk.

Similarly, we prove the risk bounds of the PNNU-AUC risk. \square

Again, $I^{\text{AUC}}(g) \leq R^{\text{AUC}}(g)$ holds in our setting. This leads to Theorem 15.

5.8.4 Proofs of Theorems 16 and 17

Note that g is independent of the data for evaluating $\widehat{R}_{\text{N-PUNU}}^{\text{MR}, \gamma}(g)$, since it is fixed in the evaluation. Thus, $\text{Var}_{\mathbf{P}}[\widehat{R}_{\mathbf{P}}(g)] = \sigma_{\mathbf{P}}^2(g)/n_{\mathbf{P}}$ and $\text{Var}_{\mathbf{N}}[\widehat{R}_{\mathbf{N}}(g)] = \sigma_{\mathbf{N}}^2(g)/n_{\mathbf{N}}$, where $\widehat{R}_{\mathbf{P}}$ and $\widehat{R}_{\mathbf{N}}$ are the sample averages of $R_{\mathbf{P}}(g) = \mathbb{E}_{\mathbf{P}}[\ell(g(\mathbf{x}))]$ and $R_{\mathbf{N}}(g) = \mathbb{E}_{\mathbf{N}}[\ell(-g(\mathbf{x}))]$, respectively. When $n_{\mathbf{U}} \rightarrow \infty$,

$$\begin{aligned}
\text{Var}[\widehat{R}_{\text{N-PUNU}}^{\text{MR}, \gamma}(g)] &= 4(1 - \gamma)^2 \theta_{\mathbf{P}}^2 \text{Var}_{\mathbf{P}}[\widehat{R}_{\mathbf{P}}(g)] + 4\gamma^2 \theta_{\mathbf{N}}^2 \text{Var}_{\mathbf{N}}[\widehat{R}_{\mathbf{N}}(g)] \\
&= 4(1 - \gamma)^2 \psi_{\mathbf{P}} + 4\gamma^2 \psi_{\mathbf{N}} \\
&= 4(\psi_{\mathbf{P}} + \psi_{\mathbf{N}})\gamma^2 - 8\psi_{\mathbf{P}}\gamma + 4\psi_{\mathbf{P}}, \tag{5.135}
\end{aligned}$$

and it is obvious that $\gamma_{\text{N-PUNU}}^{\text{MR}} \in [0, 1]$. All other claims in Theorem 16 follow from that $\text{Var}[\widehat{R}_{\text{N-PUNU}}^{\text{MR}, \gamma}(g)]$ is quadratic in γ , that $\text{Var}[\widehat{R}_{\text{N-PUNU}}^{\text{MR}, \gamma}(g)] = \text{Var}[\widehat{R}_{\mathbf{PN}}^{\text{MR}}(g)]$ at $\gamma = 1/2$, and that $\gamma_{\text{N-PUNU}}^{\text{MR}} < 1/2$ if $\psi_{\mathbf{P}} < \psi_{\mathbf{N}}$ or $\gamma_{\text{N-PUNU}}^{\text{MR}} > 1/2$ if $\psi_{\mathbf{P}} > \psi_{\mathbf{N}}$.

Likewise, when $n_U \rightarrow \infty$,

$$\text{Var}[\hat{R}_{\text{N-PNPU}}^{\text{MR},\gamma}(g)] = (1 + \gamma)^2 \psi_P + (1 - \gamma)^2 \psi_N, \quad (5.136)$$

$$\text{Var}[\hat{R}_{\text{N-PNNU}}^{\text{MR},\gamma}(g)] = (1 - \gamma)^2 \psi_P + (1 + \gamma)^2 \psi_N, \quad (5.137)$$

and $\gamma_{\text{N-PNPU}}^{\text{MR}} \geq 0$ if $\psi_P \leq \psi_N$ or $\gamma_{\text{N-PNNU}}^{\text{MR}} \geq 0$ if $\psi_P \geq \psi_N$. The rest of proof of Theorem 17 is analogous to that of Theorem 16.

5.8.5 Proof of Theorem 18

Proof. The empirical PNPU-AUC risk can be expressed as

$$\begin{aligned} \hat{R}_{\text{PNPU}}^{\text{AUC},\gamma}(g) &= (1 - \gamma) \hat{R}_{\text{PN}}^{\text{AUC}}(g) + \gamma \hat{R}_{\text{PU}}^{\text{AUC}}(g) \\ &= \frac{1 - \gamma}{n_P n_N} \sum_{i=1}^{n_P} \sum_{j=1}^{n_N} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_j^N)) \\ &\quad + \frac{\gamma}{\theta_N n_P n_U} \sum_{i=1}^{n_P} \sum_{k=1}^{n_U} \ell(g(\mathbf{x}_i^P) - g(\mathbf{x}_k^U)) \\ &\quad - \frac{\gamma \theta_P}{\theta_N n_P^2} \sum_{i=1}^{n_P} \sum_{i'=1}^{n_P} \ell(g(\mathbf{x}_i^P) - g(\bar{\mathbf{x}}_{i'}^P)). \end{aligned} \quad (5.138)$$

Assume $n_U \rightarrow \infty$, we obtain

$$\begin{aligned} \text{Var}[\hat{R}_{\text{PNPU}}^{\text{AUC},\gamma}(g)] &= \frac{(1 - \gamma)^2}{n_P n_N} \tau_{\text{PN}}^2(g) + \frac{\gamma^2 \theta_P^2}{n_P^2 \theta_N^2} \tau_{\text{PP}}^2(g) + \frac{(1 - \gamma) \gamma}{\theta_N n_P} \tau_{\text{PN,PU}}(g) \\ &\quad - \frac{\gamma^2 \theta_P}{\theta_N^2 n_P} \tau_{\text{PU,PP}}(g) - \frac{(1 - \gamma) \gamma \theta_P}{\theta_N n_P} \tau_{\text{PN,PP}}(g) \\ &= (1 - \gamma)^2 \bar{\psi}_{\text{PN}} + \gamma^2 \bar{\psi}_{\text{PU}} + (1 - \gamma) \gamma \bar{\psi}_{\text{PP}}, \end{aligned} \quad (5.139)$$

where the terms divided by n_U are disappeared. Setting the derivative with respect to γ at zero, we obtain the minimizer in Eq. (5.85).

For the empirical PNNU-AUC risk, when $n_U \rightarrow \infty$, we obtain

$$\begin{aligned}
 \text{Var}[\widehat{R}_{\text{PNNU}}^{\text{AUC},\gamma}(g)] &= \frac{(1-\gamma)^2}{n_P n_N} \tau_{\text{PN}}^2(g) + \frac{\gamma^2 \theta_N^2}{n_N^2 \theta_P^2} \tau_{\text{NN}}^2(g) + \frac{(1-\gamma)\gamma}{\theta_P n_N} \tau_{\text{PN,NU}}(g) \\
 &\quad - \frac{\gamma^2 \theta_N}{\theta_P^2 n_N} \tau_{\text{NU,NN}}(g) - \frac{(1-\gamma)\gamma \theta_N}{\theta_P n_N} \tau_{\text{PN,NN}}(g) \\
 &= (1-\gamma)^2 \bar{\psi}_{\text{PN}} + \gamma^2 \bar{\psi}_{\text{NU}} + (1-\gamma)\gamma \bar{\psi}_{\text{NN}}.
 \end{aligned} \tag{5.140}$$

Setting the derivative with respect to γ at zero, we obtain the minimizer in Eq. (5.86). \square

Chapter 6

Conclusions and Future Work

In this chapter, we present conclusions and discuss several future perspectives of this research.

6.1 Conclusions

This dissertation was devoted to expanding the possibility of semi-supervised learning.

In Chapter 3, we considered imbalanced classification from positive-unlabeled (PU) learning scenario. The problem was solved by directly maximizing AUC with respect to a classifier. We presented the risk for AUC maximization in PU learning, revealing that the existing AUC maximization method from PU data is biased. The generalization error bound was proved without assuming distributional assumptions on the data. Experimental results reported the effectiveness of our method.

In Chapter 4, we developed statistical dependency analysis tools based on squared-loss mutual information (SMI) from PU data. We proposed an SMI estimation method from only PU data and proved its optimal convergence to true SMI. The SMI estimator was applied to dimension reduction and independence test. In particular, our dimension reduction can be executed without conducting a class-prior estimation in advance, which is a significant advantage in practice. We experimentally confirmed the usefulness of the SMI-based statistical dependency analysis methods.

In Chapter 5, we explored a semi-supervised learning approach based on positive-unlabeled learning. The approach allows us to leverage unlabeled data without conventional distributional assumptions such as the cluster assumption. We provided two realizations of the risks in semi-supervised learning; the misclassification rate and AUC as the classification measures. We elucidated properties of the risks and showed the generalization error bounds and the variance reduction of the empirical risks. Finally, extensive experimental results confirmed the effectiveness of our approach.

6.2 Future Work

Finally, we discuss future directions of this dissertation.

6.2.1 Semi-Supervised SMI Estimation

The first future direction is to extend the SMI estimation from PU data to SMI estimation from PNU data.

A simple way is to combine the PN-SMI risk with the PU-SMI or NU-SMI risks, where the NU-SMI risk can be defined as:

$$R_{\text{NU}}^{\text{SMI}}(g) = \frac{1}{2} \int \left(g(\mathbf{x}) - \frac{p(\mathbf{x} \mid y = -1)}{p(\mathbf{x})} \right)^2 p(\mathbf{x}) d\mathbf{x} - C_{\text{NU}}, \quad (6.1)$$

where $C_{\text{NU}} := \frac{1}{2} \int \frac{p^2(\mathbf{x} \mid y = -1)}{p(\mathbf{x})} d\mathbf{x}$. Along with the line of PNU learning, we minimize the PNU-SMI risk defined as

$$R_{\text{PNU}}^{\text{SMI}, \eta}(g) = \begin{cases} (1 - \eta)R_{\text{PN}}^{\text{SMI}}(g) + \eta R_{\text{PU}}^{\text{SMI}}(g) & (\eta \geq 0), \\ (1 + \eta)R_{\text{PN}}^{\text{SMI}}(g) - \eta R_{\text{NU}}^{\text{SMI}}(g) & (\eta < 0), \end{cases} \quad (6.2)$$

where $-1 \leq \eta \leq 1$. Then, with the learned model $g^* = \underset{g}{\operatorname{argmin}} R_{\text{PNU}}^{\text{SMI}, \eta}$, we can compute SMI by plugging the model into the following definition of SMI:

$$\text{SMI} = \begin{cases} (1 - \eta)\text{PN-SMI} + \eta\text{PU-SMI} & (\eta \geq 0), \\ (1 + \eta)\text{PN-SMI} - \eta\text{NU-SMI} & (\eta < 0), \end{cases} \quad (6.3)$$

where PN-SMI is the expression of SMI in Eq. (4.3) and NU-SMI is defined as

$$\text{NU-SMI} := \frac{\theta_{\text{N}}}{2\theta_{\text{P}}} \int \left(\frac{p(\mathbf{x} \mid y = -1)}{p(\mathbf{x})} - 1 \right)^2 p(\mathbf{x}) d\mathbf{x}. \quad (6.4)$$

Alternatively, for SMI estimation, we can simply use unlabeled data even under the supervised learning setting unlike classification tasks. Recall the PN-SMI risk in Eq. (2.97):

$$R_{\text{PN}}^{\text{SMI}}(q) = \frac{1}{2} \sum_{y=\pm 1} \int q^2(\mathbf{x}, y) p(\mathbf{x}) p(y) d\mathbf{x} - \sum_{y=\pm 1} \int q(\mathbf{x}, y) p(\mathbf{x}, y) d\mathbf{x}. \quad (6.5)$$

To approximate the PN-SMI risk, the original SMI estimation method uses the obtained labeled sample for computing sample averages regarding $p(\mathbf{x}, y)$, $p(\mathbf{x})$, and $p(y)$. However, for computing sample averages regarding $p(\mathbf{x})$ instead of marginalized labeled samples, we can use additional unlabeled samples,¹ meaning that this procedure is SMI estimation from labeled and unlabeled data.

Since both two approaches do not assume any strong distributional assumption for utilizing unlabeled data for SMI estimation, our future work includes comparing two approaches from both empirical and theoretical viewpoints.

6.2.2 AUC Optimization from PU Data without Class-Priors

Recall that the PU-AUC risk is expressed as

$$R_{\text{PU}}^{\text{AUC}}(g) = \frac{1}{\theta_N} \mathbb{E}_{\text{P}}[\mathbb{E}_{\text{U}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\mathbf{x}^{\text{U}}))] - \frac{\theta_{\text{P}}}{\theta_N} \mathbb{E}_{\text{P}} \mathbb{E}_{\bar{\text{P}}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\bar{\mathbf{x}}^{\text{P}}))]. \quad (6.6)$$

The key observation is that the second term becomes constant with loss functions satisfying

$$\ell(m) + \ell(-m) = C, \quad (6.7)$$

where C is a constant. Since $\mathbf{x}^{\text{P}} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = +1)$ and $\bar{\mathbf{x}}^{\text{P}} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = +1)$ are interchangeable, we have

$$\begin{aligned} & \mathbb{E}_{\text{P}}[\mathbb{E}_{\bar{\text{P}}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\bar{\mathbf{x}}^{\text{P}}))] \\ &= \frac{1}{2} \mathbb{E}_{\text{P}}[\mathbb{E}_{\bar{\text{P}}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\bar{\mathbf{x}}^{\text{P}}))] + \frac{1}{2} \mathbb{E}_{\text{P}}[\mathbb{E}_{\bar{\text{P}}}[\ell(g(\bar{\mathbf{x}}^{\text{P}}) - g(\mathbf{x}^{\text{P}}))] \\ &= \frac{1}{2} \mathbb{E}_{\text{P}}[\mathbb{E}_{\bar{\text{P}}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\bar{\mathbf{x}}^{\text{P}}))] + \frac{1}{2} \mathbb{E}_{\text{P}}[\mathbb{E}_{\bar{\text{P}}}[\ell(-(g(\mathbf{x}^{\text{P}}) - g(\bar{\mathbf{x}}^{\text{P}})))] \\ &= \frac{1}{2} \mathbb{E}_{\text{P}}[\mathbb{E}_{\bar{\text{P}}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\bar{\mathbf{x}}^{\text{P}})) + \ell(-(g(\mathbf{x}^{\text{P}}) - g(\bar{\mathbf{x}}^{\text{P}})))] \\ &= \frac{C}{2}. \end{aligned} \quad (6.8)$$

¹ This is because $p(\mathbf{x}) = \sum_y p(\mathbf{x}, y)$.

In the risk minimization, we can ignore the constant and the scaling parameter $1/\theta_N$. Thus, we can show the following relation:

$$\begin{aligned} g^* &= \operatorname{argmin}_g R_{\text{PU}}^{\text{AUC}}(g) \\ &= \operatorname{argmin}_g \mathbb{E}_{\text{P}}[\mathbb{E}_{\text{U}}[\ell(g(\mathbf{x}^{\text{P}}) - g(\mathbf{x}^{\text{U}}))]]. \end{aligned} \quad (6.9)$$

The advantage of Eq. (6.9) is that a classifier can be obtained without resorting to the class-prior estimation, meaning that the empirical risk minimization is not affected by the accuracy of the estimated class-prior. This idea would be an undoubtedly interesting future direction of the practical PU learning method for imbalanced classification.

6.2.3 Multi-Class Extension of Semi-Supervised AUC Optimization

Next, we discuss a multi-class extension of semi-supervised AUC optimization. The idea is based on the AUC optimization method from PU data without a class-prior estimation discussed in the previous section.

Let

$$\mathcal{Y} := \{1, \dots, c\} \quad (6.10)$$

be the domain of class labels, where c is the number of classes. Furthermore, let us define classifiers for each class as

$$G := \{g_y \mid g_y: \mathbb{R}^d \rightarrow \mathbb{R}, \sum_{y'=1}^c g_{y'} = 1, y = 1, \dots, c\}. \quad (6.11)$$

Let us define the risk in the multi-class version of supervised AUC optimization as

$$R_{\text{MC}}^{\text{AUC}}(G) := \frac{1}{|\mathcal{Y}|} \sum_y \sum_{y' \neq y} \theta_{y'} \int \ell(g_y(\mathbf{x}) - g_{y'}(\mathbf{x}')) p(\mathbf{x} \mid y) p(\mathbf{x}' \mid y') d\mathbf{x} d\mathbf{x}'. \quad (6.12)$$

Then, the risk in semi-supervised multi-class AUC optimization is given by

$$R_{\text{SSL}}^{\text{AUC}}(G) := \frac{1}{|\mathcal{Y}|} \sum_y \int \ell(g_y(\mathbf{x}) - g_y(\mathbf{x}')) p(\mathbf{x} \mid y) p(\mathbf{x}') d\mathbf{x} d\mathbf{x}'. \quad (6.13)$$

Similarly to the calculation in the binary case shown in the previous section, with the loss functions satisfying $\ell(m) + \ell(-m) = C$, we can reformulate $R_{\text{SSL}}^{\text{AUC}}$ as

$$\begin{aligned} R_{\text{SSL}}^{\text{AUC}}(G) &= \frac{C}{2|\mathcal{Y}|} + \frac{1}{|\mathcal{Y}|} \sum_y \sum_{y' \neq y} \theta_{y'} \int \ell(g_y(\mathbf{x}) - g_{y'}(\mathbf{x}')) p(\mathbf{x} | y) p(\mathbf{x}' | y') d\mathbf{x} d\mathbf{x}' \\ &= \frac{C}{2|\mathcal{Y}|} + R_{\text{MC}}^{\text{AUC}}(G). \end{aligned} \quad (6.14)$$

Again, the first term is constant; the minimization of $R_{\text{SSL}}^{\text{AUC}}$ corresponds to the minimization of the second term, i.e., the risk in supervised multi-class AUC optimization $R_{\text{MC}}^{\text{AUC}}$. Therefore, we can obtain the classifiers from labeled and unlabeled data. Investigating theoretical and empirical properties of this multi-class semi-supervised AUC optimization risk is also an important future direction.

6.2.4 Essential Mechanism of Extracting Information from Data

The final future direction is rather abstract, but it would be an important research question. The question is revealing a mechanism of extracting information from unlabeled data.

To achieve better performance from limited labeled data by utilizing a large amount of unlabeled data, we need a carefully-designed algorithm to extract information hidden in the data. As discussed in this dissertation several times, the conventional assumption such as the low-density separation principle does not always work well. For overcoming the issue, this dissertation was devoted to the method utilizing unlabeled data without such a distributional assumption. However, we are still on the way of practical semi-supervised learning methods and feel that we are extracting an only small portion of information from the data. Note that achieving $1/\sqrt{n_U}$ is optimal in terms of statistical estimation. Towards making machine learning methods more practical, we need general and yet powerful semi-supervised learning methods.

To pursue this research question, there is no clear direction yet, but several recent advances in machine learning might become a clue and shed light in a part of our ultimate goal. There are some interesting studies, e.g., influence functions for analyzing black-box models (Koh and Liang, 2017) such as a deep neural network model, and iterative machine teaching for obtaining a better model more quickly (Liu et al., 2017). The key of these studies is to analyze the loss function and utilize the property of the loss for their purposes. In this dissertation, our focus was the form of risk functions. To go one step further for more practical semi-supervised learning, analysis of loss functions would

be a current possible approach for revealing informative data points in semi-supervised learning scenario and indicating a way of extracting more information from unlabeled data.

Bibliography

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. Software available from <https://www.tensorflow.org/>.
- Ali, S. M. and Slivey, S. D. (1966). A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society. Series B*, 28:131–152.
- Amini, M. R., Truong, T. V., and Goutte, C. (2008). A boosting algorithm for learning bipartite ranking functions with partially labeled data. In *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 99–106.
- Bartlett, P. L., Jordan, M. I., and McAuliffe, J. D. (2006). Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156.
- Basu, A., Harris, I. R., Hjort, N. L., and Jones, M. C. (1998). Robust and efficient estimation by minimising a density power divergence. *Biometrika*, 85(3):549–559.
- Belkin, M., Niyogi, P., and Sindhvani, V. (2006). Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of Machine Learning Research*, 7:2399–2434.
- Blanchard, G., Lee, G., and Scott, C. (2010). Semi-supervised novelty detection. *Journal of Machine Learning Research*, 11:2973–3009.
- Blondel, M., Seki, K., and Uehara, K. (2013). Block coordinate descent algorithms for large-scale sparse multiclass classification. *Machine Learning*, 93(1):31–52.

- Bonnans, J. F. and Cominetti, R. (1996). Perturbed optimization in Banach spaces I: A general theory based on a weak directional constraint qualification; II: A theory based on a strong directional qualification condition; III: Semiinfinite optimization. *SIAM Journal on Control and Optimization*, 34(4):1151–1171, 1172–1189, and 1555–1567.
- Bonnans, J. F. and Shapiro, A. (1998). Optimization problems with perturbations: A guided tour. *SIAM Review*, 40(2):228–264.
- Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press, New York, NY, USA.
- Chang, C.-C. and Lin, C.-J. (2011). LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Chapelle, O., Schölkopf, B., and Zien, A., editors (2006). *Semi-Supervised Learning*. MIT Press.
- Chapelle, O., Sindhwani, V., and Keerthi, S. S. (2008). Optimization techniques for semi-supervised support vector machines. *Journal of Machine Learning Research*, 9:203–233.
- Chapelle, O., Weston, J., and Schölkopf, B. (2002). Cluster kernels for semi-supervised learning. In *Advances in Neural Information Processing Systems 15*, pages 585–592.
- Chapelle, O. and Zien, A. (2005). Semi-supervised classification by low density separation. In *AISTATS*, pages 57–64.
- Chen, T., Li, M., Li, Y., Lin, M., Wang, N., Wang, M., Xiao, T., Xu, B., Zhang, C., and Zhang, Z. (2015). Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*.
- Collobert, R., Sinz, F., Weston, J., and Bottou, L. (2006). Trading convexity for scalability. In *Proceedings of the 23rd International Conference on Machine Learning*, pages 201–208.
- Cortes, C. and Mohri, M. (2003). AUC optimization vs. error rate minimization. In *Advances in Neural Information Processing Systems 16*, pages 313–320.
- Cortes, C. and Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3):273–297.

- Cover, T. M. and Thomas, J. A. (2006). *Elements of Information Theory*. John Wiley & Sons, Inc., Hoboken, NJ, USA, 2nd edition.
- Cozman, F. G., Cohen, I., and Cirelo, M. C. (2003). Semi-supervised learning of mixture models. In *Proceedings of the 20th International Conference on Machine Learning*, pages 99–106.
- Csiszár, I. (1967). Information-type measures of difference of probability distributions and indirect observation. *Studia Scientiarum Mathematicarum Hungarica*, 2:229–318.
- Dredze, M., Crammer, K., and Pereira, F. (2008). Confidence-weighted linear classification. In *Proceedings of the 25th International Conference on Machine learning*, pages 264–271.
- du Plessis, M. C., Niu, G., and Sugiyama, M. (2014). Analysis of learning from positive and unlabeled data. In *Advances in Neural Information Processing Systems 27*, pages 703–711.
- du Plessis, M. C., Niu, G., and Sugiyama, M. (2015). Convex formulation for learning from positive and unlabeled data. In *Proceedings of 32nd International Conference on Machine Learning*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 1386–1394.
- du Plessis, M. C., Niu, G., and Sugiyama, M. (2017). Class-prior estimation for learning from positive and unlabeled data. *Machine Learning*, 106(4):463–492.
- du Plessis, M. C. and Sugiyama, M. (2014). Class prior estimation from positive and unlabeled data. *IEICE Transactions on Information and Systems*, E97-D(5):1358–1362.
- Efron, B. and Tibshirani, R. J. (1994). *An introduction to the bootstrap*. CRC Press.
- Elkan, C. and Noto, K. (2008). Learning classifiers from only positive and unlabeled data. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 213–220.
- Fujino, A. and Ueda, N. (2016). A semi-supervised AUC optimization method with generative models. In *IEEE 16th International Conference on Data Mining*, pages 883–888.
- Gao, W., Wang, L., Jin, R., Zhu, S., and Zhou, Z.-H. (2016). One-pass AUC optimization. *Artificial Intelligence*, 236(C):1–29.

- Gao, W. and Zhou, Z.-H. (2015). On the consistency of AUC pairwise optimization. In *International Joint Conference on Artificial Intelligence*, pages 939–945.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press.
- Grandvalet, Y. and Bengio, Y. (2004). Semi-supervised learning by entropy minimization. In *Advances in Neural Information Processing Systems 17*, pages 529–536.
- Grandvalet, Y. and Bengio, Y. (2006). Entropy regularization. In Chapelle, O., Schölkopf, B., and Zien, A., editors, *Semi-Supervised Learning*. MIT Press, Cambridge Massachusetts.
- Hanley, J. A. and McNeil, B. J. (1982). The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology*, 143(1):29–36.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning*. Springer-Verlag, New York, 2 edition.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778.
- Herschtal, A. and Raskutti, B. (2004). Optimising area under the ROC curve using gradient descent. In *Proceedings of the 21st International Conference on Machine Learning*.
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 448–456.
- Jain, S., White, M., and Radivojac, P. (2016). Estimating the class prior and posterior from noisy positives and unlabeled data. In *Advances in Neural Information Processing Systems 29*.
- Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., and Darrell, T. (2014). Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM International Conference on Multimedia*, pages 675–678. ACM.
- Joachims, T. (2002). Optimizing search engines using clickthrough data. In *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 133–142.

- Kanamori, T., Hido, S., and Sugiyama, M. (2009). A least-squares approach to direct importance estimation. *Journal of Machine Learning Research*, 10:1391–1445.
- Kanamori, T., Suzuki, T., and Sugiyama, M. (2012). Statistical analysis of kernel-based least-squares density-ratio estimation. *Machine Learning*, 86(3):335–367.
- Kawakubo, H., du Plessis, M. C., and Sugiyama, M. (2016). Computationally efficient class-prior estimation under class balance change using energy distance. *IEICE Transactions on Information and Systems*, E99-D(1):176–186.
- Kelley, Jr, J. E. (1960). The cutting-plane method for solving convex programs. *Journal of the Society for Industrial and Applied Mathematics*, 8(4):703–712.
- Keziou, A. (2003). Dual representation of φ -divergences and applications. *Comptes Rendus Mathématique*, 336(10):857–862.
- Koh, P. W. and Liang, P. (2017). Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 1885–1894.
- Kotlowski, W., Dembczynski, K. J., and Huellermeier, E. (2011). Bipartite ranking through minimization of univariate loss. In *Proceedings of the 28th International Conference on Machine Learning*, pages 1113–1120.
- Krause, A., Perona, P., and Gomes, R. G. (2010). Discriminative clustering by regularized information maximization. In *Advances in Neural Information Processing Systems*, pages 775–783.
- Krijthe, J. H. and Loog, M. (2017). Robust semi-supervised least squares classification by implicit constraints. *Pattern Recognition*, 63:115–126.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105.
- Lang, K. (1995). NewsWeeder: Learning to filter netnews. In *Proceedings of the 12th International Machine Learning Conference*.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.

- Ledoux, M. and Talagrand, M. (1991). *Probability in Banach Spaces: Isoperimetry and Processes*. Springer.
- Lee, W. S. and Liu, B. (2003). Learning with positive and unlabeled examples using weighted logistic regression. In *Proceedings of the 20th International Conference on Machine Learning*, pages 448–455.
- Letouzey, F., Denis, F., and Gilleron, R. (2000). Learning from positive and unlabeled examples. In *Proceedings of the 11th International Conference on Algorithmic Learning Theory*, pages 71–85, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Lewis, D. D., Yang, Y., Rose, T. G., and Li, F. (2004). RCV1: A new benchmark collection for text categorization research. *Journal of Machine Learning Research*, 5:361–397.
- Li, K.-C. (1991). Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association*, 86(414):316–327.
- Li, X.-L. and Liu, B. (2005). Learning from positive and unlabeled examples with different data distributions. In *Proceedings of 16th European Conference on Machine Learning*, pages 218–229.
- Li, Y.-F., Tsang, I. W., Kwok, J. T., and Zhou, Z.-H. (2013). Convex and scalable weakly labeled SVMs. *Journal of Machine Learning Research*, 14(1):2151–2188.
- Li, Y.-F. and Zhou, Z.-H. (2015). Towards making unlabeled data never hurt. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(1):175–188.
- Lichman, M. (2013). UCI machine learning repository.
- Linsker, R. (1988). Self-organization in a perceptual network. *Computer*, 21(3):105–117.
- Liu, B., Dai, Y., Li, X., Lee, W. S., and Yu, P. S. (2003). Building text classifiers using positive and unlabeled examples. In *Third IEEE International Conference on Data Mining*, pages 179–186.
- Liu, W., Dai, B., Humayun, A., Tay, C., Yu, C., Smith, L. B., Rehg, J. M., and Song, L. (2017). Iterative machine teaching. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 2149–2158.

- Loog, M. (2016). Contrastive pessimistic likelihood estimation for semi-supervised classification. *IEEE Transactions on Pattern Analysis Machine Intelligence*, 38(3):462–475.
- Melacci, S. and Belkin, M. (2011). Laplacian support vector machines trained in the primal. *Journal of Machine Learning Research*, 12:1149–1184.
- Mendelson, S. (2008). Lower bounds for the empirical minimization algorithm. *IEEE Transactions on Information Theory*, 54(8):3797–3803.
- Mohri, M., Rostamizadeh, A., and Talwalkar, A. (2012). *Foundations of Machine Learning*. MIT Press.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. (2011). Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*.
- Nguyen, X. L., Wainwright, M. J., and Jordan, M. I. (2007). Nonparametric estimation of the likelihood ratio and divergence functionals. In *IEEE International Symposium on Information Theory*, pages 2016–2020.
- Niu, G., du Plessis, M. C., Sakai, T., Ma, Y., and Sugiyama, M. (2016). Theoretical comparisons of positive-unlabeled learning against positive-negative learning. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., and Garnett, R., editors, *Advances in Neural Information Processing Systems 29*, pages 1199–1207.
- Niu, G., Jitkrittum, W., Dai, B., Hachiya, H., and Sugiyama, M. (2013). Squared-loss mutual information regularization: A novel information-theoretic approach to semi-supervised learning. In *Proceedings of the 30th International Conference on Machine Learning*, volume 28, pages 10–18.
- Nocedal, J. and Wright, S. J. (2006). *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer.
- Pearson, K. (1990). On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arise from random sampling. *Philosophical Magazine Series 5*, 50(302):157–175.
- Quiñonero Candela, J., Sugiyama, M., Schwaighofer, A., and Lawrence, N. D. (2009). *When Training and Test Sets Are Different: Characterizing Learning Transfer*, pages 3–28. MIT Press.

- Rakhlin, A., Shamir, O., and Sridharan, K. (2012). Making gradient descent optimal for strongly convex stochastic optimization. In *Proceedings of the 29th International Conference on Machine Learning*, pages 449–456.
- Ramaswamy, H. G., Scott, C., and Tewari, A. (2016). Mixture proportion estimation via kernel embedding of distributions. In *Proceedings of the 33rd International Conference on Machine Learning*.
- Rätsch, G., Onoda, T., and Müller, K.-R. (2001). Soft margins for adaboost. *Machine learning*, 42(3):287–320.
- Rigollet, P. (2007). Generalization error bounds in semi-supervised classification under the cluster assumption. *Journal of Machine Learning Research*, 8:1369–1392.
- Sakai, T. and Sugiyama, M. (2014). Computationally efficient estimation of squared-loss mutual information with multiplicative kernel models. *IEICE Transactions on Information and Systems*, E97-D(4):968–971.
- Sakai, T., Sugiyama, M., Kitagawa, K., and Suzuki, K. (2015). Registration of infrared transmission images using squared-loss mutual information. *Precision Engineering*, 39:187–193.
- Schölkopf, B. and Smola, A. J. (2001). *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, USA.
- Sechidis, K., Calvo, B., and Brown, G. (2014). Statistical hypothesis testing in positive unlabelled data. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 66–81.
- Shalev-Shwartz, S. and Ben-David, S. (2014). *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, New York, NY, USA.
- Shannon, C. (1948). A mathematical theory of communication. *Bell Systems Technical Journal*, 27:379–423.
- Singh, S. (2000). *The Code Book: The Science of Secrecy from Ancient Egypt to Quantum Cryptography*. Anchor.
- Sokolovska, N., Cappé, O., and Yvon, F. (2008). The asymptotics of semi-supervised learning in discriminative probabilistic models. In *Proceedings of the 25th International Conference on Machine Learning*, pages 984–991.

- Sugiyama, M. (2013). Machine learning with squared-loss mutual information. *Entropy*, 15:80–112.
- Sugiyama, M., Niu, G., Yamada, M., Kimura, M., and Hachiya, H. (2014). Information-maximization clustering based on squared-loss mutual information. *Neural Computation*, 26(1):84–131.
- Sugiyama, M. and Suzuki, T. (2011). Least-squares independence test. *IEICE Transactions on Information and Systems*, E94-D:1333–1336.
- Sundararajan, S., Priyanka, G., and Selvaraj, Keerthi, S. S. (2011). A pairwise ranking based approach to learning with positive and unlabeled examples. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management*, pages 663–672.
- Suzuki, T. and Sugiyama, M. (2013). Sufficient dimension reduction via squared-loss mutual information. *Neural Computation*, 25(3):725–758.
- Suzuki, T., Sugiyama, M., Kanamori, T., and Sese, J. (2009). Mutual information estimation reveals global associations between stimuli and biological processes. *BMC Bioinformatics*, 10:S52:1–12.
- Székely, G. J. and Rizzo, M. L. (2013). Energy statistics: A class of statistics based on distances. *Journal of Statistical Planning and Inference*, 143(8):1249–1272.
- Tokui, S., Oono, K., Hido, S., and Clayton, J. (2015). Chainer: A next-generation open source framework for deep learning. In *Proceedings of Workshop on Machine Learning Systems in The Twenty-ninth Annual Conference on Neural Information Processing Systems*.
- Torkkola, K. (2003). Feature extraction by non-parametric mutual information maximization. *Journal of Machine Learning Research*, 3:1415–1438.
- Torralba, A., Fergus, R., and Freeman, W. T. (2008). 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(11):1958–1970.
- Usunier, N., Amini, M., and Patrick, G. (2006). Generalization error bounds for classifiers trained with interdependent data. In Weiss, Y., Schölkopf, P. B., and Platt, J. C., editors, *Advances in Neural Information Processing Systems 18*, pages 1369–1376.

- Vapnik, V. (1995). *The Nature of Statistical Learning Theory*. Springer-Verlag New York, Inc., New York, NY, USA.
- Vapnik, V. N. (1998). *Statistical Learning Theory*. John Wiley & Sons.
- Yamada, M., Sigal, L., Raptis, M., Toyoda, M., Chang, Y., and Sugiyama, M. (2015). Cross-domain matching with squared-loss mutual information. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1764–1776.
- Ying, Y., Wen, L., and Lyu, S. (2016). Stochastic online AUC maximization. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., and Garnett, R., editors, *Advances in Neural Information Processing Systems 29*, pages 451–459.
- Yuille, A. L. and Rangarajan, A. (2002). The concave-convex procedure (CCCP). In *Advances in Neural Information Processing Systems 14*, pages 1033–1040.
- Zhao, P., Jin, R., Yang, T., and Hoi, S. C. (2011). Online AUC maximization. In *Proceedings of the 28th International Conference on Machine Learning*, pages 233–240.
- Zhou, B., Lapedriza, A., Xiao, J., Torralba, A., and Oliva, A. (2014). Learning deep features for scene recognition using places database. In *Advances in Neural Information Processing Systems 27*, pages 487–495.